LECTURE NOTES

ON

DIGITAL SIGNAL PROCESSING

IV B.Tech I semester (JNTUH-R15)

Mr. A Naresh Kumar Assistant professor



ELECTRICAL AND ELECTRONICS ENGINEERING

INSTITUTE OF AERONAUTICAL ENGINEERING (Autonomous)

DUNDIGAL, HYDERABAD - 500043

Digital Signal Processing

I. Introduction

Introduction to Digital Signal Processing: Discrete time signals & sequences, Linear shift invariant systems, Stability, and Causality, Linear constant coefficient difference equations, Frequency domain representation of discrete time signals and systems.

Contents:

Sampling theory Discrete-time signals Transformation of the independent variable Discrete-time systems Linear constant coefficient difference equations Fourier analysis of discrete-time signals and systems Frequency response of discrete-time system Properties of the discrete-time Fourier transform (DTFT)

Sampling theory

Illustrative example A continuous-time random signal is shown. Based on this several important concepts are shown below. The signal is a continuous-time signal with continuous amplitude. Such a signal is also called an **analog signal**.



	x(t)‡										
4	8										
3	7										
2	6										
1	5										
0	4										⟨Time
-1	3										
-2	2										
-3	1										
-4	0										
	nT∢	0	1T	2T	3T	4T	5T	6T	7T	8T	<time< td=""></time<>
n<		0	1	2	3	4	5	6	7	8	
2	x(n) ≺	5.5	2.8	3.8	5.3	1.5	4.6	8.4	6.9	7.3	Sampled signal. Discrete-time
											amplitude is continuous.
		5	2	3	5	1	4	7	6	7	>Quantized. Quantization noise
								-			(error). Digital signal – both
		101	010	011	101	001	100	111	110	111	Encoded to 3 bits/sample
		101	010	011	101	001	100	111	110		Encoded to 5 bits/sample.
								\$Note this particular point exhibits saturation (out of			
								range). Rounded down to 7, not 8.			

If we were to represent every sample value with infinite precision (for example, x(1) = 2.8--, instead of being approximated as 2 or 3) then we would need registers and memory words of arbitrarily large size. However, owing to a finite word length we round off the sample values (in this case x(1) = 2.8-- will be rounded to 2). This introduces quantization noise or error.

The procedure of generating a discrete-time signal from an analog signal is shown in the following block diagram. In the digital signal processing course we are mostly dealing with *discrete-time* rather than *digital* signals and systems, the latter being a subset of the former.



The three boxes shown above can be represented by an analog to digital converter (ADC). A complete digital signal processing (DSP) system consists of an ADC, a DSP algorithm (e.g., a difference equation) and a digital to analog converter (DAC) shown below.



As the name implies discrete-time signals are defined only at discrete instants of time. Discrete-time signals can arise by sampling analog signals such as a voice signal or a temperature signal in telemetry. Discrete-time signals may also arise naturally, e.g., the number of cars sold on a specific day in a year, or the closing DJIA figure for a specific day of the year.

AT&T's *T1* **Stream** The voice signal is band limited to 3.3 kHz, sampled at 8000 Hz (8000 samples per second), quantized and encoded into 8 bits per sample. Twenty four such voice channels are combined to form the *T1* stream or signal.

Sampling interval = $\frac{1}{8000 Hz}$ = 0.125 msec. Bit rate for each channel = $8000 \frac{samples}{sec} \ge 8 \frac{bits}{sample}$ = 64000 bits/sec. Bit rate for T1 = 64000 bits/sec per channel \times 24 channels = 1 544 000 bits/sec.

Commercial examples CD, Super Audio CD (SACD), DVD Audio (DVD-A), Digital audio broadcasting - 32 kHz, and Digital audio tape (DAT) - 48 kHz.

	CD	Super Audio CD (SACD)	DVD-Audio (DVD-A)
Sampling	44.1 kHz	2.8224 MHz	44.1, 88.2 or
Rate			48, 96, 192
			kHz
Coding	16-bit PCM per sample	1-bit DSD	12-, 20-, 24-
	With 2 channels the bit rate is 1.4112	(Direct Stream	bit
	Mbits/sec, but additional error control	Digital – like Delta	
	bits etc., raise it to 4.3218 Mbits/sec.	modulation)	

Commercial Audio Examples

Pulse-train sampling For pulse-train sampling of the signal x(t) by the rectangular pulse-train p(t) resulting in the sampled signal $x_s(t)$, we have

$$x_s(t) = x(t) p(t)$$

The Fourier series of p(t) is given by $p(t) = \sum_{n = -\infty}^{\infty} e^{j n 2\pi F_s t}$,

where $F_s = 1/T$ and the Fourier coefficients, are

00

$$C_n = \frac{1}{T} \int_{-T/2} p(t) e^{-jn \, 2\pi F_s t} dt$$

Thus

$$x_s(t) = \sum_{n = -\infty} C_n x(t) e^{j n 2\pi F_s t}$$

The Fourier spectrum of $x_s(t)$ is given by

$$X_{s}(F) = \int_{-\infty}^{\infty} x_{s}(t)e^{-j2\pi Ft} dt = \int_{-\infty}^{\infty} \sum_{-\infty}^{\infty} C_{n}x(t)e^{jn2\pi Fs}e^{-j2\pi Ft} dt$$

Interchanging the order of integration and summation yields the aliasing formula

$$X_{s}(F) = \sum_{n=-\infty}^{\infty} C \int_{-\infty}^{n} \int_{-\infty}^{\infty} x(t) e_{-j2\pi(F-nF_{s})t} dt = \sum_{n=-\infty}^{\infty} C_{n}X(F-nF_{s})$$

= ...+ $C_{-2}X(F+2F_{s}) + C_{-1}X(F+1F_{s}) + C_{0}X(F) + C_{1}X(F-1F_{s})$
+ $C_{2}X(F-2F_{s}) + ...$



The sampled signal spectrum $X_s(F)$ is sketched below. For convenience of illustration we have assumed the base band spectrum, X(F), to be real valued; the maximum value of |X(F)| is taken to be 1. $X_s(F)$ consists of replicas of X(F), scaled by the Fourier coefficients C_n and repeated at intervals of F_s . Specifically, the replica at the origin is simply X(F) scaled by C_0 . Note that the magnitudes, $|C_n|$, have even symmetry. In this case, since $F_M \leq F_s - F_M$, there is no overlap among the replicas in the graph of $X_s(F)$. As a result the original signal x(t) can be

recovered by passing $x_s(t)$ through a low pass filter with a bandwidth *B* that satisfies the condition $F_M \le B \le F_s - F_M$, and a gain of $1/C_0$.



Impulse-train sampling If p(t) is an impulse-train then the Fourier coefficients are given by $C_n = \frac{1}{T} \int_{-T/2}^{T/2} \delta_{t}(t) e^{-jn2\pi F_s t} dt = \frac{1}{T} \cdot 1 = F_s \text{ for all } n$

so that the *aliasing formula* becomes

$$X_s(F) = F_s \sum_{n = -\infty} X \left(F - nF_s \right)$$

Alternative derivation It can be shown that

$$X_{s}(j\Omega) = \frac{\Omega}{2\pi} \sum_{n=-\infty}^{s} X(j(\Omega - n\Omega_{s})) = F_{s} \sum_{n=-\infty}^{s} X(j(\Omega - n\Omega_{s}))$$

Note that some use the notation $X(\Omega)$ instead of $X(j\Omega)$, so that the above equation is written as

$$X_{s}(\Omega) = \frac{\Omega_{s}}{2\pi} \sum_{n=-\infty}^{\infty} X(\Omega - n\Omega_{s}) = F_{s} \sum_{n=-\infty}^{\infty} X(\Omega - n\Omega_{s})$$

About terminology The highest frequency in the signal is called the *Nyquist frequency*. The minimum sampling rate that, in theory, allows perfect signal recovery is called the *Nyquist rate*. Thus Nyquist rate is twice the Nyquist frequency.

A signal, however, is generally sampled at more than twice the Nyquist frequency. One half the sampling rate is called the *folding frequency*. As an example, if a voice signal is band limited to 3.3 kHz and sampled at 8 kHz then the Nyquist frequency = 3.3 kHz, the Nyquist rate = 6.6 kHz, and the folding frequency = 4 kHz.

To add to the ambiguity, some references use the phrase Nyquist frequency to refer to one half the sampling rate.

Aliasing Illustration using a signal spectrum band-limited to 4 kHz and a sampling rate of 8 kHz. The signal is perfectly band-limited to 4 kHz so that there is no overlapping in the spectrum of the sampled signal.



The signal has a genuine, desirable, 1 kHz component. Since it is not perfectly band-limited to 4 kHz it has another genuine but undesirable 7 kHz component. Due to the first pair of replicas (centered at 8 kHz and -8 kHz) this 7kHz component appears as if it were a 1kHz component – in other words the 7 kHz component is an alias of 1kHz. Thus the first (lowest) alias of the 1 kHz frequency is given by 8 kHz – 1 kHz = 7 kHz.

Due to the second pair of replicas at 16 kHz and -16 kHz the 15k component in the original signal also appears as if it were a 1k component – it is another alias of the 1k. This second alias of the 1 kHz frequency is given by 16 kHz – 1 kHz = 15 kHz. The next alias is (3 x 8 - 1) kHz = 23 kHz.

In general, for any frequency F_1 within the base band (in this case any frequency from 0 to 4000 Hz) its aliases are given by

Alias =
$$kF_s - F_I$$
, k is an integer > 0
6 of 77

Aliasing example As an example, let a certain base band signal be band-limited to 8 Hz and let the sampling frequency be $F_s = 16$ Hz. We do not expect frequencies higher than 8 Hz. Then for the base band frequency of, say, $F_1 = 4$ Hz, the aliases, F_2 , are given by

$$F_2 = kF_s - F_1,$$
 $k = 1, 2, ...$

Setting k = 1 gives the lowest alias of F_1 , that is, $F_2 = 1 \ge 16 - 4 = 12$ Hz. The next alias is $2 \ge 16 - 4 = 28$ Hz. Consider two waveforms with frequencies of 4 Hz and 12 Hz given by

 $x_1(t) = 1 \cos 2\pi 4t$ and $x_2(t) = 1 \cos 2\pi 12t$

The sampled versions are

$$x_1(n) = \cos 2\pi 4nT = \cos 2\pi 4n(1/16) = \cos (n\pi/2)$$
 and
 $x_2(n) = \cos 2\pi 12nT = \cos 2\pi 12n(1/16) = \cos (n3\pi/2)$

whose first few samples are tabulated below:

n	0	1	2	3	4	5	6	7
$x_{l}(n)$	1	0	-1	0	1	0	-1	0
$x_2(n)$	1	0	-1	0	1	0	-1	0

These two sequences are seen to have the same digital frequency and cannot be distinguished from each other *as far as the frequency is concerned*. When passed through a smoothing filter, they will both appear as 4 Hz signals. This is true even if the *amplitudes* of the underlying analog waveforms are different.

In MATLAB

% Plot $x_1(t) = \cos 2\pi 4t$ as t goes from 0 to 0.5 sec (2 cycles) in steps of T = 1/160 sec. $t = 0: 1/160: 0.5; x1 = \cos(2*pi*4*t); plot(t, x1)$ % Repeat with step size of 1/16 sec. $t = 0: 1/16: 0.5; x1 = \cos(2*pi*4*t); plot(t, x1)$ % Produce samples of $x_1(t)$ as t goes from 0 to 0.5 in steps of 1/16 sec. $t = 0: 1/16: 0.5; x1 = \cos(2*pi*4*t)$ % Produce samples of $x_1(n) = \cos n\pi/2$ as n goes from 0 to 8 (2 cycles) in steps of 1 %(T = 1/16 sec.) $n = 0: 1: 8; x1 = \cos(n*pi/2)$ % Plot $x_1(n) = \cos n\pi/20$ as n goes from 0 to 80 (2 cycles) in steps of 1 (T = 1/160 sec.) n = 0: 1: 80; $x1 = \cos(n*pi/20)$; plot(n, x1, 'bo'); grid %Blue circles and grid % Stem plot $x_1(n) = \cos n\pi/2$ as n goes from 0 to 8 (2 cycles) in steps of 1 (T = 1/16 sec.) n = 0: 1: 8; x1 = cos (n*pi/2); stem(n, x1) % Stem plot $x_1(n) = \cos n\pi/20$ as n goes from 0 to 80 (2 cycles) in steps of 1 %(T = 1/160 sec.)n = 0: 1: 80; $x1 = \cos(n*pi/20)$; stem(n, x1) % Titles, labels and grid. Stem plot $x_1(n) = \cos n\pi/20$ as n goes from 0 to 80 (2 cycles) %MATLAB won't accept the kind of single quote in title ...Sampled Cosine" n = 0: 1: 80; $x1 = \cos(n*pi/20)$; stem(n, x1); title(,,Sampled Cosine"); xlabel(,,n"), ylabel(,,x1"); grid

%Plot $x_1(t) = \cos 2\pi 4t$ as t goes from 0 to 0.5 sec (2 cycles) in steps of T = 1/160 sec. t = 0: 1/160: 0.5; x1 = cos (2*pi*4*t); plot(t, x1)

 $x1 = \cos (2*pi*4*t) - 4Hz$ Cosine *plotted* at T = 1/160 sec.



% Plot $x_1(n) = \cos n\pi/20$ as n goes from 0 to 80 (2 cycles) in steps of 1 (T = 1/160 sec.) n = 0: 1: 80; x1 = cos (n*pi/20); plot(n, x1, 'bo'); grid % Blue circles and grid

 $x1 = \cos(n*pi/20) - 4$ Hz Cosine *sampled* at 160 samples per second



%Stem plot $x_1(n) = \cos n\pi/20$ as n goes from 0 to 80 (2 cycles) in steps of 1 %(T = 1/160 sec.) n = 0: 1: 80; x1 = cos (n*pi/20); stem(n, x1)

Stem plot of $x_1 = \cos(n*pi/20) - 4$ Hz Cosine *sampled* at 160 samples per second



Aliasing and digital frequency With $F_s = 16$ Hz the base band signal must be band-limited to 8 Hz. And we do not expect frequencies higher than 8 Hz. Consider the 3 signals $x_1(t)$, $x_2(t)$, and $x_3(t)$ of frequencies 4Hz, 12Hz and 28Hz, respectively, where x_2 and x_3 "s frequencies are the aliases of x_1 "s. The continuous and discrete-time signals are given in table below with a sampling rate of 16 samples/sec.

Aliasing and Digital Frequency					
Analog frequency	Analog signal	Discrete-time signal	Digital frequency		
Cycles/sec			Cycles/sample		
4	$x_I(t) = \cos 2\pi 4t$	$x_1(n) = \cos 2\pi (1/4)n$	0.25		
12	$x_2(t) = \cos 2\pi 12t$	$x_2(n) = \cos 2\pi (3/4)n$	0.75		
28	$x_3(t) = \cos 2\pi 28t$	$x_3(n)=\cos 2\pi(7/4)n$	1.75		
8	$x_4(t) = \cos 2\pi 8t$	$x_4(n) = \cos 2\pi (1/2)n$	0.5		
16/3	$x_5(t) = \cos 2\pi (16/3)t$	$x_5(n) = \cos 2\pi (1/3)n$	1/3		

If the Nyquist criterion is to be satisfied (for perfect signal reconstruction) we never expect any digital frequencies higher than 0.5 cycle/sample (or π rad./sample) in the base band signal. This corresponds to taking 2 samples for every cycle (or a digital frequency of half a cycle per sample) – the Nyquist criterion. Digital frequencies higher than 0.5 cycle/sample ($x_2(n)$ and $x_3(n)$ in this example) are actually disallowed.

When plotted $x_1(n)$, $x_2(n)$, and $x_3(n)$ cannot be distinguished from one another as far as the digital frequency is concerned. They all have a period = 4 (samples) and a frequency of 0.25 cycle per sample, though we know that $x_2(n)$ has a frequency of 0.75 cycle/sample and $x_3(n)$ has a frequency of 1.75 cycle/sample.

In general any digital frequency above 0.5 cycle/sample (π rad./sample) is an alias (or shows up as an alias of some base band frequency). It actually has more cycles per sample than is apparent in a plot of the sampled data.

We demonstrate below the phenomenon of aliasing using the three waveforms $x_1(t)$, $x_2(t)$, and $x_3(t)$ and the corresponding sequences $x_1(n)$, $x_2(n)$, and $x_3(n)$.

```
In MATLAB: %Aliasing demo
```

```
% First-----
% Plot the continuous-time waveforms x_1(t), x_2(t), and x_3(t) over a 1-second interval
t = 0: 1/500: 1;
x1 = \cos (2*pi*4*t); \%4 Hz
subplot(3, 1, 1), plot(t, x1, 'b');
                                    \$ subplot(3, 1, 1) – 3 rows, 1 column, #1
xlabel ('Time, t, seconds'), ylabel('x1(t)');
title ('4 Hz')
grid;
x^2 = \cos (2*pi*12*t); \% 12 Hz
subplot(3, 1, 2), plot(t, x2, 'k');
                                    % subplot(3, 1, 2) – 3 rows, 1 column, #2
xlabel ('Time, t, seconds'), ylabel('x2(t)');
title ('12 Hz')
x3 = \cos (2*pi*28*t); \%28 Hz
subplot(3, 1, 3), plot(t, x3, 'r');
                                    \% subplot(3, 1, 3) – 3 rows, 1 column, #3
xlabel ('Time, t, seconds'), ylabel('x3(t)');
title ('28 Hz')
%Second-----
% Plot the sequences x_1(n), x_2(n), and x_3(n)
n = 0: 1: 16;
%4 Hz sampled at 16 Hz
x1 = \cos(n*pi/2);
subplot(3, 1, 1), stem(n, x1, 'bo');
                                    \% subplot(3, 1, 1) – 3 rows, 1 column, #1
xlabel ('Sample number, n'), ylabel('x1(n)');
title ('4 Hz at 16 samples/sec')
grid;
%12 Hz sampled at 16 Hz
x^2 = \cos (3 n^* pi/2);
subplot(3, 1, 2), stem(n, x2, 'ko');
                                    % subplot(3, 1, 2) – 3 rows, 1 column, #2
xlabel ('Sample number, n'), ylabel('x2(n)');
title ('12 Hz at 16 samples/sec')
%28 Hz sampled at 16 Hz
x3 = \cos (7*n*pi/2);
subplot(3, 1, 3), stem(n, x3, 'ro'); %subplot(3, 1, 3) – 3 rows, 1 column, #3
xlabel ('Sample number, n'), ylabel('x3(n)');
title ('28 Hz at 16 samples/sec')
%_____
```



Discrete-time signals

Definition A **discrete-time signal** is a sequence, that is, a function defined on the positive and negative integers.

The sequence $x(n) = x_R(n) + j x_I(n)$ is a complex (valued) sequence if $x_I(n)$ is not zero for all *n*. Otherwise, it is a real (valued) sequence.

Examples of discrete-time signals represented in functional form are given below. (Exercise: Plot these signals for several values of n.)

$$x_1(n) = 2 \cos 3n$$

 $x_2(n) = 3 \sin (0.2\pi n)$

Alternatively, if a signal is non-zero over a finite (small enough) interval, we can list the values of the signal as the elements of a sequence. For example

$$x_3(n) = \{5, 2, -1, 1, -1/2, 4\}$$

The arrow indicates the value at n = 0. We omit the arrow when the first entry represents the value for n = 0. The above sequence is a *finite length sequence*. It is assumed that all values of the signal not listed are zero. In the above example x(0) = 1, x(1) = -1/2, x(-4) = x(3) = 0, etc.

Definition A discrete-time signal whose values are from a finite set is called a digital signal.

Classification of discrete-time signals (Along lines similar to continuous-time signals)

(*Omit*) *Discrete-time Energy and Power signals* The *energy* E of a discrete-time signal x(n) is given by

$$E = \lim_{N \to \infty} \sum_{n = -N}^{N} x(n) \ x^*(n)$$

where x^* is the complex conjugate of *x*. If x(n) is a real sequence then $x(n) x^*(n) = x^2(n)$. The above definition can also be written as

$$E = \lim_{N \to \infty} \sum_{n = -N}^{N} |x(n)|^2 \quad \text{(there are 2N+1 terms here)}$$

The *average power P* of the signal is

$$P = \lim_{N \to \infty} \frac{1}{2N + 1} \sum_{n = -N} x(n)^2 |$$

If *E* is finite but non zero (i.e., $0 < E < \infty$) the signal is an *energy signal*. It is a *power* signal if *E* is infinite but *P* is finite and nonzero (i.e., $0 < P < \infty$). Clearly, when *E* is finite, P = 0. If *E* is infinite *P* may or may not be finite.

If neither E nor P is finite, then the signal is neither an energy nor a power signal.

The terms "power" and "energy" are used independently of whether the

quantity $\sum_{n=-N} |\mathbf{x}(n)|^{\frac{1}{p}}$ (or, $\int_{t_1} |\mathbf{x}(t)|^{\frac{1}{p}} dt$, in the continuous time case) actually is related to physical

energy. Even if such a relationship exists, these quantities, Σ (.) and \int (.), may have the wrong dimensions and scaling. Still it is convenient to use these terms in a general fashion. It is helpful to imagine that x(t) or x(n) is the voltage across, or, current through, a 1-ohm resistor.

Example 1.2.1 For the signal x(n) = 1 for all n,

$$E = \lim_{N \to \infty} \underbrace{\sum_{n=-1}^{N} |x(n)|^2}_{N \to \infty} = \underbrace{\sum_{n=-\infty_2}^{\infty} 1^2}_{2n} = \infty \text{ which is infinite energy}_{2N+1}$$

$$P = \lim_{N \to \infty} \underbrace{\sum_{n=-N}^{N} |x(n)|}_{N \to \infty} = \lim_{N \to \infty} \underbrace{\sum_{n=-N}^{N} 1^2}_{N \to \infty} = 1 \text{ which is finite}$$

Thus x(n) is a power signal.

Example 1.2.2 For the signal x(n) = n both *E* and *P* are infinite. This is neither an energy nor a power signal. *(End of Omit)*

Examples of discrete-time signals In these examples w(n) and z(n) take on only a finite number of different values – hence digital. But x(n) and y(n) take on a countable infinite number of values – they are not digital. (Figure)



Important discrete-time signals If a continuous-time signal x(t) is sampled at *T*-second intervals, the result is the sequence $\{x(nT)\}$. For convenience we shall drop the *T* and the braces and use just x(n) to represent the sequence.

1) The unit sample sequence (discrete-time impulse, aka Kronecker delta)

 $\delta(n) = 1, \quad n = 0 \\ 0, \quad n \neq 0$

Whereas $\delta(n)$ is somewhat similar to the continuous-time impulse function $\delta(t)$ – the Dirac delta – we note that the magnitude of the discrete impulse is finite. Thus there are no analytical difficulties in defining $\delta(n)$. It is convenient to interpret the delta function as follows:

 δ (argument) = 1 when argument = 0 0 when argument \neq 0



2) The unit step sequence

 $u(n) = 1, n \ge 0$ 0, n < 0

 $u(\text{argument}) = 1, \text{ if argument} \ge 0$ 0, if argument < 0



a) The discrete delta function can be expressed as the first difference of the unit step function:

$$\delta(n) = u(n) - u(n-1)$$

b) The sum from $-\infty$ to *n* of the δ function gives the unit-step:



Results (a) and (b) are like the continuous-time derivative and integral respectively.

c) By inspection of the graph of u(n), shown below, we can write:

$$u(n) = \delta(n) + \delta(n-1) + \delta(n-2) + \ldots = \sum_{\lambda=0}^{\infty} \delta(n-\lambda)$$



d) For any arbitrary sequence x(n), we have

$$x(n) \ \delta(n-k) = x(k) \ \delta(n-k)$$

that is, the multiplication will pick out just the one value x(k).

If we find the infinite sum of the above we get the sifting property:

$$\sum_{n=-\infty} x(n) \,\,\delta(n-k) = x(k)$$



e) We can write *x*(*n*) as follows:

$$x(n) = \dots + x(-1) \,\delta(n+1) + x(0) \,\delta(n) + x(1) \,\delta(n-1) + x(2) \,\delta(n-2) + \dots$$

This can be verified to be true for all n by setting in turn

...,
$$n = -2$$
, $n = -1$, $n = 0$, $n = 1$, $n = 2$, etc. ...

The above can be written compactly as

$$x(n) = \sum_{k = -\infty}^{\infty} x(k) \, \delta(n-k)$$

This is a weighted-sum of delayed unit sample functions.

3) The real exponential sequence Consider the familiar continuous time signal

 $x(t) = e^{-\alpha t} = e^{-t/\tau}, \qquad t \ge 0$

The sampled version is given by setting t = nT

$$x(nT) = e^{-\alpha nT} = \left(e^{-\alpha T}\right)^n, \qquad nT \ge 0$$

Dropping the *T* from x(nT) and setting $e^{-\alpha T} = a$ we can write

$$x(n) = a^n$$
, $n \ge 0$

The sequence can also be defined for both positive and negative *n*, by simply writing $x(n) = a^n$ for all *n*.



4) The sinusoidal sequence Consider the continuous-time sinusoid x(t)

 $x(t) = A \sin 2\pi F_0 t = A \sin \Omega_0 t$

 F_0 and Ω_0 are the analog frequency in Hertz (or cycles per second) and radians per second, respectively. The sampled version is given by

$$x(nT) = A \sin 2\pi F_0 nT = A \sin \Omega_0 nT$$

We may drop the *T* from x(nT) and write

$$x(n) = A \sin 2\pi F_0 nT = A \sin \Omega_0 nT$$
, for all n

We may write $\Omega_0 T = \omega_0$ which is the digital frequency in radians (per sample), so that

$$x(n) = A \sin \omega_0 n = A \sin 2\pi f_0 n$$
, for all n

Setting $\omega_0 = 2\pi f_0$ gives $f_0 = \omega_0/2\pi$ which is the digital frequency in cycles per sample. In the analog domain the horizontal axis is calibrated in seconds; "second" is one unit of the independent variable, so Ω_0 and F_0 are in "per second". In the digital domain the horizontal axis is calibrated in samples; "sample" is one unit of the independent variable, so ω_0 and f_0 are in "per second".

Classification of discrete-time signals (cont'd)

Periodic signal The discrete-time signal x(n) is periodic if, for some integer N > 0

$$x(n+N) = x(n)$$
 for all n

The smallest value of N that satisfies this relation is the (**fundamental**) **period** of the signal. If there is no such integer N, then x(n) is an aperiodic signal.

Given that the continuous-time signal $x_a(t)$ is periodic, that is, $x_a(t) = x_a(t+T_0)$ for all t, and that x(n) is obtained by sampling $x_a(t)$ at T second intervals, x(n) will be periodic if T_0/T is a rational number but not otherwise. If $T_0/T = N/L$ for integers $N \ge 1$ and $L \ge 1$ then x(n) has exactly N samples in L periods of $x_a(t)$ and x(n) is periodic with period N.

Periodicity of sinusoidal sequences The sinusoidal sequence $\sin(2\pi f_0 n)$ has several major differences from the continuous-time sinusoid as follows:

a) The sinusoid $x(n) = \sin (2\pi f_0 n)$ or $\sin (\omega_0 n)$ is periodic if f_0 , that is, $\omega_0/2\pi$, is rational. If f_0 is not rational the sequence is not periodic. Replacing *n* with (n+N) we get

$$x(n+N) = \sin (2\pi f_0 (n+N)) = \sin 2\pi f_0 n. \cos 2\pi f_0 N + \cos 2\pi f_0 n. \sin 2\pi f_0 N$$

Clearly x(n+N) will be equal to x(n) if $f_0N = m$, an integer or $f_0 = m/N$. The fundamental period is obtained by choosing *m* as the smallest integer that yields an integer value for *N*. For example, if $f_0 = 15/25$, which in reduced fraction form is 3/5, then we can choose m = 3 and get N = 5 as the period. If f_0 is rational then $f_0 = p/q$ where *p* and *q* are integers. If p/q is in reduced fraction form then *q* is the period as in the above example.

On the other hand if f_0 is irrational, say $f_0 = \sqrt{2}$, then N will not be an integer, and thus x(n) is aperiodic.

Note: In expressions like sin ωn , sin $2\pi fn$, $e^{j\omega n}$ and $e^{j2\pi fn}$ we shall refer to ω or f as the frequency even when the signal concerned is not periodic by the definition above.

b) The sinusoidal sequences $\sin \omega_0 n$ and $\sin ((\omega_0 + 2\pi k)n)$ for $0 \le \omega_0 \le 2\pi$ are identical. This can be shown using the identity

$$\sin ((\omega_0 + 2\pi k)n) = \sin (\omega_0 n + 2\pi kn)$$

= sin \omega_0 n cos 2\pi kn + cos \omega_0 n sin 2\pi kn QED

Similarly, $\cos \omega_0 n$ and $\cos ((\omega_0 + 2\pi k)n)$ are the same. Therefore in considering sinusoidal sequences for analysis purposes, ω_0 can be restricted to the range $0 \le \omega_0 \le \pi$ without any loss of generality.

c) For $\pi < \omega_0 < 2\pi$, based on the same trigonometric identities,

sin $\omega_0 n$ is the negative of sin $((2\pi - \omega_0)n)$, and $\cos \omega_0 n$ is the same as $\cos ((2\pi - \omega_0)n)$

The sum of two discrete-time periodic sequences is also periodic. Let x(n) be the sum of two periodic sequences, $x_1(n)$ and $x_2(n)$, with periods N_1 and N_2 respectively. Let p and q be two integers such that

 $pN_1 = qN_2 = N$ (*p* and *q* can always be found)

Then x(n) is periodic with period N since, for all n,

$$x(n+N) = x_1(n+N) + x_2(n+N) = x_1(n+pN_1) + x_2(n+qN_2) = x_1(n) + x_2(n) = x(n) \text{ for all } n$$

Odd and even sequences The signal x(n) is an *even* sequence if x(n) = x(-n) for all n, and is an *odd* sequence if x(n) = -x(-n) for all n.



The *even* part of x(n) is determined as $x_e(n) = \frac{x(n) + x(-n)}{2}$ and the *odd* part of x(n) is given by $x_o(n) = \frac{x(n) - x(-n)}{2}$. The signal x(n) then is given by $x(n) = x_e(n) + x_o(n)$.

Example 1.2.3 Plot the sequences $x_1(n) = 2 \cos n$ and $x_2(n) = 2 \cos (0.2\pi n)$. What are their "frequencies"? \rightarrow hich of them is truly periodic and what is its periodicity? **Solution** The MATLAB program segment follows:

N = 21; n = 0: N-1; % %Nonperiodic x1= 2*cos(1*n); subplot(2, 1, 1), stem(n, x1); xlabel('n'), ylabel('x1'); title('x1 = 2 cos 1n'); % %Periodic x2 = 2*cos(0.2*pi*n); subplot(2, 1, 2), stem(n, x2); xlabel('n'), ylabel('x2'); title('x2 = 2 cos 0.2\pi n');



Example 1.2.4 Plot the sequence $x_3(n) = 2(0.9)^n \cos 0.2 \pi n$. **Solution** The MATLAB program segment follows:

n = [0: 30];% % ".^" stands for element-by-element exponentiation %".*" stands for element-by-element multiplication x3 = 2* ((0.9) .^n) .*cos(0.2*pi*n); stem(n, x3); xlabel('n'), ylabel('x3'); title('Sequence x3(n)');



Transformation of the independent variable

Shifting and folding (reflecting about the vertical axis) Given the sequence x(n), where *n* is the independent variable, we have the following two transformation operations:

- Shifting in time by *k* units, where *k* is an integer, is denoted by *x*(*n*–*k*). The sequence *x*(*n*–*k*) represents the sequence *x*(*n*) shifted by *k* samples, to the right if *k* is positive, or to the left if *k* is negative. Parenthesize *n* and replace it by (*n*–*k*) for *k* units of delay, or by (*n*+*k*) for *k* units of time advancement.
- Folding (a.k.a. time reversal or reflecting about the vertical axis) is denoted by x(-n). The signal x(-n) corresponds to reflecting x(n) about the time origin n = 0. Reverse the sign of *n* (replace *n* with -n).

As in the case of continuous-time signals the operations of shifting and folding are *not commutative*. In other words, the result of first shifting and then folding is not the same as that of first folding and then shifting.

A third operation is **scaling**, of which more will be said in a later unit.

Example 1.3.1 [Delta function] Given the delta function $\delta(n)$, first reflect then shift (delay) by 2 units. The other possibility is first to shift (delay) by 2 units and then reflect.

The result of "reflect then shift" is shown below left. *Reflect* means to change the sign of *n*; then *shifting* is done by replacing (*n*) by (*n*–2). The result is: (1) $\delta(n) < \delta(-n)$ and (2) $\delta(-n) = \delta(-(n-2)) = \delta(-(n+2))$.



To continue the example, the second possibility, the result of "shift then reflect" is shown above right. *Shift* means replacing (*n*) by (*n*-2); then *reflect* by changing the sign of *n*. The result is $\delta(-n-2)$. The result is: (1) $\delta(n) = \delta((n)) < \delta((n-2))$ and (2) $\delta((n-2)) = \delta(n-2) < \delta(-n-2)$.

Note that the two end results are not the same. This serves to illustrate that the two operations of shifting and reflecting are not commutative:

 $Fold(Shift(\delta(n))) \neq Shift(Fold(\delta(n)))$

Discrete-time systems

Definition A *discrete-time system* is a mapping from the set of acceptable discrete-time signals, called the input set, to a set of discrete-time signals called the output set.

Definition A discrete-time system is *deterministic* if its output to a given input does not depend upon some random phenomenon. If it does, the system is called a *random (stochastic) system*.

Definition A *digital system* is a mapping which assigns a digital output signal to every acceptable digital input signal.

A discrete-time system can be thought of as a transformation or operator, T, that maps an input sequence x(n) to an output sequence y(n) shown thus:



In what follows we focus on the presence or absence of the following properties in discrete-time systems: linearity, shift invariance, causality and stability.

Filter Some refer to a linear time-invariant (LTI) system simply as a **filter**, that is, a filter is a system *T* with a single input and a single output signal that is both linear and time-invariant.

Linearity

Definition A discrete-time system T[.] is linear if the response to a weighted sum of inputs $x_1(n)$ and $x_2(n)$ is a weighted sum (with the same weights) of the responses of the inputs separately for all weights and all acceptable inputs. Thus the system y(n) = T[x(n)] is linear if for all a_1 , a_2 , $x_1(n)$ and $x_2(n)$ we have

 $T[a_1x_1(n) + a_2x_2(n)] = a_1T[x_1(n)] + a_2T[x_2(n)]$

Another way of saying this is that if the inputs $x_1(n)$ and $x_2(n)$ produce the outputs $y_1(n)$ and $y_2(n)$, respectively, then the input $a_1x_1(n) + a_2x_2(n)$ produces the output $a_1 y_1(n) + a_2 y_2(n)$. This is called the **superposition principle**. The a_1 , a_2 , $x_1(n)$ and $x_2(n)$ may be complex-valued. The above definition combines two properties, viz.,

- 1. **Additivity**, that is, $T[x_1(n)+x_2(n)] = T[x_1(n)] + T[x_2(n)]$, and
- 2. Scaling (or homogeneity), that is, T[c x(n)] = c T[x(n)]

The procedure of checking for linearity is:

- *1.* Find outputs $y_1(n)$ and $y_2(n)$ corresponding to inputs $x_1(n)$ and $x_2(n)$
- 2. Form the sum $a_1 y_1(n) + a_2 y_2(n)$
- 3. Find output $y_3(n)$ corresponding to input $a_1x_1(n) + a_2x_2(n)$
- 4. Compare the results of steps 2 and 3

Examples of linear systems:

1. y(n) = x(n) + x(n-1) + x(n-2)2. y(n) = y(n-1) + x(n)3. y(n) = 04. y(n) = n x(n) (But time-varying)

Examples of nonlinear systems:

y(n) = x²(n)
 y(n) = 2 x(n)+3. This is a *linear equation* though! This system is made up of a linear part, 2 x(n), and a zero-input response, 3. This is called an *incrementally linear system*, for it responds linearly to *changes* in the input.

Example 1.4.1 Determine if the system y(n) = T[x(n)] = x(-n) is linear or nonlinear.



Answer Determine the outputs $y_1(.)$ and $y_2(.)$ corresponding to the two input sequences $x_1(n)$ and $x_2(n)$ and form the weighted sum of outputs:

$$y_1(n) = T[x_1(n)] = x_1(-n)$$

 $y_2(n) = T[x_2(n)] = x_2(-n)$

The weighted sum of outputs = $a_1 x_1(-n) + a_2 x_2(-n) < (A)$.

Next determine the output y_3 due to a weighted sum of inputs:

$$y_3(n) = T[a_1 x_1(n) + a_2 x_2(n)] = a_1 x_1(-n) + a_2 x_2(-n) < (B)$$

Check if (A) and (B) are equal. In this case (A) and (B) are equal; hence the system is linear.

Example 1.4.2 Examine y(n) = T[x(n)] = x(n) + n x(n+1) for linearity.



Answer The outputs due to $x_1(n)$ and $x_2(n)$ are:

$$y_1(n) = T[x_1(n)] = x_1(n) + n x_1(n+1)$$

$$y_2(n) = T[x_2(n)] = x_2(n) + n x_2(n+1)$$

The weighted sum of outputs = $a_1 x_1(n) + a_1 n x_1(n+1) + a_2 x_2(n) + a_2 n x_2(n+1) < (A)$ The output due to a weighted sum of inputs is

$$y_3(n) = T[a_1 x_1(n) + a_2 x_2(n)]$$

= $a_1 x_1(n) + a_2 x_2(n) + n (a_1 x_1(n+1) + a_2 x_2(n+1))$
= $a_1 x_1(n) + a_2 x_2(n) + n a_1 x_1(n+1) + n a_2 x_2(n+1) < (B)$

Since (A) and (B) are equal the system is linear.

Example 1.4.3 Check the system $y(n) = T[x(n)] = ne^{|x(n)|}$ for linearity.

$$x(n) \qquad \qquad y(n) = T[x(n)] = n e^{\int x(n)}$$

Answer The outputs due $x_1(n)$ and $x_2(n)$ are:

$$y_{1}(n) = T[x_{1}(n)] = n e^{|x_{1}(n)|}$$

$$y_{2}(n) = T[x_{2}(n)] = n e^{|x_{2}(n)|}$$

The weighted sum of the outputs $= a_1 n e^{|x_1(n)|} \pm a n e^{|x_2(n)|} \prec (A)$ The output due to a weighted sum of inputs is

$$y_3(n) = T[a_1 x_1(n) + a_2 x_2(n)] = n e^{a_1 x_1(n) + a_2 x_2(n)} < (B)$$

We can specify a_1 , a_2 , $x_1(n)$, $x_2(n)$ such that (A) and (B) are not equal. Hence nonlinear.

Example 1.4.4 Check the system $y(n) = T[x(n)] = a^n \cos(2\pi n/N)$ for linearity.

$$x(n) \qquad \qquad y(n) = a^n \cos(2\pi n/N)$$

Answer Note that the input is x(n). Clearly y(n) is independent of x(n). The outputs due to $x_I(n)$ and $x_2(n)$ are:

$$y_1(n) = T[x_1(n)] = a^n \cos(2\pi n/N)$$

 $y_2(n) = T[x_2(n)] = a^n \cos(2\pi n/N)$

The weighted sum of the outputs = $b_1 a^n \cos (2\pi n/N) + b_2 a^n \cos (2\pi n/N) < (A)$

The output due to a weighted sum of inputs is

 $y_3(n) = T[b_1 x_1(n) + b_2 x_2(n)] = a^n \cos((2\pi n/N)) < (B)$

(A) and (B) are not equal, so the system is not linear. (But (A) = $(b_1+b_2) a^n \cos((2\pi n/N))$ and this is equal to (B) within a constant scaling factor.)

Example 1.4.5 Check the system y(n) = T[x(n)] = n x(n) for linearity.

Answer For the two arbitrary inputs $x_1(n)$ and $x_2(n)$ the outputs are $y_1(n) = T[x_1(n)] = n x_1(n)$ $y_2(n) = T[x_2(n)] = n x_2(n)$ For the weighted sum of inputs $a_1 x_1(n) + a_2 x_2(n)$ the output is $y_3(n) = T[a_1 x_1(n) + a_2 x_2(n)] = n (a_1 x_1(n) + a_2 x_2(n))$ $= a_1 n x_1(n) + a_2 n x_2(n)$ $= a_1 y_1(n) + a_2 y_2(n)$. Hence the system is linear.

Example 1.4.6 Check $y(n) = T[x(n)] = x^2(n)$ for linearity. **Answer** For $x_1(n)$, the output is $y_1(n) = x_1^2(n)$ and for $x_2(n)$, $y_2(n) = x_2^2(n)$. The weighted sum of outputs $= a_1 x^2(n) + a_2 x^2(n) < (A)$

The output due to a weighted sum of inputs is

$$y_3(n) = T[a_1 x_1(n) + a_2 x_2(n)] = a_1^2 x_1^2(n) + a_2^2 a_2^2 x_2^2(n) + 2 a_1 a_2 x_1(n) x_2(n) < (B)$$

It is possible to specify a_1 , a_2 , $x_1(n)$, and $x_2(n)$ so that (A) and (B) are not equal. Hence the system is nonlinear.

Example 1.4.7 Determine if y(n) = T[x(n)] = 2 x(n) + 3 is linear. **Answer** For $x_1(n)$, the output is $y_1(n) = 2 x_1(n) + 3$ and for $x_2(n)$, it is $y_2(n) = 2 x_2(n) + 3$. Weighted sum of outputs $= a_1 (2 x_1(n) + 3) + a_2 (2 x_2(n) + 3)$ $= 2 a_1 x_1(n) + 3 a_1 + 2 a_2 x_2(n) + 3 a_2 < (A)$

The output due to a weighted sum of inputs is

 $y_3(n) = T[a_1 x_1(n) + a_2 x_2(n)] = 2 (a_1 x_1(n) + a_2 x_2(n)) + 3$

 $= 2 a_1 x_1(n) + 2 a_2 x_2(n) + 3 < (B)$

It is possible to specify a_1 , a_2 , $x_1(n)$, and $x_2(n)$ such that (A) and (B) are unequal. Hence the system is nonlinear.

(Since the output y(n) = 3 if x(n) = 0 we see that the system violates the "zero-in / zero-out" property of linear systems. Hence nonlinear.)

Example 1.4.8 Determine if $y(n) = T[x(n)] = \text{Re}\{x(n)\}$ is linear.

Answer Note that the input signals as well as the scaling constants a_1 and a_2 are allowed to be *complex*. Let $x_1(n) = r(n) + j s(n)$. Then $y_1(n) = Re\{x_1(n)\} = r(n)$.

The scaling / homogeneity property says that, if the response to $x_1(n)$ is $y_1(n)$, then the response to $x_2(n) = a.x_1(n)$ is $a.y_1(n)$ where a is *any* constant. Let a = j, and choose $x_2(n) = a.x_1(n) = j.x_1(n) = j.(r(n) + j s(n)) = -s(n) + j r(n)$. The corresponding output is $y_2(n) = \text{Re}\{x_2(n)\} = -s(n)$. Thus if $y_1(n) = T[.] = T[x_1(n)] = r(n)$, then

$$y_2(n) = T[.] = T[x_2(n)] = T[j x_1(n)] = -s(n).$$

This is not equal to j r(n) as would be expected from the homogeneity property. Hence the system is nonlinear.

Shift-invariance (time-invariance)

Definition A discrete time system y(n) = T[x(n)] is shift-invariant if, for all x(n) and all n_0 , we have: $T[x(n-n_0)] = y(n-n_0)$.

This means that applying a time delay (or advance) to the input of a system is equivalent to applying it to the output.



procedure for determining shift-

invariance is:

The

Step 1. Determine output y(n) corresponding to input x(n).

Step 2. Delay the output y(n) by n_0 units, resulting in $y(n-n_0)$.

Step 3. Determine output $y(n, n_0)$ corresponding to input $x(n-n_0)$.

Step 4. Determine if $y(n, n_0) = y(n-n_0)$. If equal, then the system is shift-invariant; otherwise it is time-varying.

When we suspect that the system is time-varying a very useful alternative approach is to find a counter-example to disprove time-invariance, i.e., use intuition to find an input signal for

which the condition of shift-invariance is violated and that suffices to show that a system is not shift-invariant.

Example 1.4.9 Test if y(n) = T[x(n)] = x(-n) is shift-invariant.



Answer Find output for x(n), delay it by n_0 , and compare with the output for $x(n-n_0)$. The output for x(n) is

$$\mathbf{y}(\mathbf{n}) = T[\mathbf{x}(n)] = \mathbf{x}(-\mathbf{n})$$

Delaying y(n) by n_0 gives

 $y(n-n_0) = x(-(n-n_0)) = x(-n+n_0) < (A)$

As an aside this amounts to reflecting first and then shifting.

The output for $x(n-n_0)$ is denoted $y(n, n_0)$ and is given by

 $y(n, n_0) = T[x(n-n_0)] = x(-n-n_0) < (B)$

As an aside this amounts to shifting first and then reflecting.

(A) and (B) are not equal. That is, $y(n, n_0) \neq y(n-n_0)$, so the system is time-varying.

Example 1.4.10 Examine y(n) = T[x(n)] = x(n) + n x(n+1) for time invariance. **Answer** Notice that the difference equation has a **time-varying coefficient**, n. The output y(n) corresponding to x(n) is already given above. Delaying y(n) by n_0 gives

 $y(n-n_0) = x(n-n_0) + (n-n_0) x(n-n_0+1) < (A)$

Compare with $y(n, n_0) = T[x(n-n_0)] = x(n-n_0) + n x(n-n_0+1) < (B)$

(A) \neq (B), so the system is time varying.

Example 1.4.11 Check for time invariance of the system y(n) = T[x(n)] = n x(n). **Answer** We shall do this by counterexample(s) as well as by the formal procedure. The formal procedure is:

$$\mathbf{y}(\mathbf{n}) = T[\mathbf{x}(n)] = \mathbf{n} \mathbf{x}(\mathbf{n})$$

Delay this by n_0 to get $y(n-n_0) = (n-n_0) x(n-n_0) < (A)$

Compare with $y(n, n_0) = T[x(n-n_0)] = nx(n-n_0) < (B)$

Since $(A) \neq (B)$, the system is time-varying.

Alternative We expect that it is time varying since the equation has a time varying coefficient. Find a counter example to show that the system is time varying. For input $x(n) = \delta(n)$, the output is

$$y(n) = n \delta(n) = 0$$
 for all n

For input $x(n-1) = \delta(n-1)$, the output is

$$y(n, 1) = n \delta(n-1) = 1 \delta(n-1)$$

Thus while x(n-1) is a shifted version of x(n), y(n, 1) is not a shifted version of y(n). So the system is time-varying.

Another counter-example If x(n) = u(n), then

$$y(n) = n u(n)$$

But if the input is x(n-2) = u(n-2), then

y(n, 2) = n u(n-2) = (n-2+2) u(n-2) = (n-2) u(n-2) + 2 u(n-2)

Delayed version Extra term The extra term shows that $y(n, 2) \neq y(n-2)$. So the system is time-varying.

Example 1.4.12 Check for time invariance the system $y(n) = T[x(n)] = \cos (x(n))$. **Answer** The output y(n) corresponding to input x(n) is

 $y(n) = T[x(n)] = \cos(x(n))$

Delay this output by n_0 to get $y(n-n_0) = \cos(x(n-n_0)) < (A)$

For the input $x(n-n_0)$ the output is $y(n, n_0) = T[x(n-n_0)] = \cos(x(n-n_0)) < (B)$

Since (A) and (B) are equal we have $y(n, n_0) = y(n-n_0)$. Therefore the system is time-invariant.

Example 1.4.13 The system y(n) = T[x(n)] = g(n) x(n) needs to be tested for time-invariance. **Answer** Note that the coefficient g(n) is time-varying. Hence, the system is time-varying. The output y(n) due to input x(n) is

 $\mathbf{y}(\mathbf{n}) = T[\mathbf{x}(n)] = \mathbf{g}(\mathbf{n}) \mathbf{x}(\mathbf{n})$

Delay this by n_0 to get $y(n-n_0) = g(n-n_0) x(n-n_0) \prec (A)$

The output $y(n, n_0)$ corresponding to $x(n-n_0)$ is given by

$$y(n, n_0) = T[x(n-n_0)] = g(n) x(n-n_0) < (B)$$

(A) \neq (B), i.e., y(n, n_0) \neq y(n-n_0). Hence, the system is time-varying.

Example 1.4.14 Check for time–invariance the system $y(n) = a^n \cos (2\pi n/N)$. **Answer** The output consists simply of a time-varying coefficient and is independent of the input x(n). The output y(n) due to input x(n) is $y(n) = T[x(n)] = a^n \cos(2\pi n/N)$

Delay this by n_0 to get $y(n-n_0) = a^{n-n_0} \cos(2\pi(n-n_0)/N) < (A)$

The output $y(n, n_0)$ due to input $x(n-n_0)$ is given by

 $y(n, n_0) = T[x(n-n_0)] = a^n \cos(2\pi n/N) < (B)$

(A) \neq (B), that is, y(n, n₀) \neq y(n-n₀). Hence, the system is time-varying.

Example 1.4.15 Check the system $y(n) = n e^{|x(n)|}$ for time-invariance. **Answer** Note time-varying coefficient, *n*. The output y(n) due to input x(n) is

$$y(n) = T[x(n)] = n e^{\int x(n)}$$

Delay this by n_0 to get $y(n-n_0) = (n-n_0)e^{ix(n-n_0)} < (A)$

The output $y(n, n_0)$ due to input $x(n-n_0)$ is given by

 $y(n, n_0) = T[x(n-n_0)] = ne^{x(n-n_0)} < (B)$

(A) and (B) are not equal. Hence, the system is time-varying.

Example 1.4.16 Test the system y(n) = x(2n) for time-invariance.



Answer This system represents *time scaling*. That is, y(n) is a time-compressed version of x(n), compressed by a factor of 2. For example, the value of x that occurred at 2n is occurring at n in the case of y. Intuitively, then, any time shift in the input will also be compressed by a factor of 2, and it is for this reason that the system is not time–invariant.

This is demonstrated by counter-example (Oppenheim & Willsky, p. 52). It can also be shown by following the formal procedure, which we shall do first below. For the input x(n) the output is

$$y(n) = T[x(n)] = x(2n)$$

Delay this output by n_0 to get

$$y(n-n_0) = x(2(n-n_0)) = x(2n-2n_0) < (A)$$

Next, for the input $x(n) = x(n-n_0)$ the output is

$$y(n, n_0) = x(2n) = x(2n-n_0) < (B)$$

(A) and (B) are not equal. So the system is *not* time-variant.

By counter example To show that the system y(n) = x(2n) is *not* time–invariant by way of a counter example consider the x(n) below:

$$x(n) = 1, \qquad -2 \le n \le 2$$

0, otherwise

We shall show that $y(n, 2) \neq y(n-2)$. We express x(n) in terms of unit step functions as

$$x(n) = u(n+2) - u(n-3)$$

We determine y(n-2) by first obtaining y(n) and then delaying it by 2 units:

$$y(n) < y(n) = x(2n) = u(2n+2) - u(2n-3)$$

Delay $< y(n-2) = x(2(n-2)) = u(2(n-2)+2) - u(2(n-2)-3) = u(2n-2) - u(2n-7)$

Next we determine y(n, 2) by first obtaining x(n-2) and then the corresponding output y(n, 2): $x(n-2) < x(n-2) = u((n-2)+2) - u((n-2)-3) = u(n) - u(n-5) = x_2(n)$, say $y(n, 2) < y(n, 2) = x_2(2n) = u(2n) - u(2n-5)$

We can easily sketch y(n, 2) and y(n-2) and see that $y(n, 2) \neq y(n-2)$, and therefore the system y(n) = x(2n) is *not* time–invariant.

Alternatively, this can be done entirely graphically.

Application of linearity – Convolution

An arbitrary sequence, x(n), can be written as the weighted sum of delayed unit sample functions:

$$x(n) = \dots + x(-2) \,\delta(n+2) + x(-1) \,\delta(n+1) + x(0) \,\delta(n) + x(1) \,\delta(n-1) + \dots$$
$$= \sum_{k = -\infty}^{\infty} x(k) \,\delta(n-k)$$

So the response of a linear system to input x(n) can be written down using the linearity principle, i.e., linear superposition. For a linear shift-invariant system whose impulse response is $T[\delta(n)] = h(n)$ the reasoning goes like this

- For an input $\delta(n)$ the output is h(n). For an input $x(0) \delta(n)$ the output is x(0) h(n) by virtue of scaling.
- For an input $\delta(n-1)$ the output is h(n-1) by virtue of shift-invariance. For an input $x(1) \delta(n-1)$ the output is x(1) h(n-1) by virtue of scaling.
- Therefore for an input of $x(0) \delta(n) + x(1) \delta(n-1)$ the output is x(0) h(n) + x(1)h(n-1) by virtue of additivity.

This reasoning can be extended to cover all the terms that make up x(n). In general the response to $x(k) \delta(n-k)$ is given by x(k) h(n-k).



Given that

$$h(n) = T[\delta(n)],$$
 and $x(n) = \sum_{k = -\infty} x(k) \ \delta(n-k)$

s

we have

$$y(n) = T[x(n)] = T \left[\sum_{k = -\infty}^{\infty} x(k) \, \delta(n-k) \right]$$

Since *T*[.] is linear we can apply linearity a countable infinite number of times to write

$$y(n) = \sum_{k = -\infty} T[x(k)\delta(n-k)] = \sum_{k = -\infty} x(k)T[\delta(n-k)]$$

In above equation since the system is shift-invariant we write $T[\delta(n-k)] = h(n-k)$. Else write $h_k(n)$ or h(n, k) in place of h(n-k). Thus for a linear shift-invariant system

$$y(n) = \sum_{k=-\infty}^{\infty} x(k)h(n-k)$$

Note that if the system is not specified to be shift-invariant we would leave the above result in the form

$$y(n) = \sum_{k=-\infty}^{\infty} x(k)h(n, k)$$
 or $y(n) = \sum_{k=-\infty}^{\infty} x(k)h_k(n)$

Then if shift-invariance is invoked we replace h(n, k) with h(n-k).

As in the case of continuous-time systems, the impulse response, h(n), is determined assuming that the system has no initial energy; otherwise the linearity property does not hold, so that y(n), as determined using the above equation, corresponds to only the forced response of the system.

The sum $\sum_{k=-\infty} x(k)h(n, k)$ is called the **convolution sum**, and is denoted x(n) * h(n).

A discrete-time linear shift-invariant system is completely characterized by its unit sample response h(n).

Theorem If a discrete-time system linear shift-invariant, T[.], has the unit sample response $T[\delta(n)] = h(n)$ then the output y(n) corresponding to any input x(n) is given by

$$y(n) = \sum_{k = -\infty} x(k) h(n - k) = \sum_{k = -\infty} x(n - k) h(k)$$

= $x(n) * h(n) = h(n) * x(n)$

The second summation is obtained by setting m = n-k; then for $k = -\infty$ we have $m = +\infty$, and for $k = \infty$ we have $m = -\infty$. Thus

$$\sum_{k=-\infty}^{\infty} x(k) h(n-k) = \underbrace{\sum_{m=\infty}^{-\infty} x(n-m) h(m)}_{k=-\infty} = \sum_{k=-\infty}^{\infty} x(n-k) h(k)$$

m is a dummy variable. The order of summation (forward or backward) makes no difference. Hence change m to k and switch limits **Example 1.4.17** [Linear Convolution] Given the input $\{x(n)\} = \{1, 2, 3, 1\}$ and the unit sample response $\{h(n)\} = \{4, 3, 2, 1\}$ find the response y(n) = x(n) * h(n).

Answer Since x(k) = 0 for k < 0 and h(n - k) = 0 for k > n, the convolution sum becomes

$$y(n) = \sum_{k = -\infty} x(k) h(n-k) = \sum_{k = 0} x(k) h(n-k)$$

Now y(n) can be evaluated for various values of *n*; for example, setting n = 0 gives y(0). See table below. The product terms shown in **bold italics** need not be calculated; they are zero because the signal values involved are zero.

Linear Convolution of $\{x(n)\} = \{1, 2, 3, 1\}$ and $\{h(n)\} = \{4, 3, 2, 1\}$

$$\mathbf{y}(\mathbf{n}) = \sum_{k=0} x(k) \ h(n-k)$$

	0	
n = 0	$y(0) = \sum_{k=0}^{0} x(k) h(0-k)$	= x(0) h(0) = 1 . 4 = 4
n = 1	$y(1) = \sum_{k=0}^{1} x(k) h(1-k)$	= x(0) h(1) + x(1) h(0) = 1 . 3 + 2 . 4 = 11
n = 2	$y(2) = \sum_{k=0}^{2} x(k) h(2-k)$	= x(0) h(2) + x(1) h(1) + x(2) h(0) = 1 . 2 + 2 . 3 + 3 . 4 = 20
n = 3	$y(3) = \sum_{k=0}^{3} x(k) h(3-k)$	= x(0) h(3) + x(1) h(2) + x(2) h(1) + x(3) h(0) = 1 . 1 + 2 . 2 + 3 . 3 + 1 . 4 = 18
n = 4	$y(4) = \sum_{k=0}^{4} x(k) h(4-k)$	= x(0) h(4) + x(1) h(3) + x(2) h(2) + x(3) h(1) + x(4) h(0) = 1 . 0 + 2 . 1 + 3 . 2 + 1 . 3 + 0 . 4 = 11
n = 5	$y(5) = \sum_{k=0}^{5} x(k) h(5-k)$	= x(0) h(5) + x(1) h(4) + x(2) h(3) + x(3) h(2) + x(4) h(1) + x(5) h(0) = 3 . 1 + 1 . 2 = 5
n = 6	$y(6) = \sum_{k=0}^{6} x(k) h(6-k)$	= x(0) h(6) + x(1) h(5) + x(2) h(4) + x(3) h(3) + x(4) h(2) + x(5) h(1) + x(6) h(0) = 1 . 1 = 1
n = 7	$y(7) = \sum_{k=0}^{7} x(k) h(7-k)$	= x(0) h(7) + x(1) h(6) + x(2) h(5) + x(3) h(4) + x(4) h(3) + x(5) h(2) + x(6) h(1) + x(7) h(0)
y(n) =	0 for $n < 0$ and $n > 6$	‡Product terms in <i>bold italics</i> are zero.

Sketches of the two sequences x(n) and h(n) are shown below.



To do the convolution we need the sequences x(k) and h(n-k), k being the independent variable. Of these x(k) is simply x(n) with k replacing n, shown below.



The sequence h(-k) is the reflected version of h(k). If h(-k) is delayed by *n* samples we get h(-(k-n)) that is h(n-k), shown above.

For each value of *n* the sequences x(k) and h(n-k) are multiplied point by point and the products are added, yielding the value of y(.) for the corresponding *n*.
Tabular method The following tabular method uses the second of the two forms, viz., y(n) =

 $\sum h(k) x(n-k)$, for the convolution sum.

	1 adular method of linear convolution													
	k۲	-3	-2	-1	0	1	2	3	4	5	6	7		
	h(k) <				4	3	2	1						
n\$	x(k) <				1	2	3	1					n	y(n)
0	x(0 -k)	1	3	2	1								0	4
1	x(1 –k)		1	3	2	1							1	11
2	x(2 -k)			1	3	2	1						2	20
3	x(3 -k)				1	3	2	1					3	18
4	x(4 -k)					1	3	2	1				4	11
5	x(5 –k)						1	3	2	1			5	5
6	x(6 -k)							1	3	2	1		6	1
7	x(7 –k)								1	3	2	1	7	0
•	•												•	•

The sequence y(n) is shown plotted below. Note that L_y , the length of y(n), equals the sum of the lengths of x(n) and h(n) minus 1: $L_y = L_x + L_h - 1$.



MATLAB (Linear Convolution)

xn = [1, 2, 3, 1]; hn = [4, 3, 2, 1]; yn = conv(xn, hn) M = length(xn) + length(hn) - 1; m = 0: M-1;stem(m, yn)

The result is: $y_n = 4$ 11 20 18 11 5 1



Matrix-vector multiplication [Proakis, Study Guide, p. 12] When the sequences x(n) and h(n) are of finite duration (as in this case) their linear convolution can also be implemented by using matrix-vector multiplication. We show below the formulation using the form

$$y(n) = \sum h(k) x(n-k)$$

as in the tabular method above. We arrange h(.) and y(.) as vectors (the latter with a length $L_y = L_x + L_h - 1$) and the (Toeplitz) matrix X (of size L_y rows and L_h columns) formed from the various shifted versions, x(n-k):



Note that the Toeplitz matrix may be seen (except for the zeros) in the center part of the table given under "Tabular method". \rightarrow e then evaluate *y* as the product *X*. *h*:

$$\begin{bmatrix} y(0) \\ y(1) \\ y(2) \\ y(3) \\ y(3) \\ y(4) \\ y(5) \\ y(6) \\ y(6) \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 3 & 2 & 1 & 0 \\ 1 & 3 & 2 & 1 \\ 0 & 1 & 3 & 2 \\ 0 & 1 & 3 \\ 1 & 3 \\ 0 & 1 & 3 \\ 1$$

MATLAB has a function called *toeplitz* that generates the Toeplitz matrix. We could, of course, use the alternative form $\sum x(k) h(n - k)$ in which case we evaluate $y = H \cdot x$.

Example 1.4.18 [Linear Convolution] Given the finite length sequences below find $y(n) = x_1(n)^* x_2(n)$.

 $x_1(n) = 1, \quad 0 \le n \le 3,$ $0, \quad \text{otherwise}$ $x_2(n) = 1, \quad 0 \le n \le 3$

Solution The MATLAB segment follows:

%Define sequences n = 0:3; x1 = ones(size(n)), x2 = ones(size(n)), %Length of output M = length(x1) + length(x2) - 1; m = 0: M-1; yn = conv(x1,x2); stem(m,yn) xlabel('n'), ylabel('y(n)'), title('Convolution y(n) = x1(n) * x2(n)');



Example 1.4.19 Given $x_1(n) = n \{u(n+10) - u(n-20)\}$ and $x_2(n) = \cos(0.1\pi n) \{u(n) - u(n-30)\}$, find $y(n) = x_1(n) * x_2(n)$.

Solution The intervals over which the sequences are non-zero are defined below:

$$\begin{aligned} x_{I}(n) &= n, & N_{0} \leq n \leq N_{I}, & N_{0} = -10, N_{I} = 20 \\ 0, & \text{otherwise} \end{aligned}$$

$$\begin{aligned} x_{2}(n) &= \cos \left(0.1 \, \pi n \right), & N_{2} \leq n \leq N_{3}, & N_{2} = 0, & N_{3} = 30 \\ 0, & \text{otherwise} \end{aligned}$$

$$\begin{aligned} y(n) &= \text{Non-zero,} & N_{4} \leq n \leq N_{5}, & N_{4} = N_{0} + N_{2}, & N_{5} = N_{I} + N_{3} \\ 0, & \text{otherwise} \end{aligned}$$

For (the summation limits) N_4 and N_5 see DSP-HW. The MATLAB segment and plot are given below:



40 of 77

Causality

The constraints of linearity and time-invariance define a class of systems that is represented by the convolution sum. The additional constraints of stability and causality define a more restricted class of linear time-invariant systems of practical importance.

Definition A discrete-time system is **causal** if the output at $n = n_0$ depends only on the input for $n \le n_0$.

The word "causal" has to do with cause and effect; in other words, for the system to act up there must be an actual cause. A causal system does not anticipate future values of the input but only responds to actual, present, input. As a result, if two inputs to a causal system are identical up to some point in time n_0 the corresponding outputs must also be equal up to this same time. The synonyms of "causal" are "(physically) realizable" and "non-anticipatory".

We digress below to introduce memory-less versus dynamic systems and then resume with causality.

(Aside) Systems with and without memory A system is said to be memory-less or static if its output for each value of n is dependent only on the input at that same time but not on past or future inputs.

Examples of static systems

1. y(n) = x(n) < the identity system 2. $y(n) = a x(n) - x^2(n)$ 3. A resistor *R*: y(t) = R x(t) (*y*(*t*) is voltage and *x*(*t*) is current)

In many physical systems, memory is directly associated with storage of energy. A resistor has no storage of energy. But a circuit with capacitors and/or inductors has storage of energy and is a **dynamic system**, i.e., has **memory**. However, while storage of energy has to do with past inputs only, a static system is independent not only of *past* but also of *future* inputs.

Examples of systems with memory, i.e., dynamic systems:

1. $y(n) = \sum_{k=-\infty} x(k)$. This is an accumulator or summer. The output y(n) depends on

values of x(.) prior to n such as x(n-1) etc.

2
$$y(n) = x(n-1)$$
. This is a delay element

3 A capacitor C:
$$y(t) = \frac{1}{C} \int_{-\infty}^{\infty} x(\tau) d\tau$$
, $(y(.)$ is voltage and $x(.)$ is current).

(End of Aside)

Getting back to causality, *all memory-less systems are causal* since the output responds only to the current value of the input. In addition, some dynamic systems (such as the three listed above) are also causal.

An example of a noncausal system is y(n) = x(n) + x(n+1) since the output depends on a future value, x(n+1).

Although causal systems are of great importance, they are not the only systems that are of practical importance. For example, causality is not often an essential constraint in applications in which the *independent variable is not time*, such as in image processing. Moreover, in processing data that have been *recorded previously* (non real-time), as often happens with speech,

geophysical, or meteorological signals, to name a few, we are by no means constrained to causal

processing. As another example, in many applications, including *historical* stock market analysis and demographic studies, we may be interested in determining a slowly varying trend in data that also contain higher frequency fluctuations about that trend. In this case, a commonly used approach is to average data over an interval in order to smooth out the fluctuations and keep only the trend. An example of such a noncausal averaging system is

$$y(n) = \frac{1}{2M+1} \sum_{k=-M} x(k)$$

Definition A discrete-time sequence x(n) is called causal if it has zero values for n < 0, i.e., x(n) = 0 for n < 0.

Theorem A *linear shift-invariant system* with impulse response h(n) is causal if and only if h(n) is zero for n < 0.

Proof of the "If" part By convolution the output y(n) is given by

$$y(n) = \sum_{k=-\infty} x(k) h(n-k)$$

If $h(n) = 0$ for $n < 0$, then $h(n-k) = 0$ for $n-k < 0$ or $k > n$. So

$$y(n) = \sum_{k = -\infty}^{n} x(k) h(n-k) + \sum_{k = n+1}^{\infty} x(k) h(n-k) = 0$$

$$=\sum_{k=-\infty}^{\infty} x(k) h(n-k)$$

Thus y(n) at any time *n* is a weighted sum of the values of the input x(k) for $k \le n$, that is, only the present and past inputs. Therefore, the system is causal.

Proof of the "Only If" part This is proved by contradiction, that is, if h(n) is non zero for n < 0 then the system is noncausal. Let h(n) be nonzero for n < 0:

$$h(n) \text{ is non zero for } n < 0 \qquad < h(n-k) \text{ is non zero for } n-k < 0 \text{ or } k > n$$

$$< \text{ the } 2^{\text{nd}} \text{ sum above, } \sum_{k=n+1}^{\infty} x(k) h(n-k) \text{ , is non zero}$$

$$< y(n) \text{ then depends on } x(n+1) \text{ and other future terms}$$

$$< \text{ Hence, the system is noncausal}$$

Example 1.4.20 Check the causality of the system y(n) = x(-n). **Answer** If *n* is some positive value then y(n) depends only on past values of the input x(.). But if *n* is negative, say n = -2, then

y(-2) = x(-(-2)) = x(2), a future value of input

Hence the system is noncausal.

Example 1.4.21 Check the system $y(n) = x(n) \cos(n+1)$ for causality. **Answer** Note that x(.) is the input, not the $\cos(.)$. In this system, the output at any time *n* equals the input at that same time multiplied by a number that varies with time. We can write the equation as y(n) = g(n) x(n), where $g(n) = \cos(n+1)$ is a time-varying function. Only the current value of the input x(.) determines the current output y(.). Therefore the system is causal (and also memoryless).

Example 1.4.22 Check the system $y(n) = T[x(n)] = ne^{|x(n)|}$ for causality. Answer Causal.

Example 1.4.23 Check the system $y(n) = T[x(n)] = a^n \cos((2\pi n/N))$ for causality. Answer Causal.

Example 1.4.24 Check the system $y(n) = T[x(n)] = \cos(x(n))$ for causality. Answer Causal.

Example 1.4.25 Check the system y(n) = T[x(n)] = x(-n+2) for causality. **Answer** For $n \le 0$ the argument of *x*, viz., -n+2 will be ≥ 2 . For example, if n = 0 we have

y(0) = x(2), a future value of input

Hence the system is noncausal.

Example 1.4.26 Check the system $y(n) = T[x(n)] = \sum_{k=n_0}^{n} x(k)$ for causality. Answer Causal. Example 1.4.27 Check the system $y(n) = T[x(n)] = \sum_{k=n-n_0}^{n+n_0} x(k)$ for causality. Answer $y(n) = \sum_{k=n-n_0}^{n+n_0} x(k) = \left\{ \sum_{k=n-n_0}^{n} x(k) \right\} + x(n+1) + x(n+2) + \dots + x(n+n_0)$

Future terms

The system is noncausal since y(n) depends on future values of the input.

(Aside) A streamable system is one that is either causal or *can be made causal* by adding an overall delay. The system y(n) = x(n+1) is not causal but can be made so and is therefore streamable. But y(n) = x(-n) is not streamable since it can"t be made causal.

Bounded input bounded output stability

Definition A sequence x(n) is bounded if there exists a finite M such that |x(n)| < M for all n. (Note that, as expressed here, M is a bound for negative values of x(.) as well. Another way of writing this is -M < x(n) < M.)

As an exarcmpl e, the sequence $x(n) = [1 + \cos 5\pi n] u(n)$ is bounded with $|x(n)| \le 2$. The sequence $x(n) = \begin{vmatrix} (1 & h) \sin 10n \\ 1 + (0.8)^n \end{vmatrix} u(n)$ is unbounded.

Definition A discrete-time system is bounded input-bounded output (BIBO) stable if every bounded input sequence x(n) produces a bounded output sequence. That is, if $|x(n)| \le M < \infty$, then $|y(n)| \le L < \infty$.

BIBO stability theorem A *linear shift invariant system* with impulse response h(n) is bounded input-bounded output stable if and only if *S*, defined below, is finite.

$$S = \sum_{k = -\infty}^{\infty} \left| h(k) \right| < \infty$$

i.e., the unit sample response is *absolutely summable*.

Proof of the "If" part Given a system with impulse response h(n), let x(n) be such that $|x(n)| \le M$. Then the output y(n) is given by the convolution sum:

$$y(n) = \sum_{k = -\infty} h(k) x(n-k)$$

so that

$$|y(n)| = \left| \sum_{k = -\infty}^{\infty} h(k) x(n-k) \right|$$

Using the triangular inequality that the sum of the magnitudes \geq the magnitude of the sum, we get

$$|y(n)| \leq \sum_{k=-\infty}^{\infty} h(k)x(n-k)$$

Using the fact that the magnitude of a product is the product of the magnitudes,

$$|y(n)| \le \sum_{k = -\infty} \dot{h}(k) x(n-k)$$
$$\le M \sum_{k = -\infty}^{\infty} |h(k)|$$

Thus, a sufficient condition for the system to be stable is that the unit sample response must be absolutely summable; that is,

$$\sum_{k=-\infty} |h(k)| < \infty \qquad \text{QED}$$

Proof of the "Only If" part That it is also a necessary condition can be seen by considering as input, the following bounded signal (this is the signum function),

$$x(k) = \operatorname{sgn} [h(n-k)] = \begin{array}{c} 1 & \text{where } h(n-k) > 0 \\ 0 & \text{where } h(n-k) = 0 \\ -1 & \text{where } h(n-k) < 0 \end{array}$$

Or, equivalently,

$$x(n-k) = \operatorname{sgn} [h(k)] = \begin{array}{c} 1 & \text{where } h(k) > 0 \\ 0 & \text{where } h(k) = 0 \\ -1 & \text{where } h(k) < 0 \end{array}$$

In the above we have implied that M = 1 (since *M* is some arbitrary finite number), and that $|x(n)| \le 1$. If, however, $|x(n)| \le M$ where *M* is finite but not equal to 1 we then will multiply the signum function by *M*. In either case x(n) is a bounded input. Thus

$$y(n) = \sum_{k = -\infty}^{\infty} h(k) x(n-k) = \sum_{k = -\infty}^{\infty} h(k) \operatorname{sgn}[h(k)] = \sum_{k = -\infty}^{\infty} h(k) |_{k=-\infty}^{\infty}$$

Clearly, if h(n) is not absolutely summable, y(n) will be unbounded. For a causal system the BIBO stability condition becomes QED

$$\sum_{k=0}^{\infty} \left| h(k) \right| < \infty$$

Example 1.4.28 Evaluate the stability of the linear shift-invariant system with the unit sample response $h(n) = a^n u(n)$.

Answer Evaluate

$$S = \sum_{k=-\infty}^{\infty} h(k) \models \sum_{k=-\infty}^{\infty} a^{k} u(k) \neq \sum_{k=0}^{\infty} a^{k} \models \left| \sum_{k=0}^{\infty} a^{k} \right|$$

Here we have used the fact that the magnitude of a product $(|a^k|)$ is the product of the magnitudes $(|a|^k)$. The summation on the right converges if |a| < 1 so that *S* is finite,

$$S = \frac{1}{1 - |a|}$$

and the system is BIBO stable.

Example 1.4.29 Evaluate the stability of the system y(n) = T[x(n)] = n x(n). **Answer** (Note that this is not a shift-invariant system, though it is linear as we determined elsewhere. The BIBO stability theorem applies to LSI systems. Note that if $x(n) = \delta(n)$ we get the result y(n) = h(n) = 0, for all n!)

If a system is suspected to be unstable, a useful strategy to verify this is to look for a specific bounded input that leads to an unbounded output. One such bounded input is x(n) = u(n), the unit step sequence. The output then is y(n) = nu(n), and as $n < \infty$, y(n) grows without bound. Hence the system is BIBO unstable.

Example 1.4.30 Examine $y(n) = T[x(n)] = e^{x(n)}$ for stability. **Answer** This system is nonlinear as determined elsewhere. Its unit sample response is given by setting $x(n) = \delta(n)$ and is given below

$$h(n) = e, \quad n = 0$$

1, $n \neq 0$

Here we are unable to find a bounded input that results in an unbounded output. So we proceed to verify that all bounded inputs result in bounded outputs. Let x(n) be an arbitrary signal bounded by B an arbitrary positive number,

|x(n)| < B or -B < x(n) < B for all n

Then $y(n) = e^{x(n)}$ must satisfy the condition $e^{-B} < y(n) < e^{B}$. Or, the output is guaranteed to be bounded by e^{B} . Thus the system is BIBO stable.

Example 1.4.31 Check for stability the system $y(n) = a^n \cos(2\pi n/N)$ Answer The output is independent of the input. The system is stable for $|a| \le 1$, otherwise it is unstable.

Example 1.4.32 $y(n) = x^2(n)$ Answer For any |x(n)| < B, we have $y(n) = x^2(n) < B^2$. Hence the system is stable.

Convolution - Properties

45 of 77

Property 1 The convolution operation is **commutative**, that is, x(n) * h(n) = h(n) * x(n).

A linear shift-invariant system with input x(n) and unit sample response h(n) will have the same output as a linear shift-invariant system with input h(n) and unit sample response x(n).

Property 2 The convolution of h(n) with $\delta(n)$ is, by definition, equal to h(n).

That is, the convolution of any function with the δ function gives back the original function. Stated another way: The **identity sequence** for the convolution operator is the unit sample, or

$$x(n) * \delta(n) = \delta(n) * x(n) = x(n)$$

Property 3 The convolution of a delayed unit sample sequence with x(n)

$$x(n) * \delta(n-k) = x(n-k)$$

Property 4 Associativity

$$x(n) * (y(n) * z(n)) = (x(n) * y(n)) * z(n) = x(n) * y(n) * z(n)$$

Property 5 Distributivity over sequence addition.

$$x(n) * (y(n) + z(n)) = x(n) * y(n) + x(n) * z(n)$$

Block diagram manipulation The properties of commutativity, associativity and distributivity enable us to determine the impulse response of series or parallel combinations of systems in terms of their individual impulse responses, as below:

(1) The following three LSI systems have identical unit sample response:



(2) The following two LSI systems are equivalent:



(3) A combination of the above two is the series parallel combination:



Convolution – Overlap-and-add See Unit II.

Linear constant coefficient difference equations

A subclass of linear shift-invariant systems is that for which the input x(n) and output y(n) satisfy an N^{th} order linear constant coefficient difference equation, given by

$$a_0 y(n) + a_1 y(n-1) + \dots + a_N y(n-N) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M), \qquad a_0 \neq 0$$

This may be written in the more compact, though daunting, form

$$\sum_{k=0}^{N} a_k y(n-k) = \sum_{r=0}^{M} b_r x(n-r), \qquad a_0 \neq 0$$

If the system is causal, then we can rearrange the above equation to compute y(n) iteratively from the present input x(n) and the past inputs x(n-1), x(n-2), ..., x(n-M) and the past outputs y(n-1), y(n-2), y(n-2), y(n-2).

$$y(n) = -\sum_{k=1}^{\infty} \left| \frac{a_k}{a^0} \right| y(n-k) + \dots$$
 (

If we think of the input as beginning at n = 0, then y(n) can be computed for all $n \ge 0$ once y(-1), y(-2), ..., y(-N) are specified. This is an iterative solution.

(*Aside*) Some write the difference equation with the terms y(n-1) through y(n-N) on the right hand side with positive (symbolic) coefficients and the coefficient of y(n) = 1, thus

$$y(n) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M) + a_1 y(n-1) + \dots + a_N y(n-N)$$

If we need to use this form we shall use different symbols altogether as follows:

$$y(n) = \beta_0 x(n) + \beta_1 x(n-1) + \dots + \beta_M x(n-M) + \alpha_1 y(n-1) + \dots + \alpha_N y(n-N)$$

(End of Aside)

Example 1.5.1 [First order system] For the first order system y(n) - a y(n-1) = x(n) find the output sequence y(n) assuming y(n) = 0 for all n < 0 and $x(n) = \delta(n)$. This corresponds to calculating the impulse response assuming zero initial conditions.

$$y(n) = a \ y(n-1) + x(n) = a \ y(n-1) + \delta(n)$$

$$n = 0: \ y(0) = a \ y(-1) + \delta(0) = a. \ 0 + 1 = 1$$

$$n = 1: \ y(1) = a \ y(0) + \delta(1) = a. \ 1 + 0 = a$$

$$n = 2: \ y(2) = a \ y(1) + \delta(2) = a. \ a + 0 = a^2$$

$$n = 3 \qquad \dots$$

Continuing the process we have $y(n) = a^n$, $n \ge 0$. This is also the unit sample response $h(n) = a^n u(n)$. It is not always possible to express y(n) as an analytical expression (closed form) as above.

Note It is also possible to recast the above problem as a noncausal or negative-time system with y(n) = 0 for $n \ge 0$. In this case, solving for y(n-1), we have

$$y(n-1) = \frac{1}{a} (y(n) - x(n))$$

This can be recast (by letting (n-1) = m etc.) as

$$y(n) = \frac{1}{a} (y(n+1) - x(n+1)), n < 0$$

The solution now is $y(n) = -a^n$, for n < 0, or the impulse response is $h(n) = -a^n u(-n-1)$.

Note Unless stated otherwise we shall generally assume that the difference equation represents a causal system.

Other techniques for solving difference equations Among other techniques is a method paralleling the procedure for solving linear constant coefficient differential equations, which involves finding and combining the **particular** and **homogeneous solutions**. Another method uses the *z*-transform (paralleling the Laplace transform). The state variable approach provides another *formulation* of the problem and solutions in the time as well as frequency domains.

Example 1.5.2 (Moving average filter) The three-term average $y(n) = \frac{x(n) + x(n-1) + x(n-2)}{3}$

as a lead-in to FIR and IIR, see below.

"FIR", "IIR", "Recursive" and "Nonrecursive" In the first example above the impulse response $h(n) = a^n$, $n \ge 0$ lasts for all positive time and is of infinite duration. In the second example (moving average) $h(n) = \{1/3, 1/3, 1/3\}$ which is of finite duration.

Definition If the unit sample response of a linear shift invariant system is of infinite duration, the system is said to be an **infinite impulse response (IIR)** system.

Definition If the unit sample response of a linear shift invariant system is of finite duration, the system is said to be a **finite impulse response (FIR)** system.

Theorem A causal linear shift invariant system characterized by

48 of 77

$$\sum_{k=0}^{N} a_k y(n-k) = \sum_{r=0}^{M} b_r x(n-r)$$

represents a finite impulse response (FIR) system if $a_0 \neq 0$, and $a_k = 0$ for k = 1, 2, ..., N. [Otherwise it could represent either an IIR or FIR system.] This is equivalent to saying that for an FIR system N = 0. For an FIR system then we have

$$a_0 y(n) = \sum_{\substack{m \ (b) \\ b}} b_r x(n-r), \text{ or}$$
$$y(n) = \sum_{r=0}^{M} \left| \frac{-r}{b} \right| x(n-r)$$

The above difference equation is identical to the convolution sum, and the (b_r/a_0) terms can be recognized as h(r), the value of the unit sample response at time r, i.e., we can set $(b_r/a_0) = h_r = h(r)$. So the impulse response, h(n), is given by

$$h(n) = (b_n/a_0), \ 0 \le n \le M$$

0, otherwise

which, obviously, is of finite duration.

Note: If the above difference equation were written so that $a_0 = 1$, we have $y(n) = \sum_{r=0}^{M} b_r x(n-r)$.

In this case the impulse response consists simply of the coefficients b_r^+ of the x(n-r) terms.

Iterative solution with initial conditions The LTI discrete-time system can be characterized by

$$\sum_{k=0}^{N} a_k y(n-k) = \sum_{r=0}^{M} b_r x(n-r), \qquad a_0 \neq 0, \text{ and } n \ge 0 \qquad \to (A)$$

Here N is the order of the difference equation. When written out in full the equation is

$$a_0 y(n) + a_1 y(n-1) + \dots + a_N y(n-N) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M), \qquad a_0 \neq 0, \text{ and } n \ge 0$$

The equation can be divided through by a_0 so that the coefficient of y(n) is 1 or, alternatively, we could impose the equivalent condition that $a_0 = 1$.

An alternative form of the above equation is sometimes given as

$$\sum_{k=0}^{N} a_k y(n+k) = \sum_{r=0}^{M} b_r x(n+r), \qquad n \ge 0 \quad \to (A'')$$

In this form, if the system is causal, we must have $M \leq N$.

The solution to either one of the above equations can be determined (by analogy with the differential equation) as the sum of two components: (1) the **homogeneous solution**, which depends on the **initial conditions** assumed to be known, and (2) the **particular solution**, which depends on the **input**.

Computing y(n) for successive values of *n* one after another is called an iterative solution. To obtain y(n) for $n \ge 0$ from Eq. (A) we rearrange it as

$$y(n) = -\sum_{k=1}^{\infty} \left(\frac{a_k}{a^0} \right) y(n-k) + \dots$$
 (, \to (B)

and evaluate y(n) for n = 1, 2, ... in an iterative manner. \rightarrow e need the initial conditions y(-1), y(-2), ..., y(-N). The initial conditions needed to solve for y(n) using Eq. (A^{**}) in a similar fashion are y(0), y(1), ..., y(N-1).

We may assume that the system described by Eq. (A) is in a condition of initial rest, that is, if x(n) = 0 for n < 0, then y(n) = 0 for n < 0 as well. With initial rest the system (A) is LTI *and* causal.

An equation of the form (A) or (B) is called a **recursive equation** since it specifies a recursive procedure to determine the output y(n) in terms of the input and *previous output values*.

In the special case where N = 0, Eq. (B) reduces to

Here y(n) is an explicit function of the present and previous values of the input only. Eq. (C) is called a **non-recursive equation**, since we do not recursively use previously computed values of the *output* in order to calculate the present value of the output.

Example 1.5.3 Find the solution to

$$y(n) - \frac{3}{2}y(n-1) + \frac{1}{2}y(n-2) = \frac{1}{2}$$
, $n \ge 0$

with the initial conditions y(-1) = 4 and y(-2) = 10.

Answer Note that the input is x(n) = (. This is the iterative solution in the time domain. We

can write $y(n) = \frac{3}{2}y(n-1) - \frac{1}{2}y(n-2) + \frac{1}{2}$, and solve for y(n) starting with n = 0:

y (\mathbf{n}) - ($\mathbf{y}(\mathbf{n}-1)$ + ($\mathbf{y}(\mathbf{n}-2)$ = (,	$n \ge 0$					
y(-1) = 4, and $y(-2) = 10$							
n	$y(n) = (3/2) y(n-1) - (1/2) y(n-2) + (1/4)^{n}$	y (n)					
0	$y(n) = (3/2) y(n-1) - (1/2) y(n-2) + (1/4)^{n}$	2					
	$y(0) = (3/2) y(0-1) - (1/2) y(0-2) + (1/4)^{0}$						
	$= (3/2) y(-1) - (1/2) y(-2) + (1/4)^{0}$						
	= (3/2) (4) - (1/2) (10) + (1)						
1	$y(n) = (3/2) y(n-1) - (1/2) y(n-2) + (1/4)^{n}$	5/4					
	$y(1) = (3/2) y(1-1) - (1/2) y(1-2) + (1/4)^{1}$						
	$= (3/2) y(0) - (1/2) y(-1) + (1/4)^{1}$						
	= (3/2) (2) - (1/2) (4) + (1/4)						
2	$y(n) = (3/2) y(n-1) - (1/2) y(n-2) + (1/4)^{n}$	15/16					
	$y(2) = (3/2) y(2-1) - (1/2) y(2-2) + (1/4)^2$						
	$= (3/2) y(1) - (1/2) y(0) + (1/4)^{2}$						
	= (3/2) (5/4) - (1/2) (2) + (1/16)						
3	$y(n) = (3/2) y(n-1) - (1/2) y(n-2) + (1/4)^{n}$	51/64					
	$y(3) = (3/2) y(3-1) - (1/2) y(3-2) + (1/4)^3$						
	$= (3/2) y(2) - (1/2) y(1) + (1/4)^{3}$						
	= (3/2) (15/16) - (1/2) (5/4) + (1/64)						
4	$y(n) = (3/2) y(n-1) - (1/2) y(n-2) + (1/4)^{n}$?					
	$y(4) = (3/2) y(4-1) - (1/2) y(4-2) + (1/4)^4$						
	$= (3/2) y(3) - (1/2) y(2) + (1/4)^4$						
	= (3/2) (51/64) - (1/2) (15/16) + (1/256))					
	Etc.						

50 of 77

The solution is

$$y(n) = \begin{cases} 1, \frac{5}{4}, \frac{15}{16}, \frac{51}{64}, \dots \end{cases}$$

This procedure does not, in general, yield an analytical expression for y(n), that is, a closed-form solution. But it is easily implemented on a digital computer.

Example 1.5.4 [MATLAB] [Zero initial conditions] In the context of MATLAB, we may use *filter*(*b*, *a*, *x*) to generate the sequence y(n). In MATLAB the coefficients of y(.) and x(.) are numbered slightly differently as below:

$$a_1 y(n) + a_2 y(n-1) + a_3 y(n-2) + \dots = b_1 x(n) + b_2 x(n-1) + b_3 x(n-2) + \dots$$

From the difference equation

$$y(n) - (3/2)y(n-1) + (1/2)y(n-2) = (1/4)^n, \qquad n \ge 0$$

we note that the input is $x(n) = (1/4)^n$ and the coefficients of y(.) and x(.) give us the *a* and *b* vectors, respectively: a = [1, -1.5, 0.5] and b = [1]. The following segment generates the output with initial conditions taken as zero.

%Zero initial conditions b = [1]; a = [1, -1.5, 0.5]; $n = 0.25; xn = (1/4) .^n; \%".^" means element-by-element exponentiation$ yn = filter(b, a, xn), %Or, yn = filter(1, a, xn)subplot(2, 1, 1), stem(n, xn);xlabel('n'), ylabel('x(n)'); title('Input Sequence');subplot(2, 1, 2), stem(n, yn);xlabel('n'), ylabel('y(n)'); title('Output Sequence');

The solution is: yn = $\begin{bmatrix} 1 & 1.75 & 2.1875 & 2.4219 & 2.5430 & 2.6045 & 2.6355 \dots 2.6667 & 2.6667 & 2.6667 \end{bmatrix}$



Example 1.5.5 [MATLAB] [Non-zero initial conditions] Find the solution to the difference equation

 $y(n) - (3/2)y(n-1) + (1/2)y(n-2) = (1/4)^n$, $n \ge 0$ with the initial conditions y(-1) = 4 and y(-2) = 10.

If the problem specifies non-zero initial conditions, they must first be converted to "equivalent initial conditions" for the *filter* function to work. In this problem y(-1) = 4 and y(-2) = 10. These are specified as a vector yic = [4, 10] and used to generate the equivalent initial conditions *eic* by the function *filtic(b, a, yic)*. The equivalent initial conditions are then used to generate the filter output through *filter(b, a, xn, eic)*. The MATLAB segment follows:

```
% Non-zero initial conditions

b = [1]; a = [1, -1.5, 0.5]; yic = [4, 10];

n = 0.25; xn = (1/4) .^n; \%".^" means element-by-element exponentiation

%

% Equivalent initial conditions

<math>eic = filtic(b, a, yic);

yn = filter(b, a, xn, eic)

subplot(2, 1, 1), stem(n, xn);

xlabel('n'), ylabel('x(n)'); title('Input Sequence');

subplot(2, 1, 2), stem(n, yn);

xlabel('n'), ylabel('y(n)'); title('Output Sequence');

The output is:
```

yn = [2 1.25 0.9375 0.7969 0.7305 0.6982 0.6824...0.6667 0.6667 0.6667]



(*Omit*) Recursive realization of FIR system – a simple example The moving average of the signal x(n) is given by

$$y(n) = \frac{x(n) + x(n-1) + x(n-2) + \dots + x(n-M)}{M+1}$$

Since past values of y(.) are not used in computing y(n) this is a nonrecursive implementation of the FIR filter. The impulse response is

$$n = \left\{\begin{array}{cccc} Q & & & M \\ h(n) = & & & M \\ \frac{1}{\sqrt{m+1}} & & \frac{1}{M+1} & & \frac{1}{M+1} & & \frac{1}{M+1} \\ & & & M+1 & & M+1 \end{array}\right\}$$
A total of M+1 terms

which consists of (M+1) samples or coefficients. We recognize, by replacing *n* by (n-1) in the above equation, that

$$y(n-1) = \frac{x(n-1) + x(n-2) + \dots + x(n-M) + x(n-M-1)}{M+1} \quad < (1)$$

By adding and subtracting $\frac{x(n-M-1)}{M+1}$ on the right hand side of the equation for y(n) we can arrive at a recursive implementation of the moving average:

$$y(n) = \frac{x(n) + x(n-1) + x(n-2) + \dots + x(n-M) + x(n-M-1) - x(n-M-1)}{M+1}$$

Using eq. (1) we can write y(n) as x(n) - x(n - M - 1)

which clearly is a recursive implementation since it involves y(n-1), a past value of y.

(End of Omit)

Further nomenclature If we take $a_0 = 1$ our standard difference equation becomes

$$y(n) + a_1 y(n-1) + \dots + a_N y(n-N) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M)$$

This is an N^{th} order difference equation.

A moving average (MA) filter is one with its output dependent on the present and previous inputs. The average here is a *weighted average*, the weights being the *b* coefficients. In terms of the above difference equation this corresponds to N = 0 and the difference equation becomes

 $y(n) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M)$

There is some variation on the form of this equation. As given here it has (M+1) coefficients and M delay elements. Sometimes it is more convenient to use this with M coefficients rather than (M+1). Either way its *order* is N = 0. Note that we have tied the *order* to the oldest term in y(.) and not to the oldest term in x(.). See *Aside* below.

An **autoregressive** (AR) filter is one with its output dependent on the present input (but not previous inputs) and previous outputs – this is a *purely recursive* system. In terms of the above difference equation this corresponds to M = 0 and is of the form

$$y(n) = b_0 x(n) - a_1 y(n-1) - \dots - a_N y(n-N)$$

More generally, an **autoregressive moving average** (**ARMA**) filter has its output depend on the present input, M previous inputs, and N previous outputs. In terms of the linear constant coefficient difference equation this has the form

$$y(n) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M) - a_1 y(n-1) - \dots - a_N y(n-N)$$

(*Aside*) With $a_0 = 1$ we have

$$H(z) = \frac{\sum_{i=1}^{N} b^{i} z^{-i}}{1 + \sum_{i=1}^{N} a_{i} z^{-i}}$$

This represents an IIR filter if at least one of a_1 through a_N is nonzero, and all the roots of the denominator are not canceled exactly by the roots of the numerator. For example, in the system, $H(z) = (1 - z^{-8}) (1 - z^{-1})$, the single pole at z = 1 is canceled exactly by the zero at z = 1, making H(z) a finite polynomial in z^{-1} , that is, an FIR filter.

In general, there are M finite zeros and N finite poles. There is no restriction that M should be less than or greater than or equal to N. In most cases, especially digital filters derived from analog designs, M will be less than or equal to N. Systems of this type are called N^{th} order systems.

When M > N, the order of the system is no longer unambiguous. In this case, H(z) may be taken to be an N^{th} order system in cascade with an FIR filter of order (M - N).

When N = 0, as in the case of an FIR filter, according to our convention the *order* is 0; it is more useful in this case to focus on M and call it an FIR filter of M stages or (M+1) coefficients.

(End of Aside)

Fourier analysis of discrete-time signals and systems

Note For the discrete-time Fourier transform some authors (Oppenheim & Schafer, for instance) use the symbol $X(e^{j\omega})$ while others (Proakis, for instance) use the symbol $X(\omega)$. The symbol ω is used for digital frequency (radians per sample or just radians) and the symbol Ω for the analog frequency (radians/sec). Some authors, on the other hand, use just the opposite of our convention, that is, ω for the analog frequency (radians/sec) and Ω for the digital frequency (radians).

Discrete-time Fourier transform (DTFT) For the continuous-time signal x(t), the Fourier transform is

$$F\{x(t)\} = X(\Omega) = \int_{-\infty}^{\infty} x(t) e^{-j\Omega t} dt$$

The *impulse-train* sampled version, $x_s(t)$, is given by

$$x_s(t) = x(t) \sum_{n = -\infty}^{\infty} \delta(t - nT)$$

So the Fourier transform of $\bar{x}_s(t)$ is given by

$$X_{s}(\Omega) = \int_{-\infty}^{\infty} x(t)e^{-j\Omega t} dt = \int_{-\infty}^{0} \left(\frac{\sum_{n=-\infty}^{\infty} \delta(t-nT)}{x(t)} \right) e^{-j\Omega t} dt$$
$$= \sum_{n=-\infty}^{\infty} x(nT) e^{-j\Omega nT}$$

where the last step follows from the sifting property of the δ function. Replace ΩT by ω the discrete-time frequency variable, that is, the **digital frequency**. Note that Ω has units of radians/second, and ω has units of radians (/sample). This change of notation gives the **discrete-time Fourier transform**, $X(\omega)$, of the discrete-time signal x(n), obtained by sampling x(t), as

$$X(\omega) = F\{x(n)\} = \sum_{n = -\infty} x(n) e^{-j\omega n}$$

Note that this defines the discrete-time Fourier transforms of *any* discrete-time signal x(n). The transform exists if x(n) satisfies a relation of the type

$$\sum_{n=-\infty}^{\infty} |x(n)| < \infty \qquad \text{or} \qquad \sum_{n=-\infty}^{\infty} |x(n)|^2 < \infty$$

These conditions are sufficient to guarantee that the sequence has a discrete-time Fourier transform. As in the case of continuous-time signals there are signals that neither are absolutely summable nor have finite energy, but still have a discrete-time Fourier transform. (See also p. 22, O & S.)

Discrete-time Fourier transform of (non-periodic) sequences The Fourier transform of a general discrete-time sequence tells us what the frequency content of that signal is.

Definition The Fourier transform $X(e^{j\omega})$ of the sequence x(n) is given by

$$F\{x(n)\} = X(\omega) = \sum_{n = -\infty} x(n) e^{-j\omega n} \qquad < (A)$$

The inverse Fourier transform is given by

$$F^{-1}\{X(\omega)\} = x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega}) e^{j\omega n} d\omega \qquad < (B)$$

Equations (A) and (B) are called the Fourier transform pair for a sequence x(n) with $X(\omega)$ thought of as the *frequency content* of the sequence x(n). Equation (A) is the analysis equation and equation (B) is the synthesis equation. Since $X(\omega)$ is a periodic function of ω , we can think of x(n) as the Fourier coefficients in the Fourier series representation of $X(\omega)$. That is, equation (A), in fact, expresses $X(\omega)$ in the form of a Fourier series.

The sketch below sums up the relationship between the time and frequency domains. The periodicity is 2π . From the relation $\omega = \Omega T$ we can deduce that at the point $\omega = 2\pi$ on the horizontal axis $\Omega = \omega/T = 2\pi/T = 2\pi F_s = \Omega_s$. In other words, in terms of the analog frequency variable the point $\omega = 2\pi$ corresponds to $\Omega = \Omega_s$ or $F = F_s$.



(*Omit*) *Relationship to the z-transform* The *z*-transform X(z) and the Fourier transform $X(\omega)$ are given by

$$X(z) = \sum_{n = -\infty}^{\infty} x(n) \ z^{-n} \quad \text{and} \quad X(\omega) = \sum_{n = -\infty}^{\infty} x(n) e^{-j\omega n}$$

Comparing the two we deduce the relationship as

$$X(z)\Big|_{z=e^{j\varpi}} \equiv X(\omega) = \sum_{n=-\infty} x(n) e^{-j\omega n}$$

The *z*-transform evaluation on the unit circle gives the Fourier transform of the sequence x(n). The *z*-transform of x(n) can be viewed as the Fourier transform of the sequence $\{x(n) r^{-n}\}$, that is, x(n) multiplied by an exponential sequence r^{-n} . This can be seen by setting $z = r e^{j\omega}$ in the defining equation of X(z):

$$X(z) = \sum_{n = -\infty}^{\infty} x(n) \left(r e^{j\omega} \right)^{-n} = \sum_{n = -\infty}^{\infty} x(n) r^{-n} e^{-j\omega n}$$

$$z \qquad \qquad x(n) \text{ multiplied by the exponential sequence } r$$

Z-transform of a periodic sequence Consider a sequence x(n) that is periodic with period *N* so that x(n) = x(n+kN) for any integer value of *k*. Such a sequence cannot be represented by its z-transform, since there is no value of *z* for which the z-transform will converge. (*End of Omit*)

Example 1.6.1 [Cf. Example 4.2.3, Proakis, 4th Ed.] For the exponential sequence $x(n) = a^n u(n)$, |a| < 1, the DTFT is

$$X(e^{j\omega}) = \sum_{n=0}^{\infty} a^n e^{-j\omega n} = \sum_{n=0}^{\infty} (ae^{-j\omega})^n = \frac{1}{1 - ae^{-j\omega}} = \frac{1}{1 - a(\cos \omega - j\sin \omega)}$$

We shall put this in the form $X(\omega) = Magnitude \{X\} e^{jPhase\{X\}} = |X(\omega)| e^{j \angle X(\omega)}$ from which the magnitude and phase will be extracted. The denominator (Dr.) is

Thus

Dr. =
$$1 - a\cos\omega + ja\sin\omega = \sqrt{(1 - a\cos\omega)^2 + a^2\sin^2\omega} e^{-\frac{1}{(1 - a\cos\omega)^2}}$$

 $X(\omega) = |X(\omega)| e^{j\angle X(\omega)} = 1$

$$\frac{1}{\sqrt{(1 + a^2 - 2a\cos\omega)}} e^{-\frac{1}{(1 - a\cos\omega)^2}}$$
de and phase are: 1
$$\frac{1}{|X(\omega)|} = 1$$
and $\angle X(\omega) = -\tan^{-1}(-a\sin\omega)$

The magnitude and phase ar $|X(\omega)| =$

$$\sqrt{(1+a^2-2a\cos\omega)}$$

 $1 - a \cos \omega$

= 0.

Plots of |X| and $\angle X$ are shown. Note that $X(\omega)$ is periodic and that the magnitude is an even function of ω and the phase is an odd function. (See below on the notation |X| and $\angle X$).

The value of $X(e^{j\omega})$ at $\omega = 0$ is

$$|X(\omega)|_{\omega=0} = \frac{1}{\sqrt{(1+a_{\ell}^2 - 2a\cos 0)}} = \frac{1}{1-a}$$

$$\angle X(\omega)|_{\omega=0} = -\tan^{-1}\left|\frac{a\sin 0}{1-a\cos 0}\right| = 0$$

$$(1-a\cos 0)$$

Similarly, at $\omega = \pi$ we have $|X(\omega)|_{\omega=\pi} = 1/(1+a)$ and $\angle X(\omega)|_{\omega=\pi}$

Phase angle of a complex number In calculating the phase angle of a complex number, z = Re(z) + j Im(z), a hand held calculator, typically uses the formula $\tan^{-1} \{ Im(z)/Re(z) \}$, and returns an answer in the range $-\pi/2 < \text{angle} \le \pi/2$. Thus for both (1-j) and (-1+j) the phase angle so calculated is $\pi/4$. However, (1-j) is in the 4th quadrant with $\angle (1-j) = -\pi/4$ whereas (-1+j) is in the 2nd quadrant with $\angle (-1+j) = 3\pi/4$. MATLAB has a function *angle* which takes into account the real and imaginary parts separately (instead of their ratio) and calculates the "four-quadrant inverse tangent".

Example 1.6.2 [MATLAB *fplot*] To illustrate the magnitude and phase plots of the DTFT we take a = 0.8 in the exponential sequence $x(n) = a^n u(n)$, |a| < 1, treated above. Thus $x(n) = (0.8)^n u(n)$. We need only plot over the interval $-\pi \le \omega \le \pi$ or $0 \le \omega \le 2\pi$. To show, visually, the periodicity we have plotted over $-3\pi \le \omega \le 3\pi$. We have

$$X(e^{j\omega}) = \sum_{n=0}^{\infty} a^n e^{-j\omega n} = \sum_{n=0}^{\infty} (ae^{-j\omega})^n = \frac{1}{1 - ae} = \frac{1}{1 - 0.8 e}$$

In the MATLAB program segment that follows we write an algebraic expression for $X(e^{j\omega}) = \frac{1}{-0.8 e^{-j\omega}}$ as (1) / (1-0.8*exp(-j*w)). Note that we have used ,,w" for ω and the pbt1

ranges from -2π to 2π . Both ω and the phase, $\angle X(\omega)$, are in radians. The parameter 'k' means that the plot/display is in black "color".

subplot(2,1,1);fplot('abs((1)/(1-0.8*exp(-j*w)))', [-3*pi,3*pi], 'k'); xlabel('\omega');ylabel('Magnitude'); subplot(2,1,2);fplot('angle((1)/(1-0.8*exp(-j*w)))', [-3*pi,3*pi], 'k'); xlabel('\omega');ylabel('Phase');

$$|X(\omega)|_{\omega=0} = \frac{1}{1-a} = \frac{1}{1-0.8} = 5$$
$$|X(\omega)|_{\omega=\pi} = 1/(1+a) = 1/1.8 = 0.555$$



Example 1.6.3 [MATLAB *freqz*] We repeat the frequency response plots using the *freqz* function. Taking a = 0.8 as before, we have $x(n) = (0.8)^n u(n)$ and

$$X(e^{j\omega}) = \sum_{n=0}^{\infty} a^n e^{-j\omega n} = \sum_{n=0}^{\infty} \left(ae^{-j\omega}\right)^n = \frac{1}{1 - ae} = \frac{1}{1 - 0.8 e}$$

Instead of writing an algebraic expression for $X(e^{j\omega}) = \frac{1}{1 - 0.8 e^{-j\omega}}$ as in the previous example, we

shall now specify the numerator and denominator in terms of their coefficients. We shall use the following convention to specify the parameters of the function

$$\frac{X(e^{j\omega}) = b(1) + b(2)e^{-j\omega} + b(3)e^{-j2\omega} + \dots}{a(1) + a(2)e^{-j\omega} + a(3)e^{-j2\omega} + \dots} = \frac{1}{1 - 0.8e^{-j\omega}}$$

Here the vectors *b* and *a* specify, respectively, the numerator and denominator coefficients. In our example b(1) = 1, a(1) = 1, and a(2) = -0.8. The MATLAB segment and the corresponding plots follow. Note that the plot goes from $-\pi$ to π , not -3π to 3π .

b = [1]; %Numerator coefficient a = [1, -0.8]; %Denominator coefficients w = -pi: pi/256: pi; %A total of 512 points [Xw] = freqz(b, a, w); subplot(2, 1, 1), plot(w, abs(Xw)); xlabel('Frequency \omega'), ylabel('Magnitude of X(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Xw)); xlabel('Frequency \omega'), ylabel('Phase of X(\omega)'); grid



Example 1.6.4 Find the DTFT, $X(\omega)$, for $x(n) = \{1, 2, 3, 4, 3, 2, 1\}$. **Solution** The DTFT is

$$X(e^{j\omega}) = \sum_{n = \infty} x(n) e^{-j\omega n} = 1 + 2e^{-j\omega 1} + 3e^{-j\omega 2} + 4e^{-j\omega 3} + 3e^{-j\omega 4} + 2e^{-j\omega 5} + 1e^{-j\omega 6}$$
$$= (2\cos 3\omega + 4\cos 2\omega + 6\cos \omega + 4) e^{-j3\omega}$$

The numerator vector is b = [1, 2, 3, 4, 3, 2, 1] and the denominator vector is $a = \{1\}$. The MATLAB segment plots the sequence and the frequency response.

%Sketch of sequence n = 0:1:6; xn = [1, 2, 3, 4, 3, 2, 1]; subplot (3, 1, 1), stem(n, xn) xlabel('n'), ylabel('x(n)'); grid %Frequency response b = [1, 2, 3, 4, 3, 2, 1]; %Numerator coefficients a = [1]; %Denominator coefficients w = -pi: pi/256: pi; %A total of 512 points [Xw] = freqz(b, a, w); subplot(3, 1, 2), plot(w, abs(Xw)); xlabel('Frequency \omega'), ylabel('Magnitude of X(\omega)'); grid subplot(3, 1, 3), plot(w, angle(Xw)); xlabel('Frequency \omega'), ylabel('Phase of X(\omega)'); grid



Example 1.6.5 Find the DTFT, $X(\omega)$, for (a) $x(n) = \{1, 2, 3, 4, 3, 2, 1, 0\}$ (b) $x(n) = \{1, 2, 3, 4, 3, 2, 1, 0, ..., 0\}$ Solution !?!

The following examples are repeated in Unit II – DFS & DFT.

Example 1.6.6 Obtain the 7-point DFT of the sequence $x(n) = \{1, 2, 3, 4, 3, 2, 1\}$ by taking 7 samples of its DTFT uniformly spaced over the interval $0 \le \omega \le 2\pi$.

Solution The sampling interval in the frequency domain is $2\pi/7$. From Example 4 we have $X(e^{j\omega})$ or $X(\omega) = 1 + 2e^{-j\omega 1} + 3e^{-j\omega 2} + 4e^{-j\omega 3} + 3e^{-j\omega 4} + 2e^{-j\omega 5} + 1e^{-j\omega 6}$

$$= (2\cos 3\omega + 4\cos 2\omega + 6\cos \omega + 4) e^{-j3\omega}$$

The DFT, X(k), is given by replacing ω with $k(2\pi/7)$ where k is an index ranging from 0 to 6: DFT = $X(\omega)$ = $X(2\pi k/7)$, k = 0 to 6

This is denoted XkfromDTFT in the MATLAB segment below.

MATLAB:

w = 0: 2*pi/7: 2*pi-0.001XkfromDTFT = (4+6*cos(w)+4*cos(2*w)+2*cos(3*w)) .* exp(-j*3*w)

MATLAB solution:

XkfromDTFT = [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)]

This is the 7-point DFT obtained by sampling the DTFT at 7 points uniformly spaced in $(0, 2\pi)$. It should be the same as the DFT directly obtained, for instance, by using the *fft* function in MATLAB:

MATLAB:

xn = [1 2 3 4 3 2 1]Xkusingfft = fft(xn)

MATLAB solution:

Xkusingfft = [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)]

It can be seen that " \downarrow kfromDTFT" = " \downarrow kusingfft".

Example 1.6.7 Obtain the 7-point inverse DTFT x(n) by finding the 7-point inverse DFT of X(k):

 $X (2\pi k / 7) = [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)]$

MATLAB:

 $\begin{aligned} Xk &= [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)] \\ xn &= ifft(Xk) \end{aligned}$

MATLAB solution:

xn = [1.0000 2.0000 3.0000 4.0000 3.0000 2.0000 1.0000]

This is the original sequence we started with in Example 4.

Example 1.6.8 What will be the resulting time sequence if the DTFT of the 7-*point* sequence is sampled at 6 (or fewer) uniformly spaced points in $(0, 2\pi)$ and its inverse DFT is obtained? **Solution** The sampling interval in the frequency domain now is $2\pi/6$. From Example 4 we have

$$X(e^{j\omega}) \text{ or } X(\omega) = 1 + 2e^{-j\omega 1} + 3e^{-j\omega 2} + 4e^{-j\omega 3} + 3e^{-j\omega 4} + 2e^{-j\omega 5} + 1e^{-j\omega 6}$$
$$= (2\cos 3\omega + 4\cos 2\omega + 6\cos \omega + 4)e^{-j3\omega}$$

The DFT then is given by

DFT = $X(\omega)|_{\omega = 2\pi k/6} = X(2\pi k/6), \quad k = 0 \text{ to } 5$

This is denoted Xk6point in the MATLAB segment below.

MATLAB:

w = 0: 2*pi/6: 2*pi-0.001 Xk6point = (4+6*cos(w)+4*cos(2*w)+2*cos(3*w)) .* exp(-j*3*w)

MATLAB solution:

 $\begin{aligned} Xk6point = [16, (-3.0000 - 0.0000i), (1.0000 + 0.0000i), 0, (1.0000 + 0.0000i), \\ (-3.0000 - 0.0000i)] \end{aligned}$

This is the 6-point DFT obtained by sampling the DTFT at 6 points uniformly spaced in $(0, 2\pi)$.

Example 1.6.9 Obtain the 6-point inverse DTFT x(n) by finding the 6-point inverse DFT of Xk6point:

 $X (2\pi k / 6) = [16, (-3.0000 - 0.0000i), (1.0000 + 0.0000i), 0, (1.0000 + 0.0000i), (-3.0000 - 0.0000i)]$

MATLAB:

Xk6point = [16, (-3.0000 - 0.0000i), (1.0000 + 0.0000i), 0, (1.0000 + 0.0000i), (-3.0000 - 0.0000i)] xn = ifft(Xk6point)

MATLAB solution:

 $\mathbf{xn} = \begin{bmatrix} 2 & 2 & 3 & 4 & 3 & 2 \end{bmatrix}$

Comparing with the original 7-point sequence, $xn = \begin{bmatrix} 1 & 2 & 3 & 4 & 3 & 2 & 1 \end{bmatrix}$, we see the consequence of *under-sampling* the continuous- ω function $X(\omega)$: the corresponding time domain sequence x(n) is said to suffer *time-domain aliasing*. This is similar to the situation that occurs when a continuous-time function x(t) is under-sampled: the corresponding frequency domain function $X_s(\omega)$ contains *frequency-domain aliasing*.

Example 1.6.10 What will be the resulting time sequence if the DTFT of the 7-*point* sequence is sampled at 8 (or more) uniformly spaced points in $(0, 2\pi)$ and its inverse DFT is obtained? **Solution** The sampling interval in the frequency domain now is $2\pi/8$. From Example 4 we have

$$X(e^{j\omega}) \text{ or } X(\omega) = 1 + 2e^{-j\omega 1} + 3e^{-j\omega 2} + 4e^{-j\omega 3} + 3e^{-j\omega 4} + 2e^{-j\omega 5} + 1e^{-j\omega 6}$$
$$= (2\cos 3\omega + 4\cos 2\omega + 6\cos \omega + 4)e^{-j3\omega}$$

The DFT then is given by

DFT = $X(\omega)|_{\omega = 2\pi k/8} = X(2\pi k/8), \qquad k = 0 \text{ to } 7$

This is denoted Xk8point in the MATLAB segment below.

MATLAB:

w = 0: 2*pi/8: 2*pi-0.001 Xk8point = (4+6*cos(w)+4*cos(2*w)+2*cos(3*w)) .* exp(-j*3*w)

MATLAB solution:

Xk8point = [16, (-4.8284 - 4.8284i), (0.0000 - 0.0000i), (0.8284 - 0.8284i), 0, (0.8284 + 0.8284i), (0.0000 - 0.0000i), (-4.8284 + 4.8284i)]

This is the 8-point DFT obtained by sampling the DTFT at 8 points uniformly spaced in $(0, 2\pi)$.

Example 1.6.11 Obtain the 8-point inverse DTFT x(n) by finding the 8-point inverse DFT of Xk8point:

 $X (2\pi k / 8) = [16, (-4.8284 - 4.8284i), (0.0000 - 0.0000i), (0.8284 - 0.8284i), 0, (0.8284 + 0.8284i), (0.0000 - 0.0000i), (-4.8284 + 4.8284i)]$

MATLAB:

Xk8point = [16, (-4.8284 - 4.8284i), (0.0000 - 0.0000i), (0.8284 - 0.8284i), 0, (0.8284 + 0.8284i), (0.0000 - 0.0000i), (-4.8284 + 4.8284i)] xn = ifft(Xk8point)

MATLAB solution:

 $xn = [1 \ 2 \ 3 \ 4 \ 3 \ 2 \ 1 \ 0]$

We see that the original 7-point sequence has been preserved with an *appended zero*. The original sequence and the *zero-padded sequence* (with any number of zeros) have the same DTFT. This is a case of *over-sampling* the continuous- ω function $X(\omega)$: there is no *time-domain aliasing*. This is similar to the situation that occurs when a continuous-time function x(t) is oversampled: the corresponding frequency domain function $X_s(\omega)$ is free from *frequency-domain aliasing*.

%Sketch of sequences n = 0:1:6; xn = [1, 2, 3, 4, 3, 2, 1];



Frequency response of discrete-time system

For a linear shift-invariant system with impulse response h(n), the Fourier transform $H(\omega)$ gives the *frequency response*. Consider the input sequence $x(n) = e^{j\omega n}$ for $-\infty < n < \infty$, i.e., a complex exponential of radian frequency ω and magnitude 1, applied to a linear shift-invariant system whose unit sample response is h(n). Using convolution we obtain the output y(n) as

$$y(n) = h(n) * x(n) = \sum_{k=-\infty}^{\infty} h(k) x(n-k) = \sum_{k=-\infty}^{\infty} h(k) e^{j\omega(n-k)} = e^{j\omega n} \sum_{k=-\infty}^{\infty} h(k) e^{-j\omega k}$$
$$H(\omega)$$
$$= H(\omega) e^{j\omega n}$$

Thus we see that $H(\omega)$ describes the change in complex amplitude of a complex exponential as a function of frequency. The quantity $H(\omega)$ is called the **frequency response** of the system. In

general, $H(\omega)$ is complex valued and may be expressed either in the Cartesian form or the polar from as

$$H(\omega) = H_R(\omega) + j H_I(\omega)$$
 or $H(\omega) = \hat{H}(\omega) e^{j \angle H(\omega)}$

where H_R and H_I are the real part and imaginary part respectively. $\hat{H}(\omega)$ is loosely called the magnitude and $\angle H(\omega)$ is loosely called the phase. Strictly speaking, $\hat{H}(\omega)$ is called the **zerophase frequency response**; note that $\hat{H}(\omega)$ is *real valued* but may be positive or negative. We may use the symbol $|H(\omega)|$ for the magnitude which is strictly non-negative. If $\hat{H}(\omega)$ is positive then

Magnitude =
$$|H(\omega)| = \hat{H}(\omega)$$
 & Phase = $\angle H(\omega)$

If $\hat{H}(\omega)$ is negative then

Magnitude = $|H(\omega)| = |\hat{H}(\omega)| = -\hat{H}(\omega)$ & Phase = $\angle H(\omega) \pm \pi$

We shall often loosely use the symbol $|H(\omega)|$ to refer to $\hat{H}(\omega)$ as well with the understanding that when the latter is negative we shall take its absolute value (the magnitude) and accordingly adjust $\angle H(\omega)$ by $\pm \pi$.

Example 1.7.1 [Moving average filter] The impulse response of the LTI system

3

$$y(n) = \frac{x(n) + x(n-1) + x(n-2)}{3}$$

is

h(n) = 1/3, n = 0, 1, 20. otherwise

$$h(n) = (1/3)[u(n) - u(n-3)]$$

The frequency response is obtained below.

$$H(\omega) = \sum_{k=0}^{\infty} h(k) e^{-j\omega k} = \sum_{k=0}^{-j\omega} (1/3) e^{-j\omega k} = \frac{1}{3} \left(e^{-j\omega 0} + e^{-j\omega 1} + e^{-j\omega 2} \right)$$

$$= \frac{e^{-j\omega}}{3} \left(e^{-j\omega 0} + 1 + e^{-j\omega} \right) = \frac{1}{3} \left| 1 + \frac{1}{2} \right|_{2} = \frac{1}{3} \left| 1 + \frac{1}{2} \right|_{2} = \frac{1}{3} \left| 1 + \frac{1}{3} \left| 1 + \frac{1}{3} \right|_{2} = \frac{1}{3} \left| 1 + \frac{1}{3} \left| 1 + \frac{1}{3} \right|_{2} = \frac{1}{3} \left| 1 + \frac{1}{3} \left| 1 + \frac{1}{3} \right|_{2} = \frac{1}{3} \left| 1 + \frac{1}{3} \left| 1 + \frac{1}{3} \right|_{2} = \frac{1}{3} \left| 1 + \frac{1}$$

which is already in the polar form $H(\omega) = \pm H(\omega) e^{j \angle H(\omega)}$, so that

 $|H(\omega)| = (1+2\cos\omega)/3$ and $\angle H(\omega) = -\omega$

The zero crossings of the magnitude plot occur where $H(\omega) \models (1 + 2\cos\omega) / 3 = 0$, or ω

 $=\cos^{-1}(-1/2) = 2\pi/3 = 120^{\circ}$. A frequency of $\omega = 2\pi/3$ rad./sample (f = 1/3 cycle/sample) is totally stopped (filtered out) by the filter. The corresponding digital signal is $x_5(n) = \cos \frac{1}{2}$ $2\pi(1/3)n$. The underlying continuous-time signal, $x_5(t)$, depends on the sampling frequency. If, for example, the sampling frequency is 16Hz, then $x_5(t) = \cos 2\pi (16/3)t$, and a frequency of 16/3 Hz will be totally filtered out. If the sampling frequency is 150Hz, then $x_5(t) = \cos 2\pi (150/3)t$, and a frequency of 50 Hz will be eliminated.

In calibrating the horizontal axis in terms of the cyclic frequency, F, we use the relation $\omega = \Omega T = 2\pi F T = 2\pi F/F_s$ from which the point $\omega = 2\pi$ corresponds to $F = F_s$.



(Aside) The system $y(n) = \frac{x(n)}{3} + \frac{x(n-1)}{3} + \frac{x(n-2)}{3}$ is a crude low pass filter, but the

attenuation does not increase monotonically with frequency. In fact, the highest possible frequency, $F_s/2$ Hz, (or π rad/sample) is not well attenuated at all. The following is a slight variation of the three-term moving average:

$$y(n) = \frac{x(n)}{4} + \frac{x(n-1)}{2} + \frac{x(n-2)}{4}$$

Its magnitude response is a "raised cosine" (with no zero crossing, monotonically decreasing but wider than the 3-term).

(End of Aside)

Example 1.7.2 [MATLAB *fplot*] [Moving average filter]

$$H(\omega) = \frac{1}{3} \Big(1 + e^{-j\omega \, 1} + e^{-j\omega \, 2} \Big).$$

The program follows. Note that in MATLAB we use ",w" for ω and the plot ranges over $(-2\pi, 2\pi)$. Both ω and the phase, $\angle X(\omega)$, are in radians.

 $subplot(2,1,1); fplot('abs((1/3)*(1+exp(-j*w)+exp(-j*2*w)))', [-2*pi,2*pi], 'k'); xlabel('\omega'); ylabel('Magnitude'); subplot(2,1,2); fplot('angle((1/3)*(1+exp(-j*w)+exp(-j*2*w)))', [-2*pi,2*pi], 'k'); xlabel('\omega'); ylabel('Phase');$



[Homework] Plot the output signal $y(n) = \frac{x(n) + x(n-1) + x(n-2)}{3}$ for several values of $n \ge 0$ where the input is $x(n) = \cos 2\pi(1/3)n$. Take x(-1) = x(-2) = 0.

Example 1.7.3 [MATLAB *filter*, *freqz*] [Moving average filter] Consider a signal $x(t) = \cos(2\pi 50t)$ sampled at 150 Hz. The corresponding $x(n) = \cos(2\pi n/3)$ is the input to the moving average filter. The following MATLAB segment plots x(n) vs. n and y(n) vs. n.

n=0:1:50; x=cos(2*pi*n/3); b=[1/3, 1/3, 1/3]; a=[1]; %Filter coefficients y=filter(b, a, x); subplot(3, 1, 1), stem(n, x, 'ko');title('Input x(n)=Cosine(2\pin/3)');





Example 1.7.4 [MATLAB *filter*] [Moving average filter] As an extension of above example, consider the signal consisting of a 2 Hz desirable component plus a noise component of 50 Hz (with a smaller amplitude), sampled at 150 Hz.

 $x(t) = 5 \cos (2\pi 2t) + 2 \cos (2\pi 50t)$ $x(n) = 5 \cos (2\pi n/75) + 2 \cos (2\pi n/3)$

In the following MATLAB segment we show both sequences, x(n) and y(n), on the same (multi)plot so that they have the same scale. The smoothing action of the filter is easily discernible in the multiplot.

 $\begin{array}{l} n = 0: \ 1: \ 50; \\ x = 5^* \cos(2^* p i^* n / 75) + 2^* \cos(2^* p i^* n / 3); \\ b = [1/3, \ 1/3, \ 1/3]; \ a = [1]; \ \% \ Filter \ coefficients \\ y = filter(b, a, x); \\ plot(n, x, \ b^*', n, y, \ 'ko'); \\ legend ('Input \ x(n)', \ 'Output \ y(n)'); \\ title ('Moving \ average \ filter'); \\ xlabel ('n'), \ ylabel('x(n), \ y(n)'); \end{array}$



The same plot is shown below over a longer period.



[Homework] Examples 5.1.1, 5.1.2, 5.1.3 and 5.1.4, Proakis, 4th Ed.

Sketches of ideal digital low pass and high pass filters Ref. p. 23, O&S for LP filter.

Example 1.7.5 [2003] [O & S, p. 20, and L.C. Ludeman, p. 51] A linear time-invariant system has unit sample response h(n) = u(n) - u(n-N). Find the amplitude and phase spectra.

Note that this would be an *N*-term moving average filter if h(n) = (1/N)[u(n) - u(n-N)].



$$h(n) = 1, \quad 0 \le n \le N-1$$

$$0, \quad \text{elsewhere}$$

$$H(\omega) = \sum_{k = -\infty}^{\infty} h(k) e^{-j\omega k} = \sum_{k=0}^{N-1} 1 e^{-j\omega k} = \frac{1 - e^{-j\omega N}}{1 - e^{-j\omega}}$$

$$= \frac{e^{-j\omega N/2} (e^{j\omega N/2} - e^{-j\omega N/2})}{e^{-j\omega/2} (e^{j\omega/2} - e^{-j\omega/2})} = \frac{\sin(\omega N/2)}{\sin(\omega/2)} e^{-j\omega(N-1)/2}$$

$$|H(\omega)| = \frac{\sin(\omega N/2)}{\sin(\omega/2)} \quad \text{and} \quad \angle H(e^{j\omega}) = -\omega(N-1)/2$$

See the plots in Oppenheim and Schafer.

Example 1.7.6 [MATLAB *fplot*] **[5-term moving average filter]** The program for the case where *N* = 5 follows.

$$H(\omega) = \left(1 + e^{-j\omega 1} + e^{-j\omega 2} + e^{-j\omega 3} + e^{-j\omega 4}\right) = \frac{1 - e^{-j\omega 5}}{2}$$

Note that in MATLAB we use ",w" for ω and the plot ranges from -2π to 2π . Both ω and the phase, $\angle X(\omega)$, are in radians.

 $subplot(2,1,1); fplot('abs((1-exp(-j*w*5))/(1-exp(-j*w)))', [-2*pi,2*pi], 'k'); xlabel('\omega (Rad)'); ylabel('Magnitude'); subplot(2,1,2); fplot('angle((1-exp(-j*w*5))/(1-exp(-j*w)))', [-2*pi,2*pi], 'k'); xlabel('\omega (Rad)'); ylabel('Phase (Rad)');$



Properties of the discrete-time Fourier transform (DTFT)

For the DTFT Oppenheim & Schafer use the symbol $X(e^{j\omega})$ while Proakis uses $X(\omega)$.

(1) **Periodicity** $X(\omega)$ is periodic with period 2π , that is, $X(\omega+2\pi) = X(\omega)$ for all ω . Since $e^{j\omega n}$ is periodic in ω with period 2π , it follows that $X(\omega)$ is also periodic with the same period. Replacing ω with $(\omega+2\pi)$ gives

$$X(\omega+2\pi) = \sum_{n=-\infty}^{\infty} x(n) \ e^{-j(\omega+2\pi)n} = \sum_{n=-\infty}^{\infty} x(n) \ e^{-j\omega n} e^{-j2\pi n} = \sum_{n=-\infty}^{\infty} x(n) \ e^{-j\omega n} 1$$
$$= X(\omega) \text{ for all } \omega$$

As a result, while in the continuous time case Ω ranges from $-\infty$ to ∞ , in the discrete-time case we need only consider values of ω over the range 0 to 2π (or, $-\pi$ to π , or any 2π -long interval).

(2) Linearity The discrete Fourier-transform is a linear operation. If $F\{x_1(n)\} = X_1(\omega)$ and $F\{x_2(n)\} = X_2(\omega)$, then $F\{a_1 x_1(n) + a_2 x_2(n)\} = a_1 X_1(\omega) + a_2 X_2(\omega)$ for any constants a_1 and a_2 .

(3) Time shifting Time shift results in phase shift. If $F\{x(n)\} = X(\omega)$, then $F\{x(n-k)\} = e^{-j\omega k}X(\omega)$. **Proof** We have

$$F\{x(n-k)\} = \sum_{n=-\infty}^{\infty} x(n-k) e^{-j\omega n}$$

On the right hand side set n-k = m, so that n = m+k and the limits $n = -\infty$ to ∞ change to $m = -\infty$ to $+\infty$. Then
$$F\{x(n-k)\} = \sum_{m=-\infty}^{\infty} x(m) \ e^{-j\omega(m+k)} = e^{-j\omega k} \sum_{m=-\infty}^{\infty} x(m) \ e^{-j\omega m} = e^{-j\omega k} X(\omega) \quad \text{QED}$$

(4) **Frequency shifting** Multiplication in the time domain by a complex exponential results in frequency shifting. Given $F\{x(n)\} = X(\omega)$, then $F\{e^{j\omega_0 n}x(n)\} = X(\omega-\omega_0)$. **Proof** We have

$$F\left\{e^{j\omega_0 n}x(n)\right\} = \sum_{n=-\infty}^{\infty} e^{j\omega_0 n}x(n)e^{-j\omega n} = \sum_{n=-\infty}^{\infty} x(n)e^{-j(\omega-\omega_0) n} = X(\omega-\omega_0) \quad \text{QED}$$

Alternatively, using the synthesis equation,

$$\mathcal{F}^{-1}\{X(\omega-\omega_0)\} = \frac{1}{2} \int_{2\pi} X(\omega-\omega) e^{j\omega n} d\omega$$

Set $\omega - \omega_0 = \lambda$ so that $\omega = \lambda + \omega_0$ and the limits $\omega = 0$ to 2π change to $\lambda = -\omega_0$ to $(-\omega_0 + 2\pi)$, which amounts to any interval of length 2π . Also $d\omega = d\lambda$. Then

$$\mathcal{F}^{-1}\{X(\omega-\omega_0)\} = \underbrace{\frac{1}{2}}_{2\pi} \int_{2\pi} X(\lambda) e^{j(\lambda+\omega_0)n} d\lambda$$
$$= e^{j\omega_0 n} \frac{1}{2\pi} \int_{2\pi} X(\lambda) e^{j\lambda n} d\lambda = e^{j\omega_0 n} x(n) \qquad \text{QED}$$

(5) Time reversal corresponds to frequency reversal. Given $F\{x(n)\} = X(\omega)$, then $F\{x(-n)\} = X(-\omega)$.

Proof We have

$$F\{x(-n)\} = \sum_{n = -\infty} x(-n) e^{-j\omega n}$$

x

On the right hand side set m = -n so that the limits $n = -\infty$ to ∞ change to $m = \infty$ to $-\infty$, and

$$F\{x(-n)\} = \sum_{m=\infty}^{\infty} x(m) e^{j\omega m}$$

Since this is a summation the limits can be written in reverse order, and we have

$$F\{x(-n)\} = \sum_{m = -\infty} x(m) \ e^{-j(-\omega)m} = X(-\omega) = X \ (e^{-j\omega})$$
QED

(6) Differentiation in frequency $F\{n x(n)\} = j \frac{dX(e^{j\omega})}{d\omega}$. Since

$$X(e^{j\omega}) = \sum_{n = -\infty} x(n) e^{-j\omega n}$$

we differentiate both sides with respect to ω to get $W(z^{j\omega}) \sim \int dz^{-j\omega n} dz^{-j\omega n} dz^{-j\omega n}$

$$\frac{dX(e^{j\omega})}{d\omega} = \sum_{n=-\infty}^{\infty} x(n) \frac{de^{-j\omega n}}{d\omega} = \sum_{n=-\infty}^{\infty} x(n)(-jn) e^{-j\omega}$$

or

$$j\frac{dX(e^{j\omega})}{d\omega} = \sum_{n=-\infty}^{\infty} n x(n) e^{-j\omega n} = F\{n x(n)\}$$
 QED

(7) Convolution If y(n) represents the convolution of two discrete-time signals x(n) and h(n), that is, y(n) = x(n)*h(n), then

$$Y(e^{j\omega}) = F\{x(n)*h(n)\} = X(e^{j\omega}) \cdot H(e^{j\omega})$$

From the definition of the Fourier transform

$$Y(e^{j\omega}) = \sum_{n=-\infty}^{\infty} y(n) \ e^{-j\omega n} = \sum_{n=-\infty}^{\infty} \{x(n)^* \ h(n)\} e^{-j\omega n}$$
$$= \sum_{n=-\infty}^{\infty} \{\sum_{k=-\infty}^{\infty} h(k)x(n-k)\} e^{-j\omega n}$$

Interchanging the order of summation

$$Y(e^{j\omega}) = \sum_{k = -\infty} h(k) \left\{ \sum_{n = -\infty}^{l} x(n-k) e^{-j\omega n} \right\}$$

00

The inner sum (I.S.) is handled thus: Let $(n-k) = \lambda$. Then as *n* goes from $-\infty$ to ∞ , λ goes from $-\infty$ to ∞ as well. Further $n = \lambda + k$. Thus the inner sum becomes

)

$$\text{I.S.} = \sum_{n = -\infty}^{\infty} x(n-k) \ e^{-j\omega n} = \sum_{\lambda = -\infty}^{\infty} x(\lambda) \ e^{-j\omega(\lambda+k)} = \left| \left(\sum_{\lambda = -\infty}^{\infty} x(\lambda) \ e^{-j\omega\lambda} \right) \right|_{\lambda=-\infty} e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{n=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega k}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ e^{-j\omega k} = e^{-j\omega k} X \ (e^{-j\omega}) = \sum_{\lambda=-\infty}^{\infty} x(\lambda) \ (e^{-j\omega k}) = \sum_{$$

00

Thus we have

$$Y(e^{j\omega}) = \sum_{k = -\infty} h(k) \left\{ e^{-j\omega k} X(e^{j\omega}) \right\} = X(e^{j\omega}) \quad \sum_{k = -\infty} h(k) e^{-j\omega k} = X(e^{j\omega}) \cdot H(e^{j\omega})$$

The function $H(e^{j\omega})$ is referred to as the **frequency response** of the system.

(8) Multiplication of two sequences Let y(n) be the product of two sequences $x_1(n)$ and $x_2(n)$ with transforms $X_1(e^{j\omega})$ and $X_2(e^{j\omega})$, respectively. Then

$$Y(e^{j\omega}) = F\{x_1(n), x_2(n)\} = X_1(e^{j\omega}) * X_2(e^{j\omega})$$
$$= \frac{1}{2\pi} \int_{\pi} X_1^{j\theta} X(e^{j(\omega-\theta)}) d\theta$$
$$\begin{pmatrix} e_1 & 2\\ 2\pi & 2 \end{pmatrix}$$

This is called **periodic convolution** since both $X(e^{j\omega})$ and $X(e^{j\omega})$ are periodic functions.

(9) **Parseval's Theorem** If x(n) and $X(e^{j\omega})$ are a Fourier transform pair, then

$$\sum_{n=-\infty}^{\infty} |x(n)|^2 = \frac{1}{2\pi} \int_{2\pi} |X(e^{j\omega})|^2 d\omega$$

Proof The energy *E* of the discrete-time sequence is defined as

$$E = \sum_{\substack{n = -\infty \\ n = -\infty}}^{\infty} |x(n)|^{2} = \sum_{\substack{n = 1 - \infty \\ m = 1 - \infty}}^{\infty} x(n) x^{*}(n)$$

Making the substitution $x^{*}(n) = \frac{1 - \infty}{2\pi} x^{*} (\frac{j\omega}{e}) e^{-j\omega n} d\omega$, we get
$$E = \sum_{\substack{n = -\infty \\ n = -\infty}}^{\infty} x(n) \left\{ \frac{1}{2\pi} \int_{-\pi}^{\pi} X^{*}(e^{j\omega}) e^{-j\omega n} d\omega \right\}$$

Interchanging the order of integration and summation,

73 of 77

$$E = \frac{1}{2\pi} \int_{-\pi}^{\pi} X^{*} \left(e^{j\omega} \right) \left(\sum_{n = -\infty}^{\infty} x(n) e^{-j\omega n} \right) d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(e^{j\omega} \right) X \left(e^{j\omega} \right) d\omega$$
$$= X \left(e^{j\omega} \right)$$
$$= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left| X \left(e^{j\omega} \right)^{2} d\omega \right| QED$$
$$QED$$

(10) Scaling – expansion in time For continuous-time Fourier transforms we have the property that if $x(t) \leftrightarrow X(\Omega)$, then $x(at) = \begin{bmatrix} 1 \\ a \end{bmatrix}$. This amounts to time-compression if a > 1 and to time remaining if x = 1.

time-expansion if a < 1.

The relation between time- and frequency-scaling in discrete time takes on a songewhat different form though. However, there is a result that does closely parallel $x(at) = -\frac{1}{a} \frac{X}{a} = \frac{1}{a}$.

Define y(n) = x(n/k), where *k* is a positive integer, as

y(n) = x(n/k), if *n* is a multiple of *k* 0, if *n* is not a multiple of *k*

This is time expansion. For example, if k = 3, then y(n) is obtained from x(n) by placing k-1 = 2 zeros between successive values of x(n). That is,

$$y(0) = x(0),$$
 $y(1) = y(2) = 0,$ $y(3) = x(1),$ $y(4) = y(5) = 0,$... Etc.

Since y(n) equals zero unless *n* is a multiple of *k*, i.e., unless n = rk, where r = all integers from $-\infty$ to $+\infty$, we have the Fourier transform of y(n) as:

$$Y(e^{j\omega}) = \sum_{n = -\infty}^{\infty} y(n) e^{-j\omega n} = \sum_{r = -\infty}^{\infty} y(rk) e^{-j\omega rk}$$

Since y(rk) = x(rk/k) = x(x), we have

$$Y(e^{j\omega}) = \sum_{r=-\infty} x(r) e^{-j\omega r k} = \sum_{\alpha} x(r) e^{-j(k\omega)r} = X(e^{jk\omega})$$

Thus we have the relation: Given $x(n) - X(e^{j\omega})$, then $x(n/k) - X(e^{jk\omega})$.

Example 1.8.1 [2003] (Frequency response) A discrete-time system is given by

$$y(n) - 5 y(n-1) = x(n) + 4 x(n-1)$$

where x(n) is the input and y(n) is the output. Determine its magnitude and phase response as a function of frequency.

Solution The frequency response is obtained by taking the discrete-time Fourier transform of both sides of the difference equation:

$$F\{y(n) - 5 y(n-1)\} = F\{x(n) + 4 x(n-1)\}$$

With $F\{y(n)\} = Y(e^{j\omega})$ and $F\{y(n-k)\} = e^{-j\omega k} Y(e^{j\omega})$ the above equation becomes

$$Y(e^{j\omega}) - 5 e^{-j\omega} Y(e^{j\omega}) = X(e^{j\omega}) + 4 e^{-j\omega} X(e^{j\omega})$$

$$Y(e^{j\omega}) (1 - 5 e^{-j\omega}) = X(e^{j\omega}) (1 + 4 e^{-j\omega})$$

$$Y(e^{j\omega}) = H(e^{j\omega}) = \frac{1 + 4e^{-j\omega}}{-j\omega} = \sqrt{\frac{(1 + 4\cos\omega)^2 + (4\sin\omega)^2 e^{-j(\frac{4\sin\omega}{1 + 4\cos\omega})}{2}}{\sqrt{(1 - 5\cos\omega)^2 + (4\sin\omega)^2 e^{-(\frac{4\sin\omega}{1 - 5\cos\omega})}}}}$$

$$= \sqrt{\frac{(1 + 4\cos\omega)^2 + (4\sin\omega)^2}{(1 - 5\cos\omega)^2 + (5\sin\omega)^2 e^{-(\frac{1 - 5\cos\omega}{1 - 5\cos\omega})}}}{(1 - 5\cos\omega)^2 + (5\sin\omega)^2 e^{-(\frac{1 - 5\cos\omega}{1 - 5\cos\omega})}}}$$

With the notation $H(\omega) = |H(\omega)| e^{j \angle H(\omega)}$ we can identify the magnitude and phase, respectively, as follows:

$$|H(\omega)| = \sqrt{\frac{(1+4\cos\omega)^2 + (4\sin\omega)^2}{(1-5\cos\omega)^2 + (5\sin\omega)^2}} \sqrt{\frac{17+8\cos\omega}{26-10\cos\omega}} \angle H(e^{j\omega}) = -\{\tan^{-1}|\underline{\Box}| + \tan^{-1}|\underline{\Box}|\}} \lfloor (1+4\cos\omega) (1-5\cos\omega) \rfloor$$

The MATLAB program follows for $H(\omega) = \frac{1 + 4e^{-j\omega}}{1 - 5e^{-j\omega}}$. Note that in MATLAB we use "w"

for ω and the plot ranges from -2π to 2π . Both ω and the phase, $\angle X(\omega)$, are in radians.

subplot(2,1,1);fplot('abs((1+4*exp(-j*w))/(1-5*exp(-j*w)))', [-2*pi,2*pi], 'k'); xlabel('Omega');ylabel('Magnitude'); subplot(2,1,2);fplot('angle((1+4*-exp(-j*w))/(1-5*exp(-j*w)))', [-2*pi,2*pi], 'k'); xlabel('Omega');ylabel('Phase');



Example 1.8.2 [Convolution] If $X(e^{j\omega}) = \text{DTFT}\{x(n)\}$ and y(n) = x(n)*x(-n) then from Properties 5 and 7 $Y(e^{j\omega}) = \text{DTFT}\{y(n)\} = \text{DTFT}\{x(n)*x(-n)\} = X(e^{j\omega}) X(e^{-j\omega})$. This result may also be obtained from the defining equation

$$Y(e^{j\omega}) = \sum_{n = -\infty}^{\infty} \{x(n)^* x(-n)\} e^{-j\omega n}$$

The convolution within the braces is $x(n)^*x(-n) = \sum_{k=-\infty}^{\infty} x(k)x(-(n-k))$, so that

$$Y(e^{j\omega}) = \sum_{n = -\infty} \left\{ \sum_{k = -\infty}^{n} x(k)x(k-n) \right\} e^{-j\omega n}$$
$$= \sum_{k = -\infty}^{\infty} x(k) \sum_{n = -\infty}^{\infty} x(k-n)e^{-j\omega n}$$

Set k-n = m etc.

Example 1.8.3 [2009] [Convolution] If $X(e^{j\omega}) = \text{DTFT}\{x(n)\}$ and $y(n) = x(n) * x^*(-n)$ find $Y(e^{j\omega}) = \text{DTFT}\{y(n)\}.$

Example 1.8.4 [2008] Find a difference equation to implement a filter with unit sample response $h(n) = (1/4)^n \cos(n\pi/3)u(n)$

Solution In the math manipulation it would help to write $\cos(n\pi/3)$ in its exponential form. You should also memorize the two standard *z*-transform pairs: $Z\{\cos(\omega_0 n) u(n)\}$ and $Z\{a^n \cos(\omega_0 n) u(n)\}$.

Hint Either use DTFT: Find the DTFT $F\{h(n)\} = H(e^{j\omega}) = \frac{Nr(\omega)}{\int_{j\omega} Dr(\omega)}$. *I*) Based on the relation $\frac{Y(e^{j\omega})}{X(e^{j\omega})} = H(e^{j\omega}) \operatorname{set} \frac{Y(e)}{X(e^{j\omega})} = \frac{Nr(\omega)}{Dr(\omega)}$,

2) Cross-multiply to get $Y(e^{j\omega}) Dr(\omega) = X(e^{j\omega})Nr(\omega)$,

3) Take inverse DTFT using time-shifting property (see below), and

4) Rearrange terms.

Or, use *z*-transforms: Find the *z*-transform $Z\{h(n)\} = H(z) = \frac{Nr(z)}{Dr(z)}$.

- 1) Based on the relation $\frac{Y(z)}{X(z)} = H(z) \operatorname{set} \frac{Y(z)}{X(\overline{z})} \frac{Nr(z)}{Dr(z)}$
- 2) Cross-multiply to get Y(z) Dr(z) = X(z) Nr(z),

3) Take inverse *z*-transform using time-shifting property, and

4) Rearrange terms.

Example 1.8.5 [2008] Find the inverse DTFT of $X(e^{j\omega}) = \frac{1}{1 - (1/3) e^{-j10\omega}}$.

Hint Use the result, derived in an earlier example, that for the exponential sequence $x(n) = a^n u(n)$, |a| < 1, the DTFT is

$$X(e^{j\omega}) = \sum_{n=0}^{\infty} a^n e^{-j\omega n} = \sum_{n=0}^{\infty} (ae^{-j\omega})^n = \frac{1}{1 - ae^{-j\omega}} \text{ with } a = 1/3$$
76 of 77

Then use the scaling property that if $x(n) - X(e^{j\omega})$, then $x(n/k) - X(e^{jk\omega})$ with k = 10. Alternatively, use the defining equation

$$\mathcal{F}^{-1}\{X(\omega)\} = x(n) = \frac{1}{2\pi} \int_{-\pi}^{\infty} X(e^{j\omega}) e^{j\omega n} d\omega$$

Example 1.8.6 [Scaling – compression in time] Given $x(n) - X(e^{j\omega})$ and y(n) = x(2n), find $Y(e^{j\omega})$.

Solution This is a *specific* result, not a general property. We have

$$Y(e^{j\omega}) = \sum_{\substack{n=-\infty \\ n=-\infty}}^{\infty} y(n) e^{-j\omega n} = \sum_{\substack{n=-\infty \\ n=-\infty}}^{\infty} x(2n) e^{-j\omega n} = \underbrace{\sum_{\substack{n=-\infty \\ n=-\infty}}^{\infty} x(n) e^{-j\omega n/2}}_{2n=-\infty} \sum_{\substack{n=-\infty \\ n=-\infty}}^{\infty} x(n) (-1) e^{-j\omega n/2} = \underbrace{\sum_{\substack{n=-\infty \\ n=-\infty}}^{\infty} x(n) e^{-j\omega n/2} + \frac{1}{2} \sum_{\substack{n=-\infty \\ n=-\infty}}^{\infty} x(n) (-1) e^{-j\omega n/2} = \underbrace{\sum_{\substack{n=-\infty \\ n=-\infty}^{\infty} x(n) (-1) e^{-j\omega n/2} = \underbrace{\sum_{\substack{$$

Notationally, we may write this as

$$Y\left(e^{j\omega}\right) = \frac{1}{2}X\left(e^{j\omega/2}\right) \frac{1}{2}\left(-e^{j\omega/2}\right)$$

77 of 77

Digital Signal Processing – 1(B)

I(B). The z-transform and realization of digital filters

Review of z-transforms, Applications of z-transforms, Solution of difference equations of digital filters, Block diagram representation of linear constant-coefficient difference equations, Basic structures of IIR systems, Transposed forms, Basic structures of FIR systems, System function.

Contents: Introduction Important properties of z-transforms Transforms of some useful sequences Region of convergence and stability Inverse z-transform by partial fractions Relationships among system representations Inverse z-transform by power series expansion (long division) Computation of frequency response Z-transforms with initial conditions Steady-state and transient responses for a first order system Realization of digital filters The Lattice structure – Introduction *Inverse z-transform by complex inversion integral

Introduction

For continuous-time systems the **Laplace transform** is an extension of the **Fourier transform**. The Laplace transform can be applied to a broader class of signals than the Fourier transform can, since there are many signals for which the Fourier transform does not converge but the Laplace transform does. The Laplace transform allows us, for example, to perform transform analysis of unstable systems and to develop additional insights and tools for LTI system analysis.

The *z*-transform is the discrete-time counterpart of the Laplace transform. The *z*-transform enables us to analyze certain discrete-time signals that do not have a **discrete-time Fourier transform**. The motivations and properties of the *z*-transform closely resemble those of the Laplace transform. However, as with the relationship of the continuous time versus the discrete-time Fourier transforms, there are distinctions between the Laplace transform and the *z*-transform.

Definition The two-sided (bilateral) *z*-transform, X(z), of the sequence x(n) is defined as

$$X(z) = \Im\{x(n)\} = \sum_{n = -\infty} x(n) z^{-n}$$

where $z = r e^{j\omega}$ is the complex variable. The above power series is a **Laurent series**.

The one-sided (unilateral) z-transform is defined as

$$X_+(z) = \sum_{n=0}^{\infty} x(n) z^{-n}$$

The unilateral *z*-transform is particularly useful in analyzing causal systems specified by linear constant-coefficient difference equations with nonzero initial conditions into which inputs are stepped. It is extensively used in digital control systems.



The region of convergence (ROC) is the set of z values for which the above summation converges. In general the ROC is an annular region in the complex z-plane given by

ROC =
$$R_{x-} < |z| < R_{x+}$$

Relationship between the *z***-transform and the discrete-time Fourier transform** Setting $z = r e^{j\omega}$ in the definition gives us

$$X(z)\Big|_{z=re^{j\omega}} = \sum_{n=-\infty}^{\infty} x(n) \Big(re^{j\omega}\Big)^{-n} = \sum_{n=-\infty}^{\infty} [r^{-n} x(n)]e^{-j\omega n}$$

If |z| = r = 1, then the *z*-transform, evaluated on the unit circle, gives the discrete-time Fourier transform of the sequence x(n), i.e.,

$$X(e^{j\omega}) = \sum_{n=-\infty}^{\infty} x(n)e^{-j\omega n} = X(z)\Big|_{z=e^{j\omega}}$$

Example 1.1.1 The positive-time signal

$$x(t) = e^{-\alpha t}, t \ge 0$$

0, otherwise

is sampled at *T*-second intervals resulting in the sequence x(nT) or x(n)

$$x(n) = e^{-\alpha t} \Big|_{t=nT} = e^{-\alpha nT} = \left(e^{-\alpha T}\right)^n = a^n, \ a = e^{-\alpha T} \text{ and } n \ge 0$$
$$x(n) = \begin{array}{c} a^n, & n \ge 0\\ 0, & n < 0 \end{array}$$

If a < 1 this sequence decays exponentially to 0 as $n \to \infty$. Substituting x(n) into the defining equation, the *z*-transform is

$$\Im\{x(n)\} = X(z) = \sum_{n=0}^{\infty} a^n z^{-n} = \sum_{\underline{z}=0}^{\infty} (az^{-1})^n = \underline{1-az^{-1}}, \qquad |az^{-1}| < 1$$
$$= \frac{1}{1-az^{-1}} = \frac{z}{z-a}, \qquad |z| > |a|$$



The ROC is |z| > |a|. This X(z) is a rational function (a ratio of polynomials in *z*). The roots of the numerator polynomial are the zeros of X(z) and the roots of the denominator polynomial are the poles of X(z).

This is a right-sided sequence. *Right-sided sequences have a ROC that is the exterior of a circle* with radius $R_{x-}(|z| > |a|)$ in this case). If the ROC is the exterior of a circle it is a right-sided sequence.

Definition A **right-sided sequence** x(n) is one for which x(n) = 0 for all $n < n_0$ where n_0 is positive or negative but finite. If $n_0 \ge 0$ then x(n) is a *causal* or *positive-time sequence*.

Example 1.1.2 The negative-time sequence $x(n) = -b^n u(-n-1)$. Recall that the unit step sequence u(.) = 1 if the argument of u(.) is ≥ 0 , i.e., if $(-n-1) \ge 0$ or $n \le -1$.

$$x(n) = -b^n, \quad n \le -l$$

0, otherwise

If b > 1 this sequence decays exponentially to 0 as $n \to -\infty$. The *z*-transform is,

$$\mathbf{z}\{x(n)\} = X(z) = \sum_{n = -\infty} x(n) z^{-n} = \sum_{n = -\infty} b^n z^{-n} = -\sum_{n = -\infty} (b z^{-1})^n$$

Let n = -m and change the limits accordingly to get,

$$X(z) = -\sum_{m=\infty} (bz^{-1})^{-m} = 1 - \sum_{m=0} (b^{-1}z)^{m}$$

We added 1 in the last step above to make up for the m = 0 term within the summation. The result is,

$$X(z) = 1 - \frac{1}{(1 - zb^{-1})}, \quad |z \ b^{-1}| < 1$$
$$= \frac{z}{z - b}, \quad \text{ROC is } |z| < |b|$$



This is a left-sided sequence. Such a sequence has a region of convergence which is the interior of a circle, $|z| < R_{x+}$. In this case the ROC is |z| < |b|.

Note that if b = a then the two examples above have exactly the same X(z). So what makes the difference? The region of convergence makes the difference.

Definition A left-sided sequence x(n) is one for which x(n) = 0 for all $n \ge n_0$, where n_0 is positive or negative but finite. If $n_0 \le 0$ then x(n) is an *anticausal* or *negative-time sequence*.

Example 1.1.3 [Two-sided sequence] This is the sum of the positive- and negative-time sequences of the previous two examples.

$$y(n) = a^n, \quad n \ge 0 \qquad \qquad = a^n u(n) - b^n u(-n-1)$$
$$-b^n, \quad n < 0$$

Substituting into the defining equation,

$$Y(z) = \mathbf{z}\{y(n)\} = \sum_{n=-\infty}^{\infty} [a^n u(n) - b^n u(-n-1)] z^{-n} = \sum_{n=0}^{\infty} a^n z^{-n} - \sum_{n=0}^{-1} b^n z^{-n}$$

Now, from Examples 1 and 2,

$$\sum_{n=0}^{\infty} a^n z^{-n} = \frac{z}{z-a} \quad (\text{ROC } |z| > |a|) \qquad \& \qquad b^n = \frac{z}{z-b} \quad (\text{ROC } |z| < |b|)$$

-1

So, the desired transform Y(z) has a region of convergence equal to the intersection of the two separate ROC's |z| > |a| and |z| < |b|. Thus

$$Y(z) = \frac{z}{z-a} + \frac{z}{z-b}, \text{ with ROC } \{|z| > |a|\} \cap \{|z| < |b|\}$$

= $\frac{z(2z-a-b)}{(z-a)(z-b)}, \text{ with ROC } |a| < |z| < |b|$

The ROC is the overlap of the shaded regions, that is, the annular region between |a| and |b|. The two zeros are at 0 and (a+b)/2, and the two poles at *a* and *b*.



If |b| < |a| the transform does not converge.



In the above three examples we may express the *z*-transform both as a ratio of polynomials in *z* (i.e., positive powers) and as a ratio of polynomials in z^{-1} (negative powers). From the definition of the *z*-transform, we see that for sequences which are zero for n < 0, X(z) involves only negative powers of *z*. However, reference to the poles and zeros is always in terms of the roots of the numerator and denominator expressed as polynomials in *z* (i.e., positive power of *z*), as having poles at infinity if the degree of the numerator exceeds the degree of the denominator or zeros at infinity if the numerator is of smaller degree than the denominator.

Example 4.1.4 [Finite-length sequence] Only a finite number of sequence values are non-zero, as given below.

x(n) = 0 for $n < N_1$ and for $n > N_2$, where N_1 and N_2 are finite non-zero for $N_1 \le n \le N_2$

By the defining equation we have

$$\sum_{X(z)=n=N_1}^{N_2} x(n) z^{-n} = x(N^1) z^{-N_1} + \dots + x(N^2) z^{-N_2}$$

Convergence of this expression requires simply that $|x(n)| < \infty$ for $N_1 \le n \le N_2$. Then *z* may take on all values except $z = \infty$ if N_1 is negative and z = 0 if N_2 is positive. Thus the ROC is at least $0 < |z| < \infty$ and it may include either z = 0 or $z = \infty$ depending on the sign of N_1 and N_2 .

Important properties of *z***-transforms**

The proofs are easily obtained by using the basic *z*-transform definition and transformations in the summation. [Sec 2.3 Oppenheim & S]

(1) Linearity If $\mathfrak{Z}[x(n)] = X(z)$ with ROC $r_{x1} < |z| < r_{x2}$ and $\mathfrak{Z}[y(n)] = Y(z)$ with ROC $r_{y1} < |z| < r_{y2}$ then $\mathfrak{Z}[a x(n) + b y(n)] = a X(z) + b Y(z)$ with ROC at least the overlap of the ROC's of X(z) and Y(z). If there is any pole-zero cancellation due to the linear combination, then the ROC may be larger.

(2) Translation (Time-shifting) If $\mathfrak{Z}[x(n)] = X(z)$ with ROC $r_1 < |z| < r_2$ then $\mathfrak{Z}[x(n-k)] = z^{-k}X(z)$ with the same ROC except for the possible addition or deletion of z = 0 or $z = \infty$ due to z^{-k} .

Example Given $x(n) = \{1, 2\}$ and $x_2(n) = x(n+2)$ find X(z) and $X_2(z)$ and their respective ROCs. $X(z) = 1 + 2z^{-1}$, ROC: entire *z*-plane except z = 0; $X_2(z) = z^2(1+2z^{-1}) = z^2+2z$, ROC: entire *z*-plane except $z = \infty$.

(3) Multiplication by a complex exponential sequence (Scaling in the *z*-domain) If $\underline{z}[x(n)] = X(z)$ with ROC $r_1 < |z| < r_2$ then $\underline{z}[a^n x(n)] = X(z)|_{z \to (z/a)}$ with ROC $|a| r_1 < |z| < |a| r_2$. **Example** Given $x(n) = \{1, 2\}$ and $x_2(n) = 0.5^n x(n-2)$ find X(z) and $X_2(z)$ and their respective ROCs.

(4) Multiplication by a ramp If $\Im[x(n)] = X(z)$ with ROC $r_1 </z/< r_2$ then $\Im[n x(n)] = -z$ with ROC r </z/< r.

Example Given $x(n) = \{1, 2\}$ and $x_2(n) = (1 + n + n^2) x(n)$ find X(z) and $X_2(z)$ and their respective ROCs. $\Im[x_2(n)] = \Im[1] + \Im[n x(n)] + \Im[n [n x(n)]] = ...$

(5) Time reversal If $\underline{z}[x(n)] = X(z)$ with ROC $r_1 < |z| < r_2$ then $\underline{z}[x(-n)] = X(z^{-1})$ with ROC $(1/r_2) < |z| < (1/r_1)$

Example Given $x(n) = 2^{-n} u(n)$ and $X(z) = z X(z^{-1})$ determine $x_2(n)$. $x(n) = (0.5)^n u(n), X(z) = \frac{z}{z - 0.5}, \text{ROC: } 0.5 </z/$ $\mathbf{z}^{-1}\{X(z^{-1})\} = x(-n); x_2(n) = \mathbf{z}^{-1}\{z X(z^{-1})\} = x(-(n+1)) = 2^{-(-(n+1))} u(-(n+1))$

(6) Convolution in time domain leads to multiplication in frequency domain Given $[\mathfrak{z}(n)] = X(z)$ with ROC $z \in R_x$ and $\mathfrak{z}[y(n)] = Y(z)$ with ROC $z \in R_y$ and $x(n)*y(n) = \sum_{k=-\infty}^{\infty} x(k)h(n-k)$ then

 $\mathbf{g}[x(n)^*y(n)] = X(z). Y(z) \text{ with ROC } z \in R_x \cap R_y.$

(7) Multiplication in time domain leads to convolution in frequency domain If $\mathfrak{Z}[x(n)] = X(z)$ with ROC $r_{x1} < |z| < r_{x2}$ and $\mathfrak{Z}[y(n)] = Y(z)$ with ROC $r_{y1} < |z| < r_{y2}$ then

$$\mathbf{z}[x(n).y(n)] = \frac{1}{12}\pi \oint_{C_2} X(v) Y\left(\frac{z}{v}\right) v^{-1} dv, \qquad \text{ROC } r^{x_1} r_{y_1} < |z| < r_{x_2} r_{y_2}.$$

where \oint_{C_2} is a complex contour integral and C_2 is a closed contour in the intersection of the ROCs of X(v) and Y(z/v).

(8) Initial Value Theorem If x(n) is a causal sequence with z transform X(z), then $x(0) = \lim_{z \to \infty} X(z)$

(9) Final Value Theorem If $\mathfrak{Z}[x(n)] = X(z)$ and the poles of X(z) are all inside the unit circle then the value of x(n) as $n \to \infty$ is given by

$$x(n)\Big|_{n\to\infty} = \lim_{z\to 1} [(z-1) X(z)]$$

Some also give this as $x(n) \Big|_{n \to \infty} = \lim_{z \to 1} \left[(1 - z^{-1}) X(z) \right]$

Transforms of some useful sequences

(1) The **unit sample** $\delta(n)$:

$$\mathbf{z}[\delta(n)] = \sum_{n = -\infty}^{\infty} \delta(n) z^{-n} = 1. \ z^0 = 1, \qquad \text{ROC all } z$$

Delayed unit sample $\delta(n-k)$:

$$\mathbf{J}[\delta(n-k)] = \sum_{n=-\infty}^{\infty} \delta(n-k) \ z^{-n} = 1. \ z^{-k} = z^{-k}, \qquad \text{ROC } |z| > 0 \text{ if } k > 0 \ (|z| < \infty \text{ if } k < 0)$$

(2) **Unit step** *u*(*n*) (positive time):

$$\mathbf{z}[u(n)] = \sum_{n=-\infty}^{\infty} u(n) \ z^{-n} = \sum_{n=0}^{\infty} 1 \ z^{-n} = 1 + z^{-1} + z^{-2} + \dots$$
$$= \frac{1}{1 - z^{-1}} = \frac{z}{z - 1}, \qquad \text{ROC} \ |z^{-1}| < 1 \ \text{or} \ |z| > 1$$

(3) **Unit step** -u(-n-1) (negative time):

$$\mathbf{z}[-u(-n-1)] = \sum_{n=-\infty}^{\infty} -u(-n-1)z^{-n} = \sum_{n=-\infty}^{1} (-1)z^{-n} = -\sum_{n=-\infty}^{1} z^{-n} = 1 - \sum_{n=-\infty}^{0} z^{-n}$$
$$= 1 - (1 + z + z^{2} + ...) = 1 - \frac{1}{1-z} = \frac{z}{z-1}, \qquad \text{ROC } |z| < 1$$

(4) **Exponential** $a^n u(n)$, derived in earlier example:

$$\Im[a^n u(n)] = \frac{z}{z-a}, \quad \text{ROC} |z| > |a|$$

(5) **Exponential** $-b^n u(-n-1)$; negative time; derived earlier:

$$\mathbf{J}[-b^n u(-n-1)] = \frac{z}{z-b}, \qquad \text{ROC } |z| < |b|$$

(6) **Unit ramp** n u(n). Given that $\Im[u(n)] = U(z) = \frac{z}{z^{z-1}}$ $\Im[n u(n)] = -z \quad \frac{dU(z)}{dz} = -z \quad \frac{d}{z^{z-1}} = \frac{-z[(z-1).1-z.1]}{(z-1)^{2}}$ $= \frac{z}{(z-1)^{2}} \quad \text{ROC } |z| > 1$, same as that of U(z)(7) **Sinusoid** $\sin \omega_{0} n u(n)$: $\Im\{\sin \omega_{0} n u(n)\} = \frac{z \sin \omega_{0}}{z^{2} - 2z \cos \omega_{0} + 1}$, ROC |z| > 1

$$\begin{aligned} \mathbf{J}\left\{\begin{array}{l} \boldsymbol{\omega} \\ \mathbf{\delta} \\ \mathbf{\delta} \\ \sin \ _{0}n \ u(n) \\ = \sum \limits_{n=0} \limits_{\substack{n=0 \\ j \neq n}} sin \ _{0}n \ _{z^{-n}} = \sum \limits_{\substack{n=0 \\ j \neq n}} \left[\frac{1}{j^{2}} \sum \limits_{\substack{n=0 \\ j \neq n}} \frac{1}{j^{2}} \sum \limits_{\substack{n=0 \\ j \neq n}} \frac{1}{j^{2}} \sum \limits_{\substack{n=0 \\ i \neq n}} \frac{1}{j^{2}} \sum \atop_{\substack{n=0 \\ i \neq$$

Region of convergence and stability

Suppose x(n) is a causal sequence that can be written as a sum of complex exponentials. This takes in a wide class of signals including sinusoids, exponentials, and products thereof. Let

$$x(n) = \sum_{i=1}^{N} a_i^n u(n)$$

Taking the *z* transform of x(n) gives

$$\frac{\mathbf{z}[x(n)] = X(z) = \sum_{i=1}^{N} \frac{z}{z - a_i}$$

The region of convergence R is the intersection of the regions of convergence for each exponential as follows:

$$R = \bigcap_{i=1}^{N} R_i \text{ where } R_i = \{z: |z| > |a_i|\}$$

Therefore, $R = \{z: |z| > \text{largest of } |a_i|\}$ as shown here (Figure)



Since the ROC for a translated exponential remains the same as that for the original exponential, all right-sided sequences that are sums of translated exponentials have ROCs similar to that expressed above.

By a similar argument all left-sided sequences expressible as a sum of translated complex exponentials have a ROC, *L*, given by

 $L = \{z: |z| < \text{smallest of } |b_i|\}$

If we have a combination of right- and left-sided sequences, the corresponding ROC is the intersection of R and L. Therefore the total ROC becomes an annular region as shown below and given by

$$R_{Total} = R \cap L = \{z: \text{Largest of } |a_i| < |z| < \text{smallest of } |b_i|\}$$



The stability of a system with an impulse response that is the sum of translated right- and left-sided sequences can be determined from the region of convergence. Assume that h(n) is the

unit sample response of a causal or non-causal linear shift-invariant system. Let $\underline{J}[h(n)] = H(z)$, the so-called system function. Then:

Theorem A linear shift-invariant system with system function H(z) is BIBO stable if and only if the ROC for H(z) contains the unit circle.

This theorem can be used to determine stability for a given H(z) without obtaining the impulse response or checking outputs for all bounded input signals.

Illustration of stability and causality For A system function with 2 poles at, say, z = 0.5, and z = 1.5, there are three possible regions of convergence.

(1) ROC is 0.5 < |z| < 1.5. Here the system is stable since the unit circle is inside the region of convergence. The impulse response, h(n), is two-sided, so the system is noncausal.



(2) ROC is |z| < 0.5. Here the system is not stable. The impulse response, h(n), is left-sided, so the system is noncausal.



(3) ROC is |z| > 1.5. Here the system is not stable. The impulse response, h(n), is right-sided, so the system may be causal.



Inverse z-transform by partial fractions

(Aside) Comparison of inverse z-transform methods A limitation of the power series method is that it does not lead to a closed form solution (although this can be deduced in simple cases), but it is simple and lends itself to computer implementation. However, because of its recursive nature care should be taken to minimize possible build-up of numerical errors when the number of data points in the inverse z-transform is large, for example by using double precision.

Both the **partial fraction method** and the **inversion integral method** require the evaluation of residues albeit performed in different ways. The partial fraction method requires the evaluation of the residues of X(z) or $\frac{X(z)}{z}$. The complex inversion integral requires the evaluation of the residues of $X(z) z^{n-1}$. In many instances evaluation of the complex inversion integral is needlessly difficult and involved.

Both the partial fraction method and the inversion integral method lead to closed form solutions. The main disadvantage is having to factorize the denominator polynomial of X(z) when it is of order greater than 2. Another disadvantage is multiple order poles and the resulting differentiation(s) when determining residues.

The partial fraction method directly generates the coefficients of parallel structures for digital filters. The inversion integral method is widely used in the analysis of quantization errors in discrete-time systems.

(End of Aside)

As in Laplace transforms, in order to expand a rational function into partial fractions, the degree of the numerator should be less than the degree of the denominator – proper fraction. If it is not then we perform long division as below where Q(z) is the quotient and $N_I(z)$ is the remainder.

$$X(z) = \frac{N(z)}{D(z)} = Q(z) + \frac{N_1(z)}{D(z)}$$

Long Division $Q(z) \leftarrow \text{Quotient}$ Denominator D(z) $N(z) \leftarrow \text{Numerator}$ $N_{I}(z) \leftarrow \text{Remainder}$

The long division is done until we get a remainder polynomial $N_I(z)$ whose degree is less than the degree of the denominator D(z). We then obtain x(n) as

$$x(n) = \mathbf{z}^{-1} \{ X(z) \} = \mathbf{z}^{-1} \{ Q(z) \} + \mathbf{z}^{-1} \{ \begin{matrix} N(z) \\ 1 \\ D(z) \end{matrix} \}$$

Since $N_I(z)/D(z)$ is a proper fraction it can be expanded into partial fractions. The overall inverse transform is obtained by looking up a table of *z*-transform pairs.

However, there is an alternative available in the case of *z*-transforms which is not available in Laplace transforms. This is a result of the fact that *z*-transforms are characterized by a *z* in the numerator (as can be verified by looking at a table of *z*-transforms). Therefore, instead of expanding X(z) we may, instead, expand [X(z)/z] into partial fractions giving

$$\frac{X(z)}{z} = \frac{A}{z-z_1} + \frac{B}{z-z_2} + \dots$$

D-

so that X(z) is given by A_{z}

$$X(z) = \underbrace{\begin{array}{c} A z & D z \\ z - z_1 & z - z_2 \end{array}}_{z - z_1} + \ldots$$

This can be inverted by a simple look-up of a table of transforms. Note also that in some cases X(z) = N(z)/D(z) may not be a proper fraction but [X(z)/z] is and, therefore, this method obviates the need for long division of N(z) by D(z). (In still other cases even [X(z)/z] may not be a proper fraction. See later under "General procedure for partial fraction expansion".)

Example 1.5.1 (See also long division later). Find the inverse *z*-transform, using partial fractions, of

$$\frac{X(z) = \frac{2 z^2 - 3 z}{z^2 - 3 z + 2} = \frac{N(z)}{D(z)}$$

This is not a proper fraction since the degree of the numerator is not less than the degree of the denominator. However, (X(z)/z) is a proper fraction

$$\frac{X(z)}{z} = \frac{2z-3}{z^2-3z+2} = \frac{2z-3}{(z-1)(z-2)}$$

which has the partial fraction expansion

$$\frac{X(z)}{z} = \frac{1}{(z-1)} + \frac{1}{(z-2)} \qquad \text{or} \qquad X(z) = \frac{z}{(z-1)} + \frac{z}{(z-2)}$$

By looking up a table of *z*-transforms the inverse *z*-transform is

$$\mathbf{z}^{-1}{X(z)} = x(n) = u(n) + 2^n u(n)$$

Note that we are giving here the causal solution that corresponds to a ROC |z| > 2 (not 1 < |z| < 2 or |z| < 1) so that x(n) is a right-sided sequence.

The **alternative** method is to divide N(z) by D(z) as below (as is standard practice in Laplace transforms). Note that in this long division the numerator and denominator polynomials are arranged in the order of decreasing powers of z. There are three other ways (all of them wrong) of arranging the two polynomials for the long division.

Long Division
2
$$\leftarrow$$
 Quotient
Denominator $\rightarrow z^2 - 3z + 2$ $2z^2 - 3z \leftarrow$ Numerator
 $2z^2 - 6z + 4$
 $3z - 4$ \leftarrow Remainder

Thus X(z) can be expressed as

$$X(z) = 2 + \underbrace{3z-4}_{(z-1)(z-2)} = 2 + X_{1}(z) \quad \text{where } X_{1}(z) = \underbrace{3z-4}_{(z-1)(z-2)}$$

 $X_l(z)$ is a proper fraction and can be expanded into partial fractions as below:

$$X_{I}(z) = \frac{3z - 4}{(z - 1)(z - 2)} = \frac{A}{z - 1} + \frac{B}{z - 2}$$

Solving for *A* and *B* we get A = 1 and B = 2, so that X(z) may be written

$$X(z) = 2 + \frac{1}{z^{-1}} + \frac{2}{z^{-2}}$$

Taking the inverse z-transform we get

$$x(n) = \mathbf{3}^{-1} \{X(z)\} = \mathbf{3}^{-1} \left(2 + \frac{1}{z-1} + \frac{2}{z-2} \right)$$

= $\mathbf{3}^{-1} \{2\} + \mathbf{3}^{-1} \left(\frac{1}{2} \right) + \mathbf{3}^{-1} \left(\frac{2}{z-2} \right)$
A term like $\mathbf{3}^{-1} \left(\frac{1}{z-1} \right)$ is handled by writing ______ as _____. We know that $\mathbf{3}^{-1} \left(\frac{z}{z-1} \right) = u(n)$, so $\mathbf{3}^{-1} \left(\frac{z}{z-1} \right) = u(n-1)$

Similarly
$$\mathbf{z}^{-1} \left(\frac{2}{z-2} \right) = 2 \cdot 2^{n-1} u(n-1)$$
. Thus
 $x(n) = 2 \,\delta(n) + u(n-1) + 2 \cdot 2^{n-1} u(n-1)$

This can be verified to be equivalent to $x(n) = u(n) + 2^n u(n)$ obtained earlier.

In MATLAB (Partial fractions) The partial fractions may be computed by using the *residuez* function. In this method X(z) is arranged as a ratio of polynomials in negative powers of z and, in the denominator, the leading coefficient $a_0 \neq 0$. See "General procedure for partial fraction expansion" later.

$$X(z) = \frac{2 - 3 z^{-1}}{1 - 3z^{-1} + 2z^{-2}} = K + \frac{R_1}{1 - p_1 z^{-1}} + \frac{R_2}{1 - p_2 z^{-1}}$$

We define the coefficient vectors b = [2, -3] and a = [1, -3, 2]; $R = [R_1, R_2]$ represents the residues (partial fraction constants), $p = [p_1, p_2]$ the poles and *K* a constant.

%Partial fractions b = [2, -3], a = [1, -3, 2], [R, p, K] = residuez (b, a)

The MATLAB results returned are



The MATLAB output tells us that the poles are at z = 2 and z = 1 and the corresponding residues are, respectively, 1 and 1. Further K = 0. Therefore,

$$X(z) = \frac{2 - 3z^{-1}}{1 - 3z^{-1} + 2z^{-2}} = 0 + \frac{1}{1 - 2z^{-1}} + \frac{1}{1 - 1z^{-1}} = \frac{z}{z - 2} + \frac{z}{z - 1}$$

Note that the $X(z) = \frac{1}{1 - 2z^{-1}} + \frac{1}{1 - 1z}$ tobtained by the *residuez* function and $\frac{X(z)}{z} = \frac{1}{(z - 1)^{+}} + \frac{1}{(z - 2)}$ are

the same since X(z) has no repeated poles. This won't be the case if X(z) has repeated poles.

Example 1.5.2 Find if the discrete LTI system described by

$$y(n) - y(n-1) + 0.5 y(n-2) = x(n) + x(n-1)$$

is BIBO stable or not. Find its transfer function and impulse response. Sketch its pole-zero plot. **Solution** Take the *z*-transform of both sides:

 $\mathbf{z} \{ y(n) - y(n-1) + 0.5 \ y(n-2) \} = \mathbf{z} \{ x(n) + x(n-1) \}$ $Y(z) - z^{-1} Y(z) + 0.5 \ z^{-2} Y(z) = X(z) + z^{-1} X(z)$ $Y(z) \ (1 - z^{-1} + 0.5 \ z^{-2}) = X(z) \ (1 + z^{-1})$

$$H(z) = \frac{Y(z)}{X(z)} = \frac{1+z^{-1}}{1-z^{-1}+0.5z^{-2}} = \frac{(z+1)/z}{(z^{2}-z+0.5)/z^{2}} = \frac{z(z+1)}{z^{2}-z+0.5}$$

The denominator has roots at

$$z_{1}, z_{2} = \frac{-(-1)\pm\sqrt{(-1)^{2}-4.1.(0.5)}}{2} = \frac{1\pm\sqrt{1-2}}{2} = \frac{1\pm j1}{2} = 0.5\pm j0.5$$
$$= (0.5+j0.5) \text{ and } (0.5-j0.5)$$

Thus the transfer function H(z) has zeros at z = 0, z = 1, and poles at $z = 0.5 \pm j0.5$. For a causal system (right-sided sequence, h(n)) the region of convergence is |z| > |0.5 + j0.5| or |z| > 0.707. (Figure)



The impulse response is given by $h(n) = \mathbf{z}^{-1} \{H(z)\}$. We need partial fractions for H(z); we shall instead handle H(z)/z: $\frac{H(z)}{z+1} = \frac{z+1}{z+1}$

$$\frac{1}{z} = \frac{1}{z^2 - z + 0.5} = \frac{1}{(z - 0.5 - j0.5)(z - 0.5 + j0.5)}$$
$$= \frac{A}{z - 0.5 - j0.5} + \frac{A}{z - 0.5 + j0.5}$$

Solving for A and A^* , we get

$$A = \frac{z+1}{z-0.5+j0.5} \bigg|_{z=0.5+j0.5} = \frac{0.5+j0.5+1}{0.5+j0.5-0.5+j0.5} = \frac{1.5+j0.5}{j}$$
$$= 0.5-j1.5 = \sqrt{5/2} \ e^{-j\tan^{-1}3}$$
$$A^* = 0.5+j1.5$$

Thus we have

$$\frac{H(z)}{z} = \frac{A}{z - (0.5 + j0.5)} + \frac{A}{z - (0.5 - j0.5)}$$

$$H(z) = A \frac{z}{z - (0.5 + j0.5)} + A^* \frac{z}{z - (0.5 - j0.5)}$$

$$= a = b$$

٦

where

$$a = 0.5 + j0.5 = \sqrt{0.5^2 + 0.5^2} e^{j \tan^{-1} 1} = \frac{1}{\sqrt{2}} e^{j\pi/4}$$
$$b = 0.5 - j0.5 = \frac{1}{\sqrt{2}} e^{-j\pi/4}$$

The inverse *z*-transform is

$$h(n) = A a^{n} + A^{*} b^{n}, n \ge 0$$

0, otherwise

So that for $n \ge 0$,

$$\begin{aligned} 0, \\ h(n) &= A \left(\frac{1}{\sqrt{2}} e^{j} \right)^n + A^* \left(\frac{1}{\sqrt{2}} e^{-j\pi/4} \right)^n \\ &= A \left(\frac{1}{\sqrt{2}} \right)^n e^{j\pi n/4} + A^* \left(\frac{1}{\sqrt{2}} \right)^n e^{-j\pi n/4} \\ &= \left(\frac{1}{\sqrt{2}} \right)^n \left(1 \right)^n \frac{j\pi n/4}{j\pi n/4} + A^* \left(\frac{1}{\sqrt{2}} \right)^n e^{-j\pi n/4} \\ &= \left(\frac{1}{\sqrt{2}} \right)^n \left(1 \right)^n \frac{j\pi n/4}{j\pi n/4} + \left(1 \right)^n \frac{3}{j\pi n/4} + \left(\frac{1}{\sqrt{2}} \right)^n e^{j\pi n/4} \\ &= \left(\frac{1}{\sqrt{2}} \right)^n \left(e^{j\pi n/4} + e^{-j\pi n/4} \right) + j \frac{1}{2} \left(\sqrt{2} \right) \left(-e^{j\pi n/4} + e^{-j\pi n/4} \right) \\ &= \left(\frac{1}{\sqrt{2}} \right) \left(\frac{2}{2} \right) \left(\frac{1}{2} \right)^n e^{j\pi n/4} + e^{-j\pi n/4} \\ &= \left(\frac{1}{\sqrt{2}} \right)^n \cos(\pi n/4) + 3 \left(\frac{1}{\sqrt{2}} \right)^n \left[e^{j\pi n/4} - e^{-j\pi n/4} \right] \\ &= \left(\frac{1}{\sqrt{2}} \right)^n \cos(\pi n/4) + 3 \left(\frac{1}{\sqrt{2}} \right)^n \left[e^{j\pi n/4} - e^{-j\pi n/4} \right] \end{aligned}$$

`

To sum up,

$$h(n) = \begin{pmatrix} 1 \\ \sqrt{2} \end{pmatrix}^n [\cos(\pi n/4) + 3\sin(\pi n/4)], n \ge 0$$

0, otherwise

Alternatively, since the two terms in

$$h(n) = A \frac{e^{j\pi n/4}}{(\sqrt{2})^n} + A^* \frac{e^{-j\pi n/4}}{(\sqrt{2})^n}, \qquad n \ge 0$$

are complex conjugates of each other we can write $|(\rho^{jn\pi/4})|$

$$h(n) = 2 \operatorname{Re} \left[\begin{array}{c} A \\ A \\ \sqrt{2} \end{array} \right], \qquad n \ge 0$$

Alternative In the above solution the impulse response initially contains complex numbers; these have been algebraically manipulated into sine and cosine terms. A more direct way to obtain the impulse response in a form that contains no complex numbers is to use results #7 and #8 in "Transforms of some useful sequences" and manipulate the transform

$$H(z) = \frac{z(z+1)}{z^2 - z + 0.5} = \frac{z + z}{z^2 - z + 0.5}$$

into those forms. Comparing the denominator of H(z) with the denominator of the transforms of the sine and cosine functions

$$\mathbf{z}\left\{a^{n}\sin \omega n \ u(n)\right\} = \frac{az\sin \omega_{0}}{z^{2} - 2az\cos \omega + a^{2}}, \qquad \text{ROC } |z| > a$$

and

$$\mathfrak{Z}\left\{\begin{array}{ccc} a \cos & n \ u(n) = \frac{z^2}{z - az \cos \omega_0}, & \operatorname{ROC} / z/ > a \\ & & z^2 - 2az \cos \omega + a^2 \end{array}\right\}$$

we get

$$z^2 - z + 0.5 = z^2 - 2az\cos\omega_0 + a^2$$

from which

$$a^{2} = 0.5 \rightarrow a = \frac{1}{\sqrt{2}} \qquad 2a \cos \omega_{0} = 1 \rightarrow \omega_{0} = \pi/4,$$

$$\cos \omega_{0} = \frac{1}{\sqrt{2}} \qquad a \cos \omega_{0} = \frac{1}{2}, \sin \omega_{0} = \frac{1}{\sqrt{2}} \qquad a \sin \omega_{0} = \frac{1}{2}$$

The numerators of the two transforms then are

$$az \sin \omega_0 = \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \frac{1}{2} \frac{1}{z}$$
 and $z^2 - az \cos \omega_0 = z^2 - \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \frac{1}{z}$

In light of these we manipulate the numerator of H(z) so that it will contain $\frac{1}{2}z$ and $z^2 - \frac{1}{2}z$

Numerator =
$$z^{2} + z = \begin{pmatrix} 2 & 1 \\ z - & z \\ 2 & z \end{pmatrix} + \begin{pmatrix} 3 \\ z & z \\ 2 & y \end{pmatrix}$$

Denominator = $z^{2} - z + 0.5 = z^{2} - 2 \begin{pmatrix} 1 \\ 2 & y \\ 2 & y \\ -z & z \end{pmatrix} + \begin{vmatrix} 1 \\ 2 & y \\ -z & z \\ -z & z \end{pmatrix}$

Thus

$$H(z) = \frac{z^{2} + z}{z^{2} - z + 0.5} = \frac{\begin{pmatrix} z & 1 \\ z & -2z \end{pmatrix} + \begin{pmatrix} 3 \\ 2z \end{pmatrix}}{z^{2} - z + 0.5} = \frac{z^{2} - z}{z^{2} - z + 0.5} + 3\frac{z^{2}}{z^{2} - z + 0.5}$$

We have, in effect, arranged H(z) as

$$H(z) = \frac{z^{2} - 1}{z - 2} \frac{1}{\sqrt{z}} \frac{1}{\sqrt{z}} \frac{1}{\sqrt{z}} \frac{1}{\sqrt{z}} \frac{1}{\sqrt{z}} + 3 \frac{1}{z^{2} - 2} \frac{1}{\sqrt{z}} \frac{1}{\sqrt$$

Therefore,

$$h(n) = \mathbf{z} \quad \{H(z)\} = \mathbf{z} \quad \{\frac{2}{z} - \frac{1}{2}z \\ \frac{1}{z} - \frac{1}{2}z \\ \frac{2}{z} - \frac{1}{2}z \\ \frac{1}{z} - \frac{1}{2}z \\ \frac{2}{z} \\ \frac{1}{z} - \frac{1}{2}z \\ \frac{1}{z} - \frac{1}{2}z \\ \frac{1}{z} \\ \frac{1}{z} - \frac{1}{z}z \\ \frac{1}{z} \\ \frac{1}{z$$

$$\binom{1}{\frac{1}{\sqrt{2}}}^{n} \cos(n\pi/4) u(n) + 3\left(\frac{1}{\sqrt{2}}\right)^{n} \sin(n\pi/4) u(n)$$

In MATLAB (Pole-zero plot) It is convenient to specify the transfer function as a ratio of polynomials in z^{-1}

$$H(z) = \frac{1+z^{-1}}{1-z^{-1}+0.5z^{-2}}$$

The numerator coefficients, from left to right, are $\{b_i, i = 0 \text{ to } M\}$, specifying the vector $b = [b_0, b_1] = [1, 1]$. Similarly, the denominator coefficients are $\{a_i, i = 0 \text{ to } N\}$ from left to right, (with $a_0 = 1$) specifying the vector $a = [1, a_1, a_2] = [1, -1, 0.5]$.

%Pole-zero plot b = [1, 1]; a = [1, -1, 0.5]; zplane (b, a)



In MATLAB (Partial fractions) The partial fractions may be computed by using the *residuez* function as below. Note that H(z) is arranged as a ratio of polynomials in negative powers of z. $H(z) = \frac{1+z^{-1}}{1-z^{-1}+0.5z^{-2}} = K + \frac{R_1}{1-p_2z^{-1}} + \frac{R_2}{1-p_2z^{-1}}$

We define the coefficient vectors b = [1, 1] and a = [1, -1, 0.5]; $\mathbf{R} = [\mathbf{R}_1, \mathbf{R}_2]$ represents the residues (partial fraction constants), $p = [p_1, p_2]$ the poles and *K* a constant.

%Partial fractions b = [1, 1], a = [1, -1, 0.5], [R, p, K] = residuez (b, a)

The MATLAB results returned are

 $R = 0.5 - 1.5i \\ 0.5 + 1.5i \\ p = 0.5 + 0.5i \\ 0.5 - 0.5i \\ K = []$

Therefore,

$$H(z) = \frac{1+z^{-1}}{1-z^{-1}+0.5z^{-2}} = 0 + \frac{0.5-j1.5}{1-(0.5+j0.5)z^{-1}} + \frac{0.5+j1.5}{1-(0.5-j0.5)z^{-1}}$$

Example 1.5.3 Find the inverse transform of $X(z) = \frac{z}{3z^2 - 4z + 1}$, where the ROC is (a) |z| > 1, (b) |z| < (1/3), (c) (1/3) < |z| < 1.

Solution The three possible regions of convergence are shown below. The example shows that the inverse transform, x(n), is unique only when the ROC is specified.



(a) For ROC =
$$|z| > 1$$

$$\frac{X(z)}{z} = \frac{1}{3\left(\frac{z}{z} - \frac{1}{3} + \frac{1}{3}\right)^2} = \frac{1}{-\frac{1}{z} - \frac{1}{2} - \frac{1}{3} + \frac{1}{3}\right)^2} = \frac{1}{-\frac{1}{z} - \frac{1}{2} - \frac{1}{3} + \frac{1}{3}} = \frac{1}{-\frac{1}{z} - \frac{1}{2} - \frac{1}{3} + \frac{1}{3}} = \frac{1}{-\frac{1}{z} - \frac{1}{2} - \frac{1}{3} + \frac{1}{3}} = \frac{1}{-\frac{1}{z} - \frac{1}{z} - \frac{1}{3}} = \frac{1}{2}$$

$$= \frac{1}{3(z - 1)(z - (1/3))} = \frac{A}{z - 1} + \frac{B}{z - (1/3)}$$

$$A = \frac{1}{3(z - (1/3))} = \frac{1/2}{z - 1}, \text{ and } B = \frac{1}{3(z - 1)} = -\frac{1}{z} = \frac{1/2}{z - 1/3}$$

$$\frac{X(z)}{z} = \frac{(1/2)z}{z - 1} + \frac{(-1/2)z}{z - (1/3)}$$
The inverse is
$$x(n) = \mathbf{z}^{-1} \{X(z)\} = \mathbf{z}^{-1} \left\{ \frac{(1/2)z}{z - 1} + \frac{(-1/2)z}{z - (1/3)} \right\} = \mathbf{z}^{-1} \left\{ \frac{(1/2)z}{z - 1} \right\} + \mathbf{z}^{-1} \left\{ \frac{(-1/2)z}{z - (1/3)} \right\}$$

The ROC is outside the largest pole signifying a right-sided sequence for each pole. The inverse becomes

(b) For ROC $\equiv |z| < (1/3)$. The partial fraction expansion does not change. Since the ROC is inward of the smallest pole, x(n) consists of two negative-time sequences.

$$x(n) = \mathbf{z}^{-1} \{ X(z) \} = \mathbf{z}^{-1} \left\{ \begin{array}{c} (1/2)z \\ \hline z - 1 \\ \end{array} \right\} + \mathbf{z}^{-1} \left\{ \begin{array}{c} (-1/2)z \\ \hline \Box \\ \hline \Xi \\ \hline z \\ \end{array} \right\} \\ = 2^{(-1)} u(-n-1) - \frac{-1}{2} \\ = (1/2) \left(-1 + (1/3)^n \right) u(-n-1) \end{array} \right\}$$

(c) For ROC $\equiv (1/3) < |z| < 1$. The partial fraction expansion stays the same. The pole at z = 1 corresponds to a negative-time sequence (left-sided sequence) while the pole at z = 1/3 gives a positive-time sequence (right-sided sequence).

The overall result is a two-sided sequence.

Example 1.5.4 (Sometimes there is no *z* in the numerator to factor out, but we still can divide X(z) by *z* as in this example.) Find x(n) for $X(z) = \frac{z+1}{3z^2 - 4z + 1}$ where the ROC is |z| > 1. Solution

$$\frac{X(z)}{z} = \frac{z+1}{z(3z^2-4z+1)} = \frac{z+1}{3z(z-1)(z-(1/3))} = \frac{A}{z} + \frac{B}{z-1} + \frac{C}{z-(1/3)}$$

$$A = \frac{z+1}{3(z-1)(z-(1/3))} \Big|_{z=0} = \frac{1}{3(-1)(-1/3)} = 1$$

$$B = \frac{z+1}{3z(z-(1/3))} \Big|_{z=1} = \dots = 1$$

$$C = \frac{z+1}{3z(z-1)} \Big|_{z=1/3} = \dots = -2$$

$$\frac{X(z)}{z} = \frac{1}{z} + \frac{1}{z-1} - \frac{2}{z-(1/3)}$$

$$X(z) = 1 + \frac{z}{z-1} - \frac{2}{z-(1/3)}$$

$$x(n) = \delta(n) + u(n) - 2 \begin{pmatrix} 1 \\ -1 \\ -1 \end{pmatrix} \Big|_{z=1/3} = 0$$

Example 4.5.5 [2002] Find the inverse z-transform of $X(z) = \frac{(z-1)^2}{z^2 - 0.1 z - 0.56}$.

Solution The roots of the quadratic in the denominator are given by

$$z_{1, z_{2}} = \frac{-b \pm \sqrt{b^{2} - 4ac}}{2a} = \frac{-(-0.1) \pm \sqrt{(-0.1)^{2} - 4(1)(-0.56)}}{2(1)}$$

$$= \frac{0.1 \pm \sqrt{0.01 + 2.24}}{\frac{2}{2}} = \frac{0.1 \pm \sqrt{2.25}}{2} = \frac{0.1 \pm 1.5}{2} = 0.8 \text{ and } -0.7$$

$$X(z) = \frac{(z-1)}{(z-0.8)(z+0.7)}$$

$$\frac{X(z)}{z} = \frac{(z-1)^{2}}{z(z-0.8)(z+0.7)} = \frac{A}{z} + \frac{B}{z-0.8} + \frac{C}{z+0.7}$$

$$A = \frac{(z-1)^{2}}{(z-0.8)(z+0.7)} = \dots = -1/0.56 = -1.79$$

$$B = \frac{(z-1)^{2}}{z(z+0.7)}\Big|_{z=0.8} = \dots = 1/30$$

$$C = \frac{(z-1)^{2}}{z(z-0.8)}\Big|_{z=-0.7} = \dots = 2.75$$

$$\frac{X(z)}{z} = \frac{-1.79}{z} + \frac{(1/30)}{z-0.8} + \frac{2.75}{z+0.7}$$

$$X(z) = -1.79 + \frac{(1/30)z}{z-0.8} + \frac{2.75z}{z+0.7}$$

$$x(n) = \mathbf{z}^{-1}\{X(z)\} = 1.79 \ \delta(n)^{\frac{1}{+}}(0.8)^{n}u(n) + 2.75 \ (-0.7)^{n}u(n)$$

Example 1.5.6 Find the inverse *z*-transform of $X(z) = \frac{1}{[z - (1/2)][z - (1/4)]}, |z| > \frac{1}{2}$.

Solution

$$\frac{X(z)}{z} = \frac{1}{z[z-(1/2)][z-(1/4)]} = \frac{A}{z} + \frac{B}{z-(1/2)} + \frac{C}{z-(1/4)}$$

$$A = \frac{1}{(z-(1/2))(z-(1/4))} \Big|_{z=0} = \dots = 8$$

$$B = \frac{1}{z(z-(1/4))} \Big|_{z=1/2} = \dots = 8$$

$$C = \frac{1}{z(z-(1/4))} \Big|_{z=1/4} = \dots = -16$$

$$\frac{X(z)}{z} = \frac{8}{z} + \frac{8}{z-(1/2)} - \frac{16}{z-(1/4)}$$

$$X(z) = 8 + \frac{8z}{z-(1/2)} - \frac{16z}{z-(1/4)}$$

$$x(n) = 8 \,\delta(n) + 8 \,\left(\lfloor \frac{1}{2} \rfloor^n \right| \,u(n) - 16 \,\lfloor \frac{1}{2} \rfloor^n \right| \,u(n)$$

Example 1.5.7 Find the inverse of $X(z) = \frac{z^4 + z^2}{[z - (1/2)][z - (1/4)]}$ for ROC $\frac{1}{2} < \frac{1}{2} < \frac{1}{2$

There is a pole at $z = \infty$. The numerator degree is 3 and is greater than the denominator degree. By long division we reduce the numerator degree by 2 so that the resulting numerator degree is less than that of the denominator degree. $3\left(\frac{23}{16}z - \frac{3}{32}\right)$.

$$\begin{array}{c} X(z) = 2 & z + z \\ -z & z - (3/4)z + (1/8) \end{array} = z + \underbrace{3+}_{4} \left| \left| \begin{array}{c} 23/10 \right| 2 - (3/32) \\ 2 \\ z - (3/4)z + (1/8) \end{array} \right| \\ \left| \begin{array}{c} z \\ z - (3/4)z + (1/8) \end{array} \right|$$

(Note that in the long division leading to the above result the numerator and denominator polynomials are arranged in the order of decreasing powers of z. There are three other ways (all of them wrong) of arranging the two polynomials for the long division.)

The proper fraction part can now be expanded into partial fractions:

$$\frac{(23/16)z - (3/32)}{z^2 - (3/4)z + (1/8)} = \frac{A}{(z - (1/2))} + \frac{B}{z - (1/4)}$$

$$A = \frac{(23/16)z - (3/32)}{z - (1/4)} = 5/2$$

$$B = \frac{(23/16)z - (3/32)}{z - 3(1/2)} = -17/16$$

$$\frac{X(z) = z + + (5/2)}{z - 4} + \frac{(-17/16)}{(z - (1/2))} + \frac{(-17/16)}{z - (1/4)}$$

$$X(z) = z^2 + \frac{3z}{4} + \frac{(5/2)z}{\delta(n+1) + 1} + \frac{(-17/16)z}{1z - (1/4)}$$

$$x(n) = \delta(n+2) + \frac{\delta(n+1) + 1}{\delta(n+1) + 1} + \frac{17}{16} + \frac{17}{16} + \frac{11}{16} + \frac{17}{16} + \frac{17}{16} + \frac{11}{16} + \frac{17}{16} + \frac{11}{16} + \frac{17}{16} + \frac{11}{16} + \frac{17}{16} + \frac{11}{16} +$$

The answer has values for n = -1 and n = -2 due to the pole at $z = \infty$. The resulting x(n) is not a causal sequence.

In MATLAB (Partial fractions) The transform X(z) represents a noncausal sequence.

$$X(z) = \frac{z^{4} + z^{2}}{z^{2} - (3/4)z + (1/8)} = \frac{1 + z^{-2}}{z^{-2} - (3/4)z^{-3} + (1/8)z^{-4}}$$

Partial fractions cannot be computed by using the *residuez* function directly on X(z) since $a_0 = 0$. However, $3 \sqrt{(23/16)z^2 - (3/32)z}$

$$X(z) = z^{2} + \underbrace{z}_{4} + \begin{bmatrix} (23/10)z & -(3/32)z \\ (23/10)z & -(3/32)z \\ \vdots & \vdots \\ z^{2} - (3/4)z + (1/8) \end{bmatrix} = z^{2} + \underbrace{z}_{4} + X(z) \\ \vdots & \vdots \\ 1 \end{bmatrix}$$

where

$$X_{I}(z) = \left(\begin{vmatrix} (23/16)z^{2} - (3/32)z \\ -z^{2} - (3/4)z + (1/8) \end{vmatrix} \right) \mid = (23/16) - (3/32)z^{-1} \\ -1 - (3/4)z^{-1} + (1/8)z^{-2} \end{vmatrix}$$

On which we may use the *residuez* function.

Example 4.5.8 Assuming that $H(z) = \frac{z+4}{z-5}$ is a causal system function, prove the following independently of each other

(a)
$$h(n) = -(4/5) \ \delta(n) + (9/5) \ 5^n u(n)$$

(b) $h(n) = 5^n u(n) + (4) \ 5^{n-1} u(n-1)$
(a) $h(n) = \delta(n) + (9) \ 5^{n-1} u(n-1)$
Solution (a) $\frac{H(z)}{z} = \frac{(-4/5)}{z} + \frac{(9/5)}{z-5}$, (b) $\frac{H(z) = z}{z-5} + \frac{4}{z-5}$, (c) By long division $H(z) = 1 + \frac{9}{z-5}$.

Example 4.5.9 Partial fractions can be obtained with the *z*-transform, say H(z), expressed as a ratio of polynomials in negative powers of *z*. This amounts to expanding H(z)/z into partial fractions. Here is an example:

$$H(z) = \frac{8z^{3} - 4z^{2} + 11z - 2}{(z - (1/4))(z^{2} - z + (1/2))}$$

This example is from "Parallel realization of IIR filters", towards the end of this Unit where we obtain

$$\frac{H(z)}{z} = \frac{16}{z} + \frac{8}{(z-(1/4))} + \frac{(-16)z+20}{(z^2-z+(1/2))}$$

However, we may also proceed with negative powers of z as below (we may view $z^{-1} = p$ as a new variable):

$$H(z) = \frac{8z^{3} - 4z^{2} + 11z - 2}{(z - (1/4))(z^{2} - z + (1/2))} = \frac{8 - 4z^{-1} + 11z^{-2} - 2z^{-3}}{(1 - 0.25z^{-1})(1 - z^{-1} + 0.5z^{-2})}$$
$$= \frac{8 - 4z^{-1} + 11z^{-2} - 2z^{-3}}{1 - 1.25z^{-1} + 0.75z^{-2} - 0.125z^{-3}}$$

By long division we reduce the degree of the numerator by 1 and then expand the proper fraction part into partial fractions:

$$\begin{array}{rcl} \text{Long Division} \\ \text{Denominator} \rightarrow & -0.125z^{-3} + 0.75z^{-2} - 1.25z^{-1} + 1 & \hline 16 & \leftarrow \text{Quotient} \\ \hline & -2z^{-3} + 11z^{-2} - 4z^{-1} + 8 \\ \hline & -2z^{-3} + 12z^{-2} - 20z^{-1} + 16 \\ \hline & -z^{-2} & + 16z^{-1} - 8 \end{array} \leftarrow \text{Remainder} \end{array}$$

D · · ·

$$H(z) = 16 + \frac{-8 + 16z^{-1} - z^{-2}}{1 - 1.25z^{-1} + 0.75z^{-2} - 0.125z^{-3}}$$

Let

$$\frac{-8+16z^{-1}-z^{-2}}{1-1.25z^{-1}+0.75z^{-2}-0.125z^{-3}} = \frac{A}{(1-0.25z^{-1})} + \frac{Bz^{-1}+C}{(1-z^{-1}+0.5z^{-2})}$$

Comparing coefficients of like powers of z in the numerators on both sides

z:
$$A + C = -8$$

 z^{1} : $-A + B - 0.25C = 16$
 z^{3} : $0.5A - 0.25B = -1$
which give $A = 8, B = 20$, and $C = -16$, and
 $H(z) = 16 + \frac{8}{(-0.25z^{-1})} + \frac{-16 + 20z^{-1}}{(1 - z^{-1} + 0.5z^{-2})}$

In MATLAB The partial fractions may be computed by using the residuez function:

$$H(z) = \frac{R}{1 - 1.25z^{-1} + 0.75z^{-2} - 0.125z^{-3}} = K + \frac{R}{1 - p_{z}^{-1}} + \frac{R}{1 - p_{z}^{2}z^{-1}} +$$

We define the coefficient vectors b = [8, -4, 11, -2] and a = [1, -1.25, 0.75, -0.125]; *R* represents the residues (partial fraction constants), *p* the poles and *K* a constant. Note that in the numerator z^{-3} means M = 3, and in the denominator z^{-3} means N = 3; since *M* is not less than *N* this is not a proper rational function, so that *K* will have a nonzero element(s).

%Partial fractions b = [8, -4, 11, -2], a = [1, -1.25, 0.75, -0.125], [R, p, K] = residuez (b, a)

The MATLAB results returned are

$$R = -8 - 12i - 8 + 12i 8$$

$$p = 0.5 + 0.5i - 0.5i - 0.25$$

$$K = -16$$

Therefore,

$$H(z) = 16 + \frac{-8 - j12}{1 - (0.5 + j0.5)z^{-1}} + \frac{-8 + j12}{1 - (0.5 - j0.5)z^{-1}} + \frac{8}{1 - 0.25z^{-1}}$$

Inverse *z***-transform when there are repeated roots** With repeated roots, that is, a *k*-th order pole at z = a we have X(z) in the form

$$X(z) = \frac{z}{(z-a)^k}, \qquad \text{ROC} \quad |z| > |a|$$

The table below gives the inverse *z*-transforms for several values of k and for the general case of arbitrary k.

Repeated Roots	
X(z)	$x(n) = \mathbf{z}^{-1}[X(z)]$ for ROC $ z > a $
Z	$a^n u(n)$
(z-a)	
Z	$\frac{n}{2}a^{n-1}u(n)$
$(z-a)^2$	1!
Z	$n(n-1)a^{n-2}u(n)$
$(z-a)^3$	2!
Z	$\frac{n(n-1)(n-2)}{2}a^{n-3}u(n)$
$(z-a)^4$	3!
	$n(n-1)(n-2)(n-k-2)_{n-k-1}$
$(z-a)^k$	(k-1)! a $u(n)$

General procedure for partial fraction expansion Since X(z)/z must be rational, it takes the form

$$\frac{X(z)}{z} = \frac{\beta z^{K} + \beta z^{K-1} + \dots + \beta z + \beta}{\alpha_{L} z^{L} + \alpha_{L-1} z^{L-1} + \dots + \alpha_{1} z + \alpha_{0}}$$

If K < L then no adjustment is needed. The partial fraction expansion is straightforward. If $K \ge L$ then divide until the remainder polynomial in z has a degree of L-1 or less:

$$\frac{X(z) = (c_{K-L}z^{K-L} + \dots + c z + c) + \frac{d_{L-1}z^{L-1} + \dots + d_{1}z + d_{1}}{\alpha_{L}z^{L} + \alpha_{L-1}z^{L-1} + \dots + \alpha_{1}z + \alpha_{0}}$$

The first part of the above expression, $(c_{K-L}z^{K-L} + ... + c_1z + c_0)$, will eventually contribute δ functions to the output sequence some of which are time-advanced so that the resulting x(n) will be noncausal. The second part – the proper fraction – is expanded into partial fractions. Assume we have *one repeated pole* of order *m*, call it z_1 , and that all the rest are distinct, call them z_{m+1} , z_{m+2} , ..., z_L . Then let

$$\Psi(z) = \frac{\begin{array}{c} d \\ z^{L-1} + \dots + d \\ z + d \\ \hline \alpha \\ z^{L-1} + \alpha \\ z^{L-1} \\ A \\ \hline \end{array} + \frac{\begin{array}{c} 1 \\ \alpha \\ z + \alpha \\ z$$

The coefficients A_j (*m* of them) and B_j ($\overline{L-m}$ of them) are found as follows:

$$A_{j} = \frac{1}{(m-j)!} \left\{ \frac{d^{m-j}}{dz^{m-j}} \left[(z-z)^{m} \Psi(z) \right] \right\}_{z=z_{1}}, \quad j = 1, 2, ..., m$$

$$B_j = [(z - z_j) \Psi(z)]\Big|_{z = z_j}, \qquad j = m + 1, m + 2, ..., L$$

In the resulting x(n) the contribution of the A_j terms is a number of exponentials multiplied by n, (n-1), (n-2), etc., and the contribution of the B_j terms is a number of complex exponentials.

Example 1.5.10 Find the inverse of

$$H(z) = \frac{z}{z^3 + 2z^2 + 1.25z + 0.25} = \frac{z}{1 + 2z^{-1} + 1.25z^{-2} + 0.25z^{-3}}, |z| > 1$$

Solution This transform is a proper rational function. We shall use this example to give a summary of the three styles of obtaining partial fractions: (1) Expanding H(z) directly, (2) Expanding H(z)/z and (3) Expanding $H(z^{-1})$ as in MATLAB (also Mitra).

When the poles of H(z) are distinct the partial fraction coefficients returned by the MATLAB function *residuez* are the same as in expanding H(z)/z. However, when there are repeated poles it makes a difference in the coefficients as well as in the final analytical forms of the inverse transforms in these two methods. In addition, directly expanding H(z) results in an analytical form that is still different from the other two. In any event the three inverse transforms are the same as far as the actual sequence values are concerned.

(1) **Expanding** *H*(*z*) We have

$$H(z) = \frac{z}{\frac{3}{z+2z}+1.25z+0.25} = \frac{z}{(z+1)(z+0.5)} = \frac{A}{(z+1)} + \frac{B}{\frac{2}{(z+0.5)^2}} + \frac{B}{(z+0.5)^2} + \frac{B}$$

-2

There is a pole at z = -1 and a repeated pole at z = -0.5. The coefficients *A*, *B*₂ and *B*₁ are given by

$$A = \frac{z}{(z+0.5)^2} \Big|_{z=-1} = \frac{-1}{(-1+0.5)^2} = -4$$

$$B_2 = \frac{z}{(z+1)} \Big|_{z=-0.5} = \frac{-0.5}{(-0.5+1)} = -1$$

$$B_1 = \{-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1 - |-1$$

Thus

$$H(z) = \frac{(-4)}{(z+1)} + \frac{(-1)}{(z+0.5)^2} + \frac{4}{(z+0.5)}$$

Taking the inverse *z*-transform, $h(n) = -4 \mathbf{z}^{-1} \int_{-\infty}^{\infty} \mathbf{z}^{-1} d\mathbf{z}^{-1}$

$$\begin{aligned} h(n) &= -4 \ \mathbf{\bar{3}}^{-1} \left\{ \begin{array}{c} 1\\ \Box \\ z + 1 \end{array} \right\} - 1 \ \mathbf{\bar{3}}^{-1} \left\{ \begin{array}{c} 1\\ \Box \\ z + 0.5 \end{array} \right\} + 4 \ \mathbf{\bar{3}}^{-1} \left\{ \begin{array}{c} 1\\ \Box \\ z + 0.5 \end{array} \right\} \\ &= -4 \ \mathbf{\bar{3}}^{-1} \left\{ \begin{array}{c} z^{-1} & z\\ z + 1 \end{array} \right\} - 1 \ \mathbf{\bar{3}}^{-1} \left\{ \begin{array}{c} z^{-1} & z\\ z^{-1} & \overline{z} \\ (z + 0.5)^2 \end{array} \right\} + 4 \ \mathbf{\bar{3}}^{-1} \left\{ \begin{array}{c} z^{-1} & z\\ \overline{z + 0.5} \end{array} \right\} \\ &= -4 \left((-1)^n \ u(n) \right)_{n \to n-1} - 1 \left| \begin{array}{c} (-0.5)^{n-1} \ u(n) \\ 1! \end{array} \right|_{n \to n-1} + 4 \left((-0.5)^n \ u(n) \right)_{n \to n-1} \\ &= 4 \ (-1)^n \ u(n-1) - 4 \ (n-1)(-0.5)^n \ u(n-1) - 8 \ (-0.5)^n \ u(n-1) \end{aligned}$$

(2) Expanding H(z)/z We have

$$\frac{H(z)}{z} = \frac{1}{z + 2z + 1.25z + 0.25} = \frac{z}{(z+1)(z+0.5)^2} = \frac{C}{(z+1)} + \frac{D_2}{(z+0.5)^2} + \frac{D_$$

The coefficients C, D_2 and D_1 are given by

$$C = \frac{1}{(z+0.5)^2} \bigg|_{z=-1} = \frac{1}{(-1+0.5)^2} = 4$$

$$D_2 = \frac{1}{(z+1)} \bigg|_{z=-0.5} = \frac{1}{(-0.5+1)} = 2$$

$$D_1 = \left\{ -\left| \frac{1}{(z+1)} \right|_{z=-0.5} = \frac{(z+1).0 - 1(1)}{(z+1)^2} \right|_{z=-0.5} = \frac{-1}{(-0.5+1)^2} = -4$$

Thus

$$\frac{H(z)}{z} = \frac{4}{(z+1)} + \frac{2}{(z+0.5)^2} + \frac{(-4)}{(z+0.5)}$$
$$H(z) = 4\frac{z}{(z+1)} + 2\frac{z}{(z+0.5)^2} - 4\frac{z}{(z+0.5)}$$

Taking the inverse *z*-transform,

$$h(n) = 4 \mathbf{z}^{-1} \left\{ \begin{array}{c} z \\ \Box \\ z + 1 \end{array} \right\} + 2 \mathbf{z}^{-1} \left\{ \begin{array}{c} z \\ \Box \\ (z + 0.5)^2 \end{array} \right\} - 4 \mathbf{z}^{-1} \left\{ \begin{array}{c} z \\ \Box \\ z + 0.5 \end{array} \right\}$$
$$= 4 (-1)^n u(n) + 2 \frac{n}{1!} (-0.5)^{n-1} u(n) - 4 (-0.5)^n u(n)$$
$$= 4 (-1)^n u(n) - 4 n(-0.5)^n u(n) - 4 (-0.5)^n u(n)$$

(3) Expanding $H(z^{-1})$ as in MATLAB We start with H(z) expressed as a ratio of polynomials in negative powers of z. However, for the sake of continuity we have

$$H(z) = \frac{z}{+1.25z + 0.25} z^{3} + 2z^{2} = \frac{z}{1 + 2z^{-1} + 1.25z^{-2} + 0.25z^{-3}}$$
$$H(z) = \frac{z}{(z + 1z + 0.5)} = \frac{z}{(1 + 1z^{-1})(1 + 0.5z^{-1})^{2}}$$
$$= \frac{E}{(1 + 1z^{-1})} + \frac{F_{1}}{(1 + 0.5z^{-1})} + \frac{F_{2}}{(1 + 0.5z^{-1})^{2}}$$

(We have ordered the coefficients in the order in which MATLAB displays them). We can define $z^{-1} = v$ so that z = -1 corresponds to v = -1 and z = -0.5 to v = -2. The transform now appears as V^2

$$H(v) = \frac{1}{(1+1v)(1+0.5v)^2} = \frac{1}{(1+1v)} + \frac{1}{(1+0.5v)} + \frac{2}{(1+0.5v)^2}$$
ts E E and E are given by

The coefficients E, F_2 and F_1 are given by

$$E = \frac{z^{-2}}{(1+0.5z^{-1})^2} \bigg|_{z^{-1}=-1} = \frac{1^2}{(1+0.5(-1))^2} = 4$$
$$F_2 = \frac{z^{-2}}{(1+1z^{-1})} \bigg|_{z^{-1}=-2} = \frac{(-2)^2}{(1+1.(-2))} = -4$$

$$F_{I} = \left\{ \frac{d}{dv} \left(\frac{v^{2}}{1+1v} \right) \right\} \bigg|_{v=-2} = \left\{ \frac{(1+v) \cdot 2v - v^{2} \cdot (1)}{(1+v)^{2}} \right\} \bigg|_{v=-2} = \frac{(1-2)2(-2) - (-2)^{2}}{(1-2)^{2}} = 0$$

Therefore,

$$H(z) = \frac{4}{(1+1z^{-1})^{+}} + \frac{0}{(1+0.5z^{-1})^{+}} + \frac{(-4)}{(1+0.5z^{-1})^{2}}$$
$$= 4\frac{z}{z+1} - 4\frac{z^{2}}{(z+0.5)^{2}} = 4\frac{z}{z+1} - 4z\frac{z}{(z+0.5)^{2}}$$

Taking the inverse *z*-transform,

$$h(n) = 4 \ \mathbf{z}^{-1} \left\{ \begin{array}{c} z \\ \Box \\ z+1 \end{array} \right\} - 4 \ \mathbf{z}^{-1} \left\{ \begin{array}{c} z \\ \Xi \\ z+1 \end{array} \right\} \left[\begin{array}{c} z \\ (z+0.5)^2 \end{array} \right] \\ = 4 \ (-1)^n u(n) - 4 \left\{ \begin{array}{c} -(-0.5)^{n-1} u(n) \\ 1! \end{array} \right\} \right]_{n \to n+1} \\ = 4 \ (-1)^n u(n) - 4 \ (-0.5)^{-1} \left\{ \begin{array}{c} n \\ -(-0.5)^n u(n) \\ 1! \end{array} \right\} \right]_{n \to n+1} \\ = 4 \ (-1)^n u(n) + 8 \ (n+1)(-0.5)^{n+1} u(n+1) \end{array}$$

In MATLAB This particular set of partial fractions may be computed by using the *residuez* function:

$$H(z) = \frac{z^{-2}}{1 + 2z^{-1} + 1.25z^{-2} + 0.25z^{-3}} = K + \frac{E}{1 - p \cdot z^{-1}} + \frac{F_1}{1 - p \cdot z^{-1}} + \frac{F_2}{\left(1 - p_2 z\right)^2}$$

We define the coefficient vectors b = [0, 0, 1] and a = [1, 2, 1.25, 0.25]; $R = [E, F_1, F_2]$ represents the residues (keyed to the above partial fraction coefficients), *p* the poles and *K* a constant.

%Partial fractions b = [0, 0, 1], a = [1, 2, 1.25, 0.25], [R, p, K] = residuez (b, a)

The MATLAB results returned are

 $R = 4 \\ -0 + 0i \\ -4 - 0i \\ p = -1 \\ -0.5 + 0i \\ -0.5 - 0i \\ K = []$

Therefore,

$$H(z) = \frac{z^{-2}}{1 + 2z^{-1} + 1.25z^{-2} + 0.25z^{-3}} = 0 + \frac{4}{1 - (-1)z^{-1}} + \frac{0}{1 - (-0.5)z^{-1}} + \frac{(-4)}{(1 - (-0.5)z^{-1})^2}$$
$$\underline{-\frac{4}{1+1 z^{-1}}} + \frac{(-4)}{(1+0.5 z^{-1})^2}$$

which agrees with the hand-calculated results.

Example 1.5.11 Find the inverse of

$$\begin{aligned} \mathbf{X} (z) &= \frac{z^3 - z^2 + z - 1}{16} = \frac{z^3 - z^2 + z - 1}{16}, \text{ for } |z| > \frac{1}{2} \\ \frac{z}{z^3 - (5/4)z^2 + (1/2)z - (1/16)}, \frac{z}{z^3 - (5/4)z^2 + (1/2)z - (1/16)}, \frac{z}{z^3 - (5/4)z^2 + (1/2)z - (1/16)}, \frac{z}{z^3 - (2-z)z + z - (1/16)}, \frac{z}{z^3 - (1/2)^2} + \frac{z}{z^3 - (1/2)^2} + \frac{z}{z^3 - (1/2)^2} + \frac{z}{z^3 - (1/2)^2} + \frac{z}{z^3 - (1/2)^2}, \frac{z}{z^3 - (1/2)^2} = \frac{(1/4)^3 - (1/4)^2 + (1/4) - (1/16)}{(1/4)[(1/4) - (1/2)]} = 1 \\ C &= \frac{z^3 - z^2 + z - (1/16)}{z[z - (1/2)]^2} \Big|_{z = 1/4} = \frac{(1/2)^3 - (1/2)^2 + (1/2) - (1/16)}{(1/2)[(1/2) - (1/4)]} = 9 \\ B_2 &= \frac{z^3 - z^2 + z - (1/16)}{z[z - (1/4)]} \Big|_{z = 1/2} = \frac{(1/2)^3 - (1/2)^2 + (1/2) - (1/16)}{(1/2)[(1/2) - (1/4)]} = 5/2 \\ B_1 &= \left\{ \frac{d}{dz} \left[\frac{z^3 - z^2 + z - (1/16)}{z[z - (1/4)]} \right]_{z = 1/2} \right\} \right|_{z = 1/2} = \left\{ \frac{d}{dz} \left[\frac{z^3 - z^2 + z - (1/16)}{(1/2)[(1/2) - (1/4)]} \right]_{z = 1/2} \right\} \\ &= \left\{ \frac{d}{dz} \left[\frac{(z^3 - (z - (z/4))](3z^2 - 2z + 1) - [z^3 - z^2 + z - (1/16)][2z - (1/4)]}{(z^2 - (z/4))^2} \right]_{z = 1/2} \right\} \right\} \right\} \\ &= \frac{z}{z} = 0 \end{aligned}$$

$$\frac{X(z)}{z} = \frac{-9}{z} + \frac{(5/2)}{(z - (1/2))^2} + \frac{(-9)}{(z - (1/2))} + \frac{9}{z - (1/4)}$$
$$X(z) = 1 + \frac{5}{2} \frac{z}{(z - (1/2))^2} - 9 \frac{z}{(z - (1/2))} + 9 \frac{z}{z - (1/4)}$$

Taking the inverse z-transform,

$$\begin{aligned} x(n) &= \mathbf{\mathfrak{Z}}^{-1} \{ X(z) \} \\ &= \mathbf{\mathfrak{Z}}^{-1} \{ 1 \} + (5/2) \mathbf{\mathfrak{Z}}^{-1} \{ \left[\begin{array}{c} z \\ z - (1/2)^2 \end{array} \right]^2 - 9 \mathbf{\mathfrak{Z}}^{-1} \left[\begin{array}{c} z \\ (z - (1/2))^2 \end{array} \right] + 9 \mathbf{\mathfrak{Z}}^{-1} \{ \underline{z} \\ (z - (1/2)) \end{bmatrix} \\ &= \delta(n) + \begin{array}{c} n \\ n \\ - n \\ - u(n) - 9 \\ - u(n) - 9 \\ - u(n) + 9 \\ - u(n) \\ - 2 \\ -$$

2

Other possibilities If we choose to expand X(z), rather than X(z) / z, into partial fractions, we need to perform long division to reduce the degree of the numerator by 1 resulting in

$$X(z) = 1 + \frac{(z / 4) + (z / 2)}{z^{3} - (5 / 4)z^{2} + (1 / 2)z - (1 / 16)} = 1 + X (z)$$

where $X_1(z)$ is the proper fraction part of the above $(z^2/4) + (z/2)$

$$X_1(z) = \frac{(z^2/4) + (z/2)}{z^3 - (5/4)z^2 + (1/2)z - (1/16)}$$

Either $X_1(z)$ itself or $X_1(z) / z$ may now be expanded into partial fractions.

In MATLAB The partial fractions may be computed by using the *residuez* function:

$$X(z) = \frac{1 - z^{-1} + z^{-1} - (1/16)z^{-1}}{1 - (5/4)z^{-1} + (1/2)z^{-2} - (1/16)z^{-3}} = K + \frac{K_1}{1 - p_1 \overline{z}} + \frac{K_1}{1 - p_2 \overline{z}^{-1}} + \frac{K_1}{(1 - p_2 \overline{z}_{-1})}$$

We define the coefficient vectors b = [1, -1, 1, -1/16] and a = [1, -5/4, 1/2, -1/16]; *R* represents the residues (partial fraction coefficients), *p* the poles and *K* a constant. Note that in the numerator z^{-3} means M = 3, and in the denominator z^{-3} means N = 3; since *M* is not less than *N* this is not a proper rational function, so that *K* will have a nonzero element(s).

%Partial fractions b = [1, -1, 1, -1/16], a = [1, -5/4, 1/2, -1/16], [R, p, K] = residuez (b, a)

The MATLAB results returned are

 $R = -14.0000 \\ 5.0000 \\ 9.0000 \\ p = 0.5000 \\ 0.5000 \\ 0.2500 \\ K = 1$

Therefore,

$$X(z) = \underbrace{1 - z^{-1} + z^{-2} - (1/16)z^{-3}}_{1 - (5/4)z^{-1} + (1/2)z^{-2} - (1/16)z^{-3}} = 1 + \underbrace{9}_{1 - 0.25z^{-1}} + \underbrace{(-14)}_{1 - 0.5z^{-1}} + \underbrace{5}_{(1 - 0.5z^{-1})^2}$$

$$z^{2} + z$$

Example 1.5.12 Find the inverse of
$$X(z) = \frac{z+1}{[z-(1/2)]^3 [z-(1/4)]A} + B$$

$$\frac{\overline{X(z)}}{z} = \frac{z+1}{[z-(1/2)]^3 [z-(1/4)]} = (z-(1/2))^3 + (z-(1/2))^2 + (z-(1/2)) + B$$

$$A_3 = \frac{z+1}{[z-(1/4)]} \Big|_{z=V^2} = \frac{(1/2)+1}{[(1/2)-(1/4)]} = 6$$

$$A_2 = (3-2)! dz \Big|_{z=(1/4)} \int_{z=1/2} dz \Big|_{z=1/2} = ... = -20$$

$$A_1 = (3-1)! dz^2 \Big|_{z=(1/4)} \int_{z=1/2} dz \Big|_{z=1/2} = ... = 80$$

$$B = \frac{z+1}{[z-(1/2)]^3} \Big|_{z=1/4} = ... = -80$$

$$\frac{X(z)}{z} = \frac{6}{(z-(1/2))^3} + \frac{(-20)}{(z-(1/2))} + \frac{80}{(z-(1/2))} + \frac{(-80)}{(z-(1/4))}$$

$$X(z) = 6 \frac{z}{(z - (1/2))^3} 20 \frac{z}{(z - (1/2))^2} + 80 \frac{z}{(z - (1/2))} 80 \frac{z}{z - (1/4)}$$

$$x(n) = \mathbf{\underline{3}}^{-1} \{X(z|)\} z \\ = 6 \mathbf{\underline{3}}^{-1} \{\underbrace{(z - (1/2))^3}_{z} \} - 20 \mathbf{\underline{3}}^{-1} \{\underbrace{\underline{\Box}}_{z} \} \\ - 20 \mathbf{\underline{3}}^{-1} \{\underbrace{\underline{\Box}}_{z} \} \\ - 80 \mathbf{\underline{3}}^{-1} \\ - 80 \mathbf{\underline{3}}^{-1} \end{bmatrix}$$

$$= 6 \frac{n(n-1)}{2!} \underbrace{1}_{2}^{n-2} u(n) - 20 n ||^{n-1} u(n) + 80||^{n} u(n) - 80||^{n} u(n) \\ (\underline{1}) \\ (\underline{1}$$

The u(n) may be factored out etc.

MATLAB

$$X(z) = \frac{z^{-2} + z^{-3}}{1 - (7/4)z^{-1} + (9/8)z^{-2} - (5/16)z^{-3} + (1/32)z^{-4}}$$

%Partial fractions b = [0, 0, 1, 1], a = [1, -7/4, 9/8, -5/16, 1/32], [R, p, K] = residuez (b, a)

The MATLAB results returned are

$$R = 1.0e+002 *$$

$$1.44 + 0.0i$$

$$-0.88 - 0.0i$$

$$0.24$$

$$-0.80$$

$$P = 0.5 + 0.0i$$

$$0.5 - 0.0i$$

$$0.5 - 0.0i$$

$$0.5$$

$$0.25$$

$$K = []$$

$$X(z) = \frac{144}{(1-(1/2)z^{-1})} + \frac{(-88)}{(1-(1/2)z^{-1})^2} + \frac{24}{(1-(1/2)z^{-1})^3} + \frac{(-80)}{1-(1/4)z^{-1}}$$

Relationships among system representations

A discrete-time linear shift-invariant system can be characterized by its unit sample response, a difference equation, a system function, or a frequency response. Assume that a system is described by the linear constant coefficient difference equation

$$\sum_{k=0}^{N} a_k y(n-k) = \sum_{r=0}^{M} b_r x(n-r)$$

System function Take the *z*-transform of both sides of the above equation $\begin{bmatrix} N \\ M \end{bmatrix}$

$$\mathbf{J} \left\{ \sum_{k=0}^{N} a_{k} y(n-k) \right\} = \mathbf{J} \left\{ \sum_{r=0}^{M} b_{r} x(n-r) \right\}, \text{ or }$$

$$\sum_{k=0}^{N} a_{k} \mathbf{J} \left\{ y(n-k) \right\} = \sum_{r=0}^{M} b_{r} \mathbf{J} \left\{ x(n-r) \right\}, \text{ or }$$

$$\sum_{k=0}^{N} a_{k} z^{-k} Y(z) = \sum_{r=0}^{M} b_{r} z^{-r} X(z), \text{ or }$$

$$Y(z) \sum_{k=0}^{N} a_{k} z^{-k} = X(z) \sum_{r=0}^{M} b_{r} z^{-r}$$

$$Y(z) \sum_{k=0}^{N} b_{r} z^{-r}$$

The system function is $H(z) = \frac{Y(z)}{X(z)} = \frac{\sum_{r=0}^{T} b_r z^{-r}}{\sum_{k=0}^{T} a_k z^{-k}}$

Unit sample response If $x(n) = \delta(n)$ then $X(z) = \mathfrak{Z}[x(n)] = \mathfrak{Z}[\delta(n)] = 1$. The corresponding y(n) is the unit sample response h(n). We have

$$\frac{I'(z)}{X(z)} = H(z), \text{ or } Y(z) = H(z).X(z) = H(z).1 = H(z)$$

So, given H(z), the system function, the unit sample response is $h(n) = \mathbf{z}^{-1}[H(z)]$.

The difference equation from the H(z) The system function H(z) is first written in terms of negative powers of z and set equal to $\frac{Y(z)}{X(z)}$. Then cross-multiply and take the inverse z-transform to get the difference equation.

Frequency response of the system is the Fourier transform (DTFT) of the unit sample response h(n):

$$H(e^{j\omega}) = \sum_{n=-\infty}^{\infty} h(n)e^{-j\omega n}$$

Compare this with the system function H(z) defined as the z-transform of the unit sample response h(n)

$$H(z) = \sum_{n = -\infty}^{\infty} h(n) \ z^{-n}$$

Thus the frequency response, if it exists, can be obtained by replacing the z in H(z) by $e^{j\omega}$ as follows:

$$H(e^{j\omega}) \Delta \sum_{n=-\infty}^{\infty} h(n) e^{-j\omega n} = H(z) \Big|_{z=e^{j\omega}}$$

The system is implicitly BIBO-stable.

The above relationships for a *stable, causal system* represented by a *linear constant coefficient difference equation* are summarized in diagram below.



Example 1.6.1 Find the impulse response of y(n) = a y(n-1) + x(n). **Solution** Note that we have solved this in the time domain earlier. Taking the *z*-transform of both sides (with zero initial conditions).

$$Y(z) = a \ z^{-1} \ Y(z) + X(z), \text{ or}$$
$$H(z) = \frac{Y(z)}{X(z)} = \frac{1}{1 - az^{-1}} = \frac{z}{z - a}$$

Assume causality. Then from the table of transforms, $h(n) = \mathbf{z}^{-1}[H(z)] = a^n u(n)$.

Causality in terms of the *z***-transform,** H(z), and the ROC A causal LTI system has an impulse response h(n) = 0 for n < 0, and is, therefore, a right-sided sequence. This also implies that the ROC of H(z) is the exterior of a circle in the *z*-plane. For a causal system the power series

$$H(z) = \sum_{n=0}^{\infty} h(n) z^{-n} = h(0) + h(1) z^{-1} + h(2) z^{-2} + \dots \quad \to \text{Eq. (1)}$$

does not include any positive powers of z. Consequently, the ROC includes $z = \infty$. Therefore, we have the principle:

A discrete-time LTI system is causal if and only if the ROC of its system function is the exterior of a circle, and includes $z = \infty$.

The initial value theorem says that for a causal sequence, h(n), the initial value is given by $h(0) = \lim_{z \to \infty} H(z)$

This may be seen by setting $z \to \infty$ in Eq. (1) making all terms go to zero except the term h(0). Thus, for a causal sequence, h(n), if h(0) is finite, then, lim H(z) is finite. Consequently, with

H(z) expressed as a ratio of polynomials in z (positive powers of z), the order of the numerator polynomial cannot be greater than the order of the denominator polynomial (if it were there

would be positive powers of z in the power series of H(z), corresponding to non-zero h(n) for negative n; also $z = \infty$ would not be included in the ROC); or, equivalently, the number of finite zeros of H(z) cannot be greater than the number of finite poles.

The above discussion is summed up as follows: A discrete-time LTI system with rational system function H(z) is causal if and only if

- 1. The ROC is the exterior of a circle outside the outermost pole, and,
- 2. With H(z) expressed as a ratio of polynomials in z, (positive powers of z), the order of the numerator is not greater than the order of the denominator.

Condition 1 alone is not enough because the sequence may be right-sided but not causal. If H(z) is represented as a ratio of polynomials in z as

$$H(z) = \frac{\beta_{K} z^{K} + \beta_{K-1} z^{K-1} + \dots + \beta_{1} z + \beta_{0}}{\alpha z_{L}^{L} + \alpha z_{L-1}^{L-1} + \dots + \alpha z + \alpha} \longrightarrow \text{Eq. (2)}$$

then $L \ge K$ if the system is causal – in other words denominator degree \ge numerator degree. On the other hand, if we write H(z) as the ratio of polynomials in z^{-1} (negative powers of z) as

$$H(z) = \frac{b + b z + \dots + b}{a_0 + 1} z^{-M+1} + b_M z^{-M}}{a_0 + a z^{-1} + \dots + a_{N-1z}^{-N+1} + b_M z^{-M}} = \frac{b_0 + (b_1/z) + \dots + (b_{M-1}/z) + (b_M/z)}{a_0 + (a_1/z) + \dots + (a_{N-1}/z^{N-1}) + (a_N/z^N)}$$

then, if the system is (to be) causal, $a_0 \neq 0$. This is seen by setting $z \rightarrow \infty$, and requiring that $h(0) = (b_0/a_0)$ be finite. This is illustrated with an example where $a_0 = 0$, e.g.,

$$H(z) = \frac{1+z^{-1}+z^{-2}}{0+z^{-1}+z^{-2}} = \frac{z^{2}+z+1}{z+1}$$

which, by long division, can be seen to contain z^{1} – a positive power of z – hence non-casual.)

Note When H(z) is written as a ratio of polynomials in z (positive power of z), as in Eq. (2), we have required that $L \ge K$ for causality. These L and K are not to be confused with the N and M contained in the difference equation. Consider, for example, the system

$$y(n) + a y(n-1) = x(n) + b x(n-3)$$

where, according to the notation of the difference equation, N = 1 and M = 3. Apparently M is greater than N and this is allowable. In other words, there is no restriction on the relative values of N and M. For, the transfer function is given by

$$\frac{H(z) = Y(z)}{\overline{X(z)}} = \frac{1 + b z^{-3}}{1 + a z^{-1}} = \frac{z^{-3}(z^{3} + b)}{z^{-1}(z + a)} = \frac{z^{3} + b}{z^{2}(z + a)}$$

and it is seen that the numerator degree (K = 3) is not greater than the denominator degree (L = 3). Thus the system is causal.

As another example consider y(n) + a y(n-1) = x(n) + b x(n+1) which is non-causal because of the x(n+1) term. The transfer function is

$$\dot{H}(z) \stackrel{'}{=} \frac{Y(z)}{X(z)} = \frac{1+\frac{1}{1+az^{-1}}}{\frac{bz}{1+az^{-1}}} = \frac{z(b+z^{-1})}{1+az^{-1}} = \frac{(b+z^{-1})}{z^{-1}+az^{-2}} = \frac{(b+z^{-1})}{0+z^{-1}+az^{-2}}$$

Note that, when the numerator and denominator are expressed in terms of negative powers of z, " a_0 " = 0. On the other hand, when the numerator and denominator are expressed in terms of positive powers of z, we have

$$H(z) = \frac{Y(z)}{X(z)} = \frac{b z}{z+a}^{2}$$

with the numerator degree greater than the denominator degree.

(*Omit*) *Rational transfer function; LTI system* Given the system with the N^{th} order difference equation,

$$a_0 y(n) + a_1 y(n-l) + \dots + a_N y(n-N) = b_0 x(n) + b_1 x(n-l) + \dots + b_M x(n-M), \qquad a_0 \neq 0$$

we may write it in the more compact form

$$\sum_{k=0}^{N} a_k \ y(n-k) = \sum_{r=0}^{M} b_r \ x(n-r), \qquad a_0 \neq 0$$

(Note that some authors take the coefficient of y(n), a_0 , to be 1. In the above difference equation we may divide through by a_0 so that the coefficient of y(n) is 1).

We can find the transfer function of the system by taking the *z*-transform on both sides of the equation. We note that in finding the impulse response of a system and, consequently, in finding the transfer function, the system must be initially relaxed ("zero initial conditions"). Thus, if we assume zero initial conditions, we can use the linearity and time-shift properties to get

$$Y(z) \sum_{k=0}^{N} a_k z^{-k} = X(z) \sum_{r=0}^{M} b_r z^{-r}$$

so that

$$\frac{H(z) = Y(z)}{X(z)} = \frac{\sum_{r \ge 0}^{M} b^{r} z^{-r}}{\sum_{k=0}^{T} a_{k} z_{-k}}$$
Eq. (1)

The corresponding impulse response can be found as $h(n) = \mathbf{z}^{-1}{H(z)}$. The poles of the system transfer function are the same as the characteristic values of the corresponding difference equation. For the system to be stable, the poles must lie within the unit circle in the *z*-plane. Consequently, for a stable, causal function, the ROC includes the unit circle.

М

The system function, H(z), is a rational function:

$$H(z) = \frac{N(z)}{D(z)} = \frac{b + b z^{-1} + b z^{-2} + \dots + b z^{-M}}{a^{0^+} a^{1} z^{-1} + a^{\frac{2}{2}} + \dots + a^{\frac{M}{2}}} = \frac{\sum_{k=0}^{k} b^{k} z^{-k}}{\sum_{k=0}^{N} a z}$$

Here N(z) and D(z) stand for numerator and denominator respectively. If $a_0 \neq 0$ and $b_0 \neq 0$, we can avoid the negative powers of z by factoring out $b_0 \overline{z_1}^{-M}$ and $a_0 \overline{z_1}^{-N}$ as follows:

$$H(z) = \underbrace{b_0 z}_{z} = \underbrace{b_0 z}_{z^{-N}} \cdot \underbrace{z}_{z^{-N}}^{N} + \underbrace{(b_1/b_0) z}_{z^{-N}}_{z^{-N}} + \underbrace{(a/a) z}_{z^{-N}}^{N} + \underbrace{(a/a) z}_{z^{-N}}_{z^{-N}} + \underbrace{(a/a) z}_{z^{-N}} + \underbrace{(a/a) z}_{z^{-N}}$$

Since N(z) and D(z) are polynor fliats in z, they can be expressed in factored form as

$$H(z) = N(z) = \begin{vmatrix} b \\ -D(z) \\ -D(z) \end{vmatrix} = \begin{vmatrix} \frac{0}{-1} \\ \frac{1}{2} \\ \frac{1$$

$$= C z \xrightarrow{N-M}_{k=1}^{M} \frac{(z - z_k)}{\sum_{k=1}^{N} (z - p_k)}, \text{ where } C = (b_0/a_0)$$

Thus H(z) has *M* finite zeros at $z = z_1, z_2, ..., z_M$, and *N* finite poles at $z = p_1, p_2, ..., p_N$, and |N-M| zeros (if N > M) or poles (if N < M) at the origin z = 0. Poles and zeros may also occur at $z = \infty$. A pole exists at $z = \infty$ if $H(\infty) = \infty$, and a zero exists at $z = \infty$ if $H(\infty) = 0$. If we count the poles and zeros at z = 0 and $z = \infty$ as well as the *N* poles and *M* zeros, we find that H(z) has exactly the same number of poles and zeros.

By definition the ROC of H(z) should not (can not) contain any poles.

Proper rational function Taking $a_0 = 1$, we have $H(z) = {N(z) \atop b + b z^{-1} + b z^{-2} + \dots + b z^{-M}}$

$$\overline{D(z)} = \frac{1^{0} + a_{z}^{2} + b_{z}^{2} + b_{z}^{2} + \dots + a_{z}^{N}}{1^{0} + a_{z}^{2} + a_{z}^{2} + \dots + a_{z}^{N}}$$

This is called a **proper** rational function if $a_N \neq 0$ and M < N. This amounts to saying that the number of finite zeros is less than the number of finite poles. (*Finite* zeros and poles exclude those at z = 0). This condition is related to partial fraction expansion and has nothing to do with causality.

(End of Omit)

Example 1.6.2 Give the pole-zero plot for $H(z) = \frac{z}{z^2 - z^1 - 1} = \frac{z^{-1}}{1 - z^{-1} - z^{-2}}$

Solution The denominator has roots (poles) at

$$z_1, z_2 = \frac{-(-1) \pm \sqrt{(-1)^2 - 4(1)(-1)}}{2} = \frac{1 \pm \sqrt{1 \pm 4}}{2} = \frac{1 \pm \sqrt{5}}{2} = 1.62 \text{ and } -0.62$$

There is a zero at z = 0. Further, since the denominator degree is greater than the numerator degree by 1 it is clear that $H(\infty) = 0$, so that there is an additional zero at $z = \infty$.

In MATLAB the transfer function is specified as a ratio of polynomials in z^{-1}

$$H(z) = \frac{0+1.z}{1-1.z^{-1}-1.z^{-2}}$$

The numerator coefficients, $\{b_i, i = 0 \text{ to } M\}$ and the denominator coefficients $\{a_i, i = 0 \text{ to } N\}$ are specified as the two vectors b = [0, 1] and a = [1, -1, -1].

%Pole-zero plot b = [0, 1]; a = [1, -1, -1]; zplane (b, a)



Example 1.6.3 Give the pole-zero plot for

$$H(z) = z^{-1} + 2z^{-2} + 3z^{-3} + 4z^{-4} + 5z^{-5} + 6z^{-6} + 7z^{-7} + 8z^{-8} + 9z^{-6}$$

Solution From

$$H(z) = \frac{z^8 + 2z^7 + 3z^6 + 4z^5 + 5z^4 + 6z^3 + 7z^2 + 8z + 9}{z^9}$$

we can see that there are 9 poles at z = 0 and 8 zeros at sundry places and an additional zero at $z = \infty$ owing to the denominator degree being greater than the numerator degree by 1. For the MATLAB segment the numerator and denominator coefficients are taken from

the MATLAB segment the numerator and denominator coefficients are taken from

$$H(z) = \frac{0 + z^{-1} + 2z^{-2} + 3z^{-3} + 4z^{-4} + 5z^{-5} + 6z^{-6} + 7z^{-7} + 8z^{-8} + 9z^{-9}}{1}$$

%Pole-zero plot b = [0: 9]; a = [1, 0]; zplane (b, a)



Example 1.6.4 Give the pole-zero plot for

$$y(n) = x(n) + 0.81 x(n-1) - 0.81 x(n-2) - 0.45 y(n-2)$$

Solution The zeros are

$$z_2 = \frac{-0.81 \pm \sqrt{(0.81)^2 - 4(1)(-0.81)}}{2(1)} = 0.5819 \text{ and } -1.3919$$

The poles are given by

Z1,

 z_1 , $z_2 = \pm j0.67082$

For the MATLAB program the coefficient vectors are b = [1, 0.81, -0.81] and a = [1, 0, 0.45].

%Pole-zero plot b = [1, 0.81, -0.81]; a = [1, 0, 0.45]; zplane (b, a)



From the general form H(z) in Eq.(1) we can obtain two important special forms: (1) the all-zero system, and (2) the all-pole system. (There are, of course, trivial poles or zeros present.)

The all-zero system If $a_k = 0$ for $1 \le k \le N$, we have $H(z) = \frac{\sum b^r z^{-r}}{a_0}$. Either take $a_0 = 1$ or consider that a_0 is absorbed in the b_r coefficients, so that $b_z^M + b_z^{M-1} + ... + b$

$$H(z) = b_{0} + b_{1}z_{0} + b_{1}z_{0} + \dots + b_{M}z^{-M} = \begin{bmatrix} 0 & 1 & M \\ 0 & 1 & 2 & M \end{bmatrix}$$

In this case, H(z) contains M zeros and an M^{th} order pole at the origin z = 0. Since the system contains only *trivial poles* (at z = 0) and M non-trivial zeros, it is called an all-zero system. Such a system has a finite-duration impulse response (FIR), and is called an FIR system or a moving average (MA) system. Note that the corresponding deference equation is

$$a_0 y(n) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M)$$

The all-pole system On the other hand, if $b_k = 0$ for $1 \le k \le M$, we have

$$H(z) = \frac{b_0}{\sum_{k=0}^{N} a_k z^{-k}} = \frac{b_0}{a_0 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_N z^{-N}}$$
$$= \left| \begin{vmatrix} \sum_{k=0}^{N} \frac{b_0}{a_0 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_N z^{-N}} \\ \sum_{k=0}^{N} \frac{z^{N}}{z^{N} + (a_1 / a_0) z^{N-1} + (a_2 / a_0) z^{N-2} + \dots + (a_N / a_0)} \end{vmatrix}$$

Here again, either take $a_0 = 1$ or imagine that it is absorbed in the other coefficients viz., b_0 , a_1 , a_2 , ..., a_N . Thus

$$H(z) = \frac{b_{0} z^{N}}{z^{N} + a z^{N-1} + a z^{N-2} + \dots + a_{N}}$$

Here H(z) has N poles and an Nth order zero at the origin z = 0. We usually do not make reference to these *trivial zeros*. As a result this system function contains only *non-trivial poles* and the corresponding system is called an all-pole system. Due to the presence of the poles, the impulse response of such system is infinite in duration, and hence it is an IIR system. (We can divide the numerator into the denominator and thereby expand H(z) into an infinite series from which it is evident that h(n) is of infinite duration). Note that the corresponding deference equation is

$$a_0 y(n) + a_1 y(n-1) + ... + a_N y(n-N) = b_0 x(n)$$

The pole-zero system The general form, though, contains both poles and zeros and the system is called a pole-zero system with *N* poles and *M* zeros,

$$H(z) = \frac{b_0 + b_z^{-1} + b_z^{-2} + \dots + b_z^{-M}}{a_0 + a_z^{-1} + a_z^{-2} + \dots + a_N^{-N}} + a_N^{-N}$$

Poles and/or zeros at z = 0 and $z = \infty$ are implied but are not counted explicitly. Due to the presence of poles, the pole-zero system is an IIR system.

Inverse *z***-transform by power series expansion (long division)**

If the z-transform is expressed as a *rational function* (a ratio of polynomials in z or z^{-1}) we can use long division to expand it into a power series. If the transform is expressed as an *irrational function* we can use the appropriate power series expansion formula available in mathematical

tables such as the CRC Tables. Note that if the transform is expressed as an irrational function then the partial fraction expansion method of inversion won't work.

By definition the *z*-transform of the sequence x(n) is given by

$$X(z) = \sum_{n = -\infty} x(n) \ z^{-n} = \dots + x(-2) \ z^{2} + x(-1) \ z^{1} + x(0) \ z^{0} + x(1) \ z^{-1} + \dots$$

This is a power series (Laurent series). So by long division we obtain the power series expansion of X(z) and then, by comparison with the power series definition given above, we can identity the sequence x(n). In particular the coefficient of z^{-k} is the sequence value x(k).

The method is useful in obtaining a quick look at the first few values of the sequence x(n). This approach does not assure an analytical solution. The ROC will determine whether the series has positive or negative exponents. For right-sided sequences the X(z) will be obtained with primarily negative exponents, while left-sided sequences will have primarily positive exponents. For an annular ROC, a Laurent expansion would give both positive- and negative-time terms. This last possibility is illustrated in the example below by taking a little help from partial fractions.

Example 1.7.1 Find the inverse transform, by long division, of

$$X(z) = \frac{2 z^{2} - 3 z}{(z - 1)(z - 2)} = \frac{2 z^{2} - 3 z}{z^{2} - 3 z + 2}$$

where the ROC is (a) |z| > 2, (b) |z| < 1, (c) 1 < |z| < 2

Solution (a) ROC is |z| > 2. We expect a right-sided sequence, with predominantly negative exponents of z. For the long division arrange numerator and denominator as *decreasing powers*



of z and then divide; or as increasingly negative power of z i.e., z^{-1} and then divide.

$$D(z) \rightarrow z^{2} - 3z + 2 \qquad \begin{array}{c} 2 & +3z^{-1} + 5z^{-2} + 9z^{-3} + \dots \\ 2z^{2} - 3z \\ 2z^{2} - 3z \\ 2z^{2} - 6z + 4 \\ 3z - 4 \\ 3z - 4 \\ 3z - 9 + 6z^{-1} \\ \hline 5 - 6z^{-1} \\ 5 - 15z^{-1} + 10z^{-2} \\ 9z^{-1} - 10z^{-2} \\ 9z^{-1} - 27z^{-2} + 18z^{-3} \\ \hline 17z^{-2} - 18z^{-3} \\ \dots \end{array}$$

Thus $X(z) = 2 + 3 z^{-1} + 5 z^{-2} + 9 z^{-3} + \dots$ By comparison with the defining equation

$$X(z) = \dots x(-1)z^{1} + x(0) + x(1)z^{-1} + x(2)z^{-2} + \dots$$

we see that the sequence values are

x(-2) = x(-1) = 0, or x(n) = 0 for n < 0, and x(0) = 2, x(1) = 3, x(2) = 5, etc.

Alternatively, it is also possible to write X(z) as a ratio of polynomials in z^{-1}

$$X(z) = \frac{2 - 3 z^{-1}}{1 - 3 z^{-1} + 2 z^{-2}}$$

Note that the polynomials are written in the order of increasing negative powers of z, that is, z^{-1} . Long division gives (the same answer as obtained earlier):

$$D(z) \rightarrow 1 - 3z^{-1} + 2z^{-2} \xrightarrow{\begin{array}{c}2 + 3z^{-1} + 5z^{-2} + 9z^{-3} + \dots \\ 2 - 3z^{-1} & \leftarrow Q(z) \\ 2 - 3z^{-1} & \leftarrow N(z) \\ 2 - 6z^{-1} + 4z^{-2} & \leftarrow N(z) \\ 3z^{-1} - 4z^{-2} & 5z^{-2} - 6z^{-3} \\ 3z^{-1} - 9z^{-2} + 6z^{-3} & \\ \end{array}}$$

Solution (b) The ROC is |z| < 1. We expect a left-sided sequence with predominantly positive exponents of z. For the long division the polynomials are written in the order of increasing powers of z (or decreasingly negative powers of z, i.e., z^{-1}).



Thus $X(z) = -(3/2)z - (5/4)z^2 - (9/8)z^3 - \dots = \dots - (9/8)z^3 - (5/4)z^2 - (3/2)z$. By comparing with the defining equation

$$X(z) = \dots x(-3)z^{3} + x(-2)z^{2} + x(-1)z + x(0) + x(1)z^{-1} + \dots$$

we see that the sequence is given by

$$x(-1) = -3/2$$
, $x(-2) = -5/4$, $x(-3) = -9/8$, ... etc., and $x(n) = 0$ for $n \ge 0$

The other way of long division is shown below:

$$D(z) \rightarrow 2z^{-2} - 3z^{-1} + 1 \qquad \begin{array}{c} -(3/2)z - (5/4)z^2 - (9/8)z^3 - \dots & \leftarrow Q(z) \\ \hline -3z^{-1} + 2 \\ -3z^{-1} + (9/2) - (3/2)z \end{array} \qquad \leftarrow N(z) \\ \hline \dots \end{array}$$

(*Omit*) Solution (c) The ROC is 1 < |z| < 2. We expect a two-sided sequence with both positive and negative exponents of z. Looking at the pole-zero configuration, the pole at z = 1 implies a right-sided sequence and the pole at z = 2 a left-sided sequence. Obviously just a single long division cannot give both the left-sided and the right-sided sequences simultaneously. We shall obtain the partial fraction expansion first and then proceed with the division to obtain the sequences separately. These two sequences are then combined into one sequence to get the solution. Note that we do this only to illustrate the method of long division. But once we use partial fractions the utility of long division is nullified.



For the term $\frac{z}{z-1}$ we have a right-sided sequence given by long division thus:

$$D(z) \rightarrow z - 1 \xrightarrow{z - 1} z \xrightarrow{z - 1} z^{-1} + z^{-2} + z^{-3} + \dots + Q(z)$$

$$(z) \rightarrow z - 1 \xrightarrow{z - 1} z^{-1} \xrightarrow{z - 1} z^{-1}$$

$$(z) \rightarrow z - 1 \xrightarrow{z - 1} z^{-1} \xrightarrow{z - 2} z^{-2}$$

$$(z) \rightarrow z - 1 \xrightarrow{z - 1} z^{-2} \xrightarrow{z - 2} z^{-2}$$

The corresponding sequence is $x_R(n) = 1, n \ge 0$ 0, otherwise

For the term $\frac{z}{z-2}$ we have a left sided sequence

$$D(z) \rightarrow -2 + z \xrightarrow{\begin{array}{c} -(1/2)z - (1/4)z^2 - (1/8)z^3 \dots + Q(z) \\ z - (1/2)z^2 \\ (1/2)z^2 \\ (1/2)z^2 - (1/4)z^3 \\ \hline (1/4)z^3 - (1/8)z^4 \\ \hline (1/8)z^4 \end{array}} \leftarrow O(z)$$

The corresponding sequence is $x_L(n) = -2^{-n}$, n < 00, otherwise

The complete sequence is then

$$x(n) = x_R(n) + x_L(n) = 1, \quad n \ge 0$$

 $-2^{-n}, n < 0$

(End of Omit)

Computation of frequency response

Let the system function be given by M

$$H(z) = \frac{r \sum b r z^{-r}}{\sum_{k=0}^{r} a_k z^{-k}}$$

The frequency response is $H(e^{j\omega})$ or $H(\omega) = H(z)|_{z=e^{j\omega}}$. Thus

$$H(e^{j\omega}) = |H(\omega)| e^{j\angle H(\omega)} = \frac{r}{N} \frac{ae^{-j\omega k}}{\sum_{k=0}^{N} ae^{-j\omega k}} = \frac{\sum_{r=0}^{M} b_r \cos \omega r - j \sum_{r=0}^{M} b_r \sin \omega r}{\sum_{k=0}^{N} a_k \cos \omega k - j \sum_{k=0}^{N} a_k \sin \omega k}$$
$$= \frac{\sum_{r=0}^{M} b_r \cos \omega r - j \sum_{r=0}^{M} b_r \sin \omega r}{\sum_{k=0}^{N} a_k \cos \omega k - j \sum_{k=0}^{N} a_k \sin \omega k} = \frac{A - jB}{C - jD}$$

where

$$A = \sum_{r=0}^{M} b_r \cos \omega r, \qquad B = \sum_{r=0}^{M} b_r \sin \omega r, \qquad C = \sum_{k=0}^{N} a_k \cos \omega k, \qquad D = \sum_{k=0}^{N} a_k \sin \omega k.$$

The magnitude and phase of $H(e^{j\omega})$ are given, respectively, by

$$|H(\omega)| = \sqrt{\frac{A^2 + B^2}{C^2 + D^2}}$$
 and $\angle H(\omega) = \tan^{-1} \left(\frac{-B}{A}\right) - \tan^{-1} \left(\frac{-D}{C}\right)$

Theorem The frequency response $H(e^{j\omega})$ for a BIBO-stable system will always converge.

Accordingly every BIBO-stable system will have a frequency response and a describable steadystate response to sinusoidal inputs. But, the converse of this statement is not true, that is, the fact that $H(e^{j\omega})$ exists does not imply that the system is stable.

Example 1.8.1 [The ideal low pass filter] For the $H(\omega)$ given in figure below find h(n), the unit sample response.



Solution The unit sample response is the inverse DTFT of $H(\omega)$

$$h(n) = \frac{1}{2\pi} \int H(\omega) e^{j\omega n} d\omega = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{1}{2\pi} e^{j\omega n} d\omega = \frac{1}{2\pi} \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{1}{2\pi} e^{j\omega n} d\omega = \frac{1}{2\pi} \frac{1}{2\pi} \frac{1}{2\pi} e^{j\omega n} e^{-j\omega n} d\omega = \frac{1}{2\pi} \frac{1}{2\pi} e^{j\omega n} e^{-j\omega n} e^{-j\omega n} d\omega = \frac{1}{2\pi} \frac{1}{2\pi} e^{j\omega n} e^{-j\omega n} e^$$

It is seen that $h(n) \neq 0$ for negative *n* so that the ideal low pass filter is *noncausal*. Moreover,

although h(n) tails off as n goes from 0 to ∞ and from 0 to $-\infty$, it can be shown that $\sum_{n=\infty} h(n)$ is

not finite. This means that the ideal low pass filter is not BIBO-stable either.

Example 4.8.2 [2002] A discrete system is given by the difference equation

y(n) - 5 y(n-1) = x(n) + 4 x(n-1)

where x(n) is the input and y(n) is the output. Determine the magnitude and phase response as a function of frequency for $\omega \le \pi$. (Note that the system is not stable since it has a pole at z = 5, which is outside the unit circle. The fact that the steady state frequency response exists does not mean that the system is stable.)

Solution [See also Unit I] Taking the z-transform and with a dose of algebra we find the transfer function

$$H(z) = \frac{Y(z)}{X(z)} = \frac{z+4}{z-5}$$

The frequency response is given by

$$H(\omega) = H(z)\Big|_{z=e^{j\omega}} = \frac{z+4}{z-5}\Big|_{z=e^{j\omega}} = \frac{e^{j\omega}+4}{e^{j\omega}-5} = \frac{N(\omega)}{D(\omega)}$$

$$N(\omega) = e^{j\omega}+4 = \cos \omega + 4 + j \sin \omega = |N(\omega)| e^{j\angle N(\omega)}$$

$$|N(\omega) = e^{j\omega}+4 = \cos \omega + 4 + j \sin \omega = |N(\omega)| e^{j\angle N(\omega)}$$

$$|N(\omega) = tan^{-1}|| || || || (\cos \omega + 4)^{2} + (\sin \omega)^{2}$$
and
$$\angle N(\omega) = tan^{-1}|| || || || (\cos \omega + 4)^{2}$$

$$|D(\omega)| = \sqrt{(\cos \omega - 5)^{2} + (\sin \omega)^{2}}$$
and
$$\angle D(\omega) = tan^{-1}\left(\frac{\sin \omega}{(\cos \omega - 5)^{2} + (\sin \omega)^{2}}\right)$$

$$|H(\omega)| = \frac{|N(\omega)|}{|D(\omega)|} = \frac{\sqrt{(\cos \omega + 4)^{2} + (\sin \omega)^{2}}}{\sqrt{(\cos \omega - 5)^{2} + (\sin \omega)^{2}}}$$

$$\angle H(\omega) = \angle N(\omega) - \angle D(\omega) = tan^{-1}|| || || + tan^{-1}|| (\cos \omega - 5)||$$

The frequency response can be plotted. Note that $|H(\omega)|$ is an even function and $\angle H(\omega)$ is an odd function of ω .

Using MATLAB:

$$H(\omega) = \frac{e^{-j\omega} + 4}{e^{-j\omega} - 5} = \frac{1 + 4e^{-j\omega}}{1 - 5e^{-j\omega}} = \frac{b(1) + b(2)e^{--j\omega} + b(3)e^{-j2\omega} + \dots}{a(1) + a(2)e^{-j\omega} + a(3)e^{-j2\omega} + \dots}$$

Here the vectors *b* and *a* specify, respectively, the numerator and denominator coefficients. In our example b(1) = 1, b(2) = 4, a(1) = 1, and a(2) = -5. The MATLAB segment and the corresponding plots follow. Note that the plot goes from -2π to 2π . Compare with the solution obtained in Unit I using a different function.

b = [1, 4]; %Numerator coefficients a = [1, -5]; %Denominator coefficients w = -2*pi: pi/256: 2*pi; [h] = freqz(b, a, w); subplot(2, 1, 1), plot(w, abs(h)); xlabel('Frequency \omega'), ylabel('Magnitude'); grid





Example 1.8.3 Assume $H(z) = \frac{12z^2 - 1}{6z^2 - z - 1}$ is a causal system. Find the difference equation and

the frequency response. **Solution** Arrange H(z) is

rrange
$$H(z)$$
 in terms of negative powers of z

$$H(z) = \frac{Y(z)}{X(z)} = \frac{z^2(12 - z^{-2})}{z^2(6 - z^{-1} - z^{-2})} = \frac{(12 - z^{-2})}{(6 - z^{-1} - z^{-2})}$$

Cross multiplying

$$Y(z) (6 - z^{-1} - z^{-2}) = X(z) (12 - z^{-2})$$

$$6Y(z) - z^{-1}Y(z) - z^{-2}Y(z) = 12X(z) - z^{-2}X(z)$$

Taking the inverse z-transform

$$6y(n) - y(n-1) - y(n-2) = 12x(n) - x(n-2)$$
$$y(n) = \frac{1}{6}y(n-1) + \frac{1}{6}y(n-2) + 2x(n) - \frac{1}{6}x(n-2)$$

The poles of H(z) are located at

$$z_{1,} z_{2} = \frac{-(-1) \pm \sqrt{(-1)^{2} - 4(6)(-1)}}{2(6)} = \frac{1 \pm \sqrt{+24}}{12} = \frac{1 \pm 5}{12} = 0.5 \text{ and } -1/3$$

and are inside the unit circle. This being a causal system, the ROC is $|z| > \frac{1}{2}$ and contains the unit circle. The system is stable, and the frequency response is meaningful. It is given by

$$H(\omega) = H(z)\Big|_{z=e} \frac{12z^{2}-1}{6z^{2}-z-1}\Big|_{z=e^{j\omega}} \frac{12(e^{j\omega})^{2}-1}{6(e_{j\omega})^{2}-e^{j\omega}-1} = \frac{N(\omega)}{D(\omega)}$$

where

$$N(\omega) = 12(e^{j\omega})^{2} - 1 = 12e^{j2\omega} - 1 = 12\cos 2\omega + j12\sin 2\omega - 1$$

= $\sqrt{(12\cos 2\omega - 1)^{2} + (12\sin 2\omega)^{2}} e^{j\tan^{-1}\left(\frac{j2\sin 2\omega}{12\cos 2\omega - 1}\right)^{2}}$
 $D(\omega) = 6(e^{j\omega})^{2} - e^{j\omega} - 1 = 6e^{j2\omega} - e^{j\omega} - 1$
= $6\cos 2\omega + j6\sin 2\omega - \cos \omega - j\sin \omega - 1$
= $\sqrt{(6\cos 2\omega - \cos \omega - 1)^{2} + (6\sin 2\omega - \sin \omega)^{2}} e^{j\tan^{-1}\left(\frac{6\sin 2\omega - \sin \omega}{6\cos 2\omega - \cos \omega - 1}\right)^{2}}$

The magnitude response is given by

$$|H(\omega)| = \frac{\sqrt{12\cos 2\omega - 1)^2 + (12\sin 2\omega)^2}}{\sqrt{(6\cos 2\omega - \cos \omega - 1)^2 + (6\sin 2\omega - \sin \omega)^2}}$$

The phase response is given by
$$\angle H(\omega) = \tan^{-1} \left(\begin{array}{c} 12\sin 2\omega \\ -\tan^{-1} \left(\begin{array}{c} 0 \\ -1 \end{array} \right) & -\left(\begin{array}{c} 6\sin 2\omega - \sin \omega \\ 6\cos 2\omega - \cos \omega - 1 \end{array} \right) \\ (6\cos 2\omega - \cos \omega - 1) \end{array} \right)$$

Using MATLAB:

$$H(z) = \frac{12z - 1}{6z^{2} - z - 1} = \frac{12 - z}{6z^{2} - z - 1} = \frac{12 - z}{6z^{2} - z - 1}$$

$$H(\omega) = \frac{12 - e^{-j2\omega}}{6 - e^{-j\omega} - e^{-j2\omega}} = \frac{b(1) + b(2)e^{-j\omega} + b(3)e^{-j2\omega} + \dots}{a(1) + a(2)e^{-j\omega} + a(3)e^{-j2\omega} + \dots}$$

Here the vectors *b* and *a* specify, respectively, the numerator and denominator coefficients. In our example b(1) = 12, b(2) = 0, b(3) = -1, a(1) = 6, a(2) = -1 and a(3) = -1. The MATLAB segment and the corresponding plots follow. Note that the plot goes from $-\pi$ to π .

b = [12, 0, -1]; %Numerator coefficients a = [6, -1, -1]; %Denominator coefficients w = -pi: pi/256: pi; [h] = freqz(b, a, w); subplot(2, 1, 1), plot(w, abs(h)); xlabel('Frequency \omega (Rad)'), ylabel('Magnitude'); grid subplot(2, 1, 2), plot(w, angle(h)); xlabel('Frequency \omega (Rad)'), ylabel('Phase - Radians'); grid



Example 1.8.4 Discuss the stability of $H(z) = \frac{z^{-1}}{1 - z^{-1} - z^{-2}}$ assuming it is a causal

system. Find the difference equation and the frequency response.

(b) Determine the frequency, magnitude and phase responses and time delay for the system y(n) + (1/4) y(n-1) = x(n) - x(n-1).

Solution

(a) Find the ROC and the poles:

$$H(z) = \frac{z^{-1}}{1 - z^{-1} - z^{-2}} = \frac{z}{z^2 - z^1 - 1}$$

There is a zero at z = 0. The denominator has roots at

$$z_{1}, z_{2} = \frac{-(-1) \pm \sqrt{(-1)^{2} - 4(1)(-1)}}{2} = \frac{1 \pm \sqrt{1 + 4}}{2} = \frac{1 \pm \sqrt{5}}{2} = 1.62 \text{ and } -0.62$$



The pole locations are shown here. For the system to be causal the ROC is the exterior of a circle with radius = 1.62. In this case ROC does not include the unit circle. (Equivalently, all the poles do not lie within the unit circle). Hence the system is not stable.

(b) Taking the *z*-transform on both sides of y(n) + (1/4) y(n-1) = x(n) - x(n-1) we get

$$Y(z) + (1/4) z^{-1} Y(z) = X(z) - z^{-1} X(z), \text{ or}$$

$$Y(z) \{1 + (1/4) z^{-1}\} = X(z) \{1 - z^{-1}\}, \text{ or}$$

$$H(z) = \frac{Y(z)}{X(z)} = \frac{1 - z}{1 + (1/4)z^{-1}} = \frac{z - 1}{z + 0.25}$$

There is a single zero at z = 1 and a single pole at z = -0.25 which is inside the unit circle – hence stable. The frequency response is given by

$$H(\omega) = H(z)\Big|_{z=e^{j\omega}} = \frac{z-1}{z+0.25}\Big|_{z=e^{j\omega}} = \frac{e^{-j\omega}-1}{e^{-j\omega}+0.25} = \frac{N(\omega)}{D(\omega)}$$

where

$$N(\omega) = e^{j\omega} - 1 = \cos\omega + j\sin\omega - 1 = \sqrt{(\cos\omega - 1)^2 + \sin^2\omega} e^{j\tan^{-1}\left(\frac{\sin\omega}{\cos\omega - 1}\right)^2}$$
$$D(\omega) = e^{j\omega} - 0.25 = \cos\omega + j\sin\omega - 0.25$$
$$= \sqrt{(\cos\omega - 0.25)^2 + \sin^2\omega} e^{j\tan^{-1}\left(\frac{\sin\omega}{\cos\omega - 0.25}\right)^2}$$

The magnitude response is given by

$$|H(\omega)| = \frac{\sqrt{(\cos \omega - 1)^2 + \sin^2 \omega}}{\sqrt{(\cos \omega - 0.25)^2 + \sin^2 \omega}}$$

The phase response is given by

$$\angle H(\omega) = \tan^{-1} \begin{pmatrix} \sin \\ \Box & -i \\ \cos \omega & -1 \end{pmatrix} - \begin{pmatrix} \omega \sin \\ \cos \omega & -i \\ \cos \omega & -1 \end{pmatrix} - \begin{pmatrix} \omega \sin \\ \cos \omega & -i \\ \cos \omega & -0.25 \end{pmatrix}$$

The time (group) delay is given by $-\frac{d\omega}{d\omega} LH(\omega)$.

Using MATLAB:

$$H(\omega) = e^{j\omega} - 1$$

 $= \frac{1 - e^{-j\omega}}{e^{-j\omega} + 0.25} = \frac{1 - e^{-j\omega}}{1 + 0.25e} = \frac{b(1) + b(2)e^{-j\omega} + b(3)e^{-j2\omega} + \dots}{a(1) + a(2)e^{-j\omega} + a(3)e^{-j2\omega} + \dots}$

Here the vectors *b* and *a* specify, respectively, the numerator and denominator coefficients. In our example b(1) = 1, b(2) = -1, a(1) = 1, a(2) = 0.25. The MATLAB segment and the corresponding plots follow. Note that the plot goes from $-\pi$ to π .

b = [1, -1]; %Numerator coefficients a = [1, 0.25]; %Denominator coefficients w = -pi: pi/256: pi; [h] = freqz(b, a, w); subplot(2, 1, 1), plot(w, abs(h)); xlabel('Frequency \omega (Rad)'), ylabel('Magnitude'); grid subplot(2, 1, 2), plot(w, angle(h)); xlabel('Frequency \omega (Rad)'), ylabel('Phase - Radians'); grid



Z-transforms with initial conditions

To solve the N^{th} order difference equation

$$y(n) = -\sum_{k=1}^{N} a_k y(n-k) + \sum_{r=0}^{M} b_r x(n-r)$$

with (non-zero) initial conditions we need *N* initial conditions on the output y(n) and *M* initial conditions on the input x(n). Usually the input is applied suddenly (i.e., it is stepped into the system) at n = 0, so that no initial conditions are needed for it, that is, x(n) = 0 for n < 0. The output y(n), however, in general, will have non-zero initial conditions for n = -1 to -N.

We are solving for y(n) for $n \ge 0$, so that $Y(z) = \Im\{y(n)\}$ is the one-sided z-transform. The difference equation contains other terms like y(n-1), y(n-2), etc. which are delayed versions of y(n). Suppose N = 3, then we shall have y(n-1), y(n-2), and y(n-3) to deal with. The transform of y(n-1) is handled as follows. First, for the sequence y(n) as shown below we define

 $Y_{+}(z) = \Im\{y(n), n \ge 0\} = A + Bz^{-1} + Cz^{-2} + \dots$

We shall refer to this loosely as just Y(z) when there is no possibility of confusion.



We then obtain y(n-1) by delaying the sequence by one unit, shown below.



As can be seen from the graph

$$\Im\{y(n-1), n \ge 0\} = y(-1)z^0 + Az^{-1} + Bz^{-2} + Cz^{-3} + \dots = z^{-1}Y_+(z) + y(-1)$$

 $z^{-1}Y_+(z)$

In a similar fashion

$$\mathbf{z}\{y(n-2), n \ge 0\} = y(-2)z^0 + y(-1) z^{-1} + Az^{-2} + Bz^{-3} + Cz^{-4} + \dots$$
$$z^{-2}Y_+(z)$$

$$= z^{-2}Y(z) + y(-2) + y(-1)z^{-1}$$

and by extension

$$\mathbf{z}\{y(n-3), n \ge 0\} = z^{-3}Y_+(z) + y(-3) + y(-2)z^{-1} + y(-1)z^{-2}$$

would be the last. But we can generalize

For N = 3 this would be the last. But we can generalize

$$\mathbf{z}\{y(n-k), n \ge 0\} = z^{-k}Y_{+}(z) + y(-k) + y(-k-1)z^{-1} + \dots + y(-1)z^{-(k-1)}$$

In the case of the input x(n), since it is applied suddenly at n = 0, the initial conditions are zero, that is, x(-1) = x(-2) = ... = x(-M) = 0, so that

$$\{x(n-k) u(n)\} = z^{-k}X_{+}(z)$$

With this intuitive background we give below the mathematical derivation of the *z*-transform of the delayed truncated sequence.

Z-transform of delayed truncated sequence (initial conditions) The one-sided *z*-transform of x(n) is

$$X_{+}(z) = \mathbf{z}\{x(n) \ u(n)\} = \sum_{n=0}^{\infty} x(n) \ z^{-n}$$

Given the sequence x(n), we delay it by k units, and then truncate it to the left of n = 0 to get x(n-k) u(n). We want find the z-transform of x(n-k) u(n).

$$\Im\{x(n-k) \ u(n)\} = \sum_{n=-\infty}^{\infty} x(n-k) \ u(n) \ z^{-n} = \sum_{n=0}^{\infty} x(n-k) \ z^{-n}$$

If we let n-k = r, then n = r+k, and the summation limits n = 0 to ∞ become r = -k to ∞ . Then

$$\mathbf{z}\{x(n-k) \ u(n)\} = \sum_{r=-k} x(r) \ z^{-(r+k)} = z^{-k} \sum_{r=-k} x(r) \ z^{-r}$$
$$= z^{-k} \left[\sum_{r=0}^{\infty} x(r) \ z^{-r} + \sum_{r=-k}^{-1} x(r) \ z^{-r} \right] = z^{-k} \left[X_{+}(z) + \sum_{r=-k}^{-1} x(r) \ z^{-r} \right]$$

IC

 $X_+(z)$

$$= z^{-k} \left[X_{+}(z) + x(-k) z^{k} + x(-k-1) z^{k-1} + \dots + x(-1) z^{1} \right]$$

= $z^{-k} X_{-k}(z) + x(-k) + x(-k-1) z^{-1} + \dots + x(-1) z^{-(k-1)}$

Due to Initial Conditions

We shall loosely refer to $X_+(z)$ as X(z) and write the result as

$$\mathbf{z}\{x(n-k) \ u(n)\} = z^{-k} X(z) + x(-k) + x(-k-1) \ z^{-1} + \dots + x(-1) \ z^{-(k-1)}$$

Due to Initial Conditions

The above result is used to solve linear constant coefficient difference equations with *inputs that are stepped into a system*. Suppose we want the solution of

$$\sum_{k=0}^{m} a_k y(n-k) = \sum_{r=0}^{m} b_r x(n-r), \qquad n \ge 0$$

subject to the initial conditions

 $\{y(i), i = -1, -2, ..., -N\}$ and $\{x(i), i = -1, -2, ..., -M\}$

We take the *z*-transform of the equation using the result derived above for delayed-truncated sequences

$$\mathbf{z} \left\{ \sum_{k=0}^{N} a_{k} y(n-k) \right\} = \mathbf{z} \left\{ \sum_{r=0}^{M} b_{r} x(n-r) \right\}, \qquad n \ge 0$$

$$\sum_{k=0}^{N} a_{k} Z\{y(n-k)\} = \sum_{r=0}^{M} b_{r} Z\{x(n-r)\}, \qquad n \ge 0$$

where we have used Z to mean \underline{z} the z-transformation operation. The left hand side is

LHS =
$$a_0 \Im\{y(n)\} + a_1 \Im\{y(n-1)\} + a_2 \Im\{y(n-2)\} + \dots + a_N \Im\{y(n-N)\}$$

= $a_0 Y(z) + a_1 \{z^{-1} Y(z) + y(-1)\} + a_2 \{z^{-2} Y(z) + y(-2) + y(-1) z^{-1}\} + \dots + a_N \{z^{-N} Y(z) + y(-N) + y(-(N-1)) z^{-1} + \dots + y(-1) z^{-N+1}\}$

(Note that in terms of the derivation earlier all of the Y(z)'s are $Y_+(z)$'s, i.e., one-sided transforms). All the Y(z) terms can be grouped together under a summation, and all the remaining terms, due to the initial conditions $\{y(i), i = -1, -2, ..., -N\}$, can be grouped together so that the above can be written as

LHS =
$$\sum_{k=0}^{\infty} a_k z^{-k} Y(z) + g\{z^{-1}, y(-1), y(-2), ..., y(-N)\}$$

Initial condition terms

By following a similar procedure the right hand side can also be written as follows (here again the X(z)'s are $X_+(z)$'s, i.e., one-sided transforms):

RHS =
$$\sum_{r=0}^{\infty} b z^{-r} X(z) + h\{z^{-1}, x(-1), x(-2), ..., x(-M)\}$$

Initial condition terms

Writing out in full, LHS = RHS becomes

$$\sum_{k=0}^{N} a_k z^{-k} Y(z) + g\{\ldots\} = \sum_{r=0}^{M} b_r z^{-r} X(z) + h\{\ldots\}$$

Factoring out Y(z) and X(z) and rearranging we have

$$Y(z) \sum_{k=0}^{N} a_{k} z^{-k} = X(z) \sum_{r=0}^{M} b_{r} z^{-r} + h\{...\} - g\{...\}$$
$$Y(z) = X(z) \frac{\sum_{r=0}^{M} b_{r} z^{-r}}{\sum_{k=0}^{N} a_{k} z^{-k}} + \frac{h\{...\} - g\{...\}}{\sum_{k=0}^{N} a_{k} z^{-k}}$$

Taking the inverse *z*-transform we get

$$y(n) = \mathbf{z} \left\{ X(z) \left[\sum_{k=0}^{M} b_{r} z^{-r} \right] + \mathbf{z} \left\{ \frac{h\{...\} - g\{...\}}{\sum_{k=0}^{N} a_{k,7} - k} \right\} + \mathbf{z} \left\{ \frac{h\{...\} - g\{...\}}{\sum_{k=0}^{N} q_{k,7} - k} \right\}$$

To summarize: to solve for y(n) we take the *z*-transform of the linear constant coefficient difference equation using initial conditions, manipulate in the *z*-domain to get Y(z) and then take the inverse *z*-transform of Y(z) to get y(n).

Example 1.9.1 Find the solution to

$$y(n) - \frac{3}{2}y(n-1) + \frac{1}{2}y(n-2) = \begin{vmatrix} 1 \\ -1 \end{vmatrix}, \qquad n \ge 0$$

with initial conditions y(-1) = 4, y(-2) = 10. Solution There are three methods of solution:

- 1. Find the iterative solution in the discrete-time domain. In general this will not give an analytical (closed) form of solution.
- 2. Solve in the discrete-time domain (homogeneous solution + particular solution).
- 3. Solve in the frequency domain as we do below.

For an *input sequence* x(n) that is stepped into a system, specified in words like x(n) = 0 for n < 0, the initial conditions are clearly zero and do not matter. But for an *output sequence* y(n) where the initial conditions y(-1), y(-2) are explicitly given to be non-zero we need to use the above derived "*z*-transform for delayed truncated sequence". In particular we have

$$\begin{aligned} \mathbf{\mathfrak{Z}}\{y(n)\} &= Y(z) \\ \mathbf{\mathfrak{Z}}\{y(n-1)\} &= z^{-1} \Big[Y(z) + y(-1) z^1 \Big] \\ \mathbf{\mathfrak{Z}}\{y(n-2)\} &= z^{-2} \Big[Y(z) + y(-2) z^2 + y(-1) z^1 \Big] \end{aligned}$$

Taking the z-transform of the difference equation we get

$$\mathbf{J}\left(\begin{array}{c} y(n) -\frac{3}{2}y(n-1) + \frac{1}{2}y(n-2)\right) = \mathbf{J}\left[\left(\left(\frac{1}{4}\right)^{n}\right)\right] \qquad n \ge 0$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right) = \frac{1}{2} \mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)\right] \qquad n \ge 0$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right) = \frac{1}{2} \mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)^{n}\right)$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right) = \frac{1}{2} \mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)^{n}\right)$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right) = \frac{1}{2} \mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)^{n}\right)$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right) = \frac{1}{2} \mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)^{n}\right)$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)^{n}\right)$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)^{n}\right)$$

$$\mathbf{J}\left(\left(\frac{1}{4}\right)^{n}\right)$$

$$\mathbf{J$$

Plugging in the initial conditions
$$y(-1) = 4$$
 and $y(-2) = 10$

$$\begin{array}{c} Y(z) - \frac{3}{2} \left[z^{-1}Y(z) + 4\right] + \frac{1}{2} \left[z^{-2}Y(z) + 10 + 4z^{-1}\right] = \frac{z}{z - (1/4)} \\ \left[1 - 2z^{-1} + \frac{1}{2}z^{-2}\right] - 6 + 5 + 2z = \overline{z - (1/4)} \\ Y(z) \left[1 - 1.5z^{-1} + 0.5z^{-2}\right] = \frac{z}{z - (1/4)} + 1 - 2z^{-1} = \frac{2z^{2} - (9/4)z + (1/2)}{z(z - (1/4))} \\ Y(z) (1 - z^{-1}) (1 - 0.5z^{-1}) = \frac{2z^{2} - (9/4)z + (1/2)}{z(z - (1/4))} \\ Y(z) \left(\frac{z - 1)(z - (1/2)}{z^{2}}\right) = \frac{2z^{2} - (9/4)z + (1/2)}{z(z - (1/4))} \\ Y(z) \left(\frac{z - 1)(z - (1/2)}{z^{2}}\right) = \frac{2z^{2} - (9/4)z + (1/2)}{z(z - (1/4))} \\ Y(z) \left(\frac{z - (1/4))(z - 1)(z - (1/2)}{z(z - (1/4))}\right) \\ Y(z) = \frac{(2z^{2} - (9/4)z + (1/2))}{(z - (1/2)(z - (1/2))} = \frac{A}{z - (1/4)} + \frac{B}{z - 1} + \frac{C}{z - (1/2)} \\ A = \frac{(2z^{2} - (9/4)z + (1/2))}{(z - (1/2)(z - (1/2))} = 1/3 \\ B = \frac{(2z^{2} - (9/4)z + (1/2))}{(z - (1/4))(z - (1/2))} \\ z = 1 \\ C = \frac{(2z^{2} - (9/4)z + (1/2))}{(z - (1/4))(z - (1/2))} \\ z = 1 \\ \frac{Y(z)}{z - (1/4)(z - (1/2))} \\ z = \frac{1}{z^{2} - (1/4)} + \frac{2z}{z^{2} - 1} + \frac{1}{z - (1/2)} \\ Y(z) = \frac{1}{z^{2} - (1/4)} + \frac{2}{3z^{2} - 1} + \frac{z}{z - (1/2)} \\ y(n) = \begin{vmatrix} 1 \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \frac{2}{1^{n}} + \begin{vmatrix} 1 \\ 1 \end{pmatrix} \\ y(n) = \begin{vmatrix} 1 \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \frac{2}{1^{n}} + \begin{vmatrix} 1 \\ 1 \end{pmatrix} \\ y(n) \end{vmatrix}$$

The iterative solution for this problem was obtained in Unit I. The time-domain solution was covered in HW (Extra). The solution is repeated below

$$y(n) = \left\{2, \frac{5\ 15\ 51}{4\ 16\ 64}, \dots\right\}$$

Example 4.9.2 [2002] Solve the following linear difference equation

$$y(n) + \frac{1}{2}y(n-1) - \frac{1}{4}y(n-2) = 0$$

given y(-1) = y(-2) = 1.

Solution Note that the output is a natural response corresponding to the specified initial conditions. There is no forced response since x(n) = 0.

For the solution in the frequency domain we take the z-transform of the difference equation

$$\begin{aligned} \mathbf{J}\left(y(n) + \frac{1}{2}y(n-1) - \frac{1}{4}y(n-2)\right) &= \mathbf{J}\left(0\right) \qquad n \ge 0\\ \mathbf{J}\left\{y(n)\right\} + \frac{1}{2} \ \mathbf{J}\left\{y(n-1)\right\} - \frac{1}{4} \ \mathbf{J}\left\{y(n-2)\right\} &= 0\\ Y(z) + \frac{1}{2} \ z^{-1}\left[Y(z) + y(-1)z^{1}\right] - \frac{1}{4} \ z^{-2}\left[Y(z) + y(-2)z^{2} + y(-1)z^{1}\right] &= 0\\ Y(z) + \frac{1}{2} \ z^{-1}\left[Y(z) + 1z^{1}\right] - \frac{1}{4} \ z^{-2}\left[Y(z) + 1z^{2} + 1z^{1}\right] &= 0\\ Y(z) \left[1 + 0.5 \ z^{-1} - 0.25 \ z^{-2}\right] &= -0.25 + 0.25 z^{-1}\\ Y(z) &= \frac{(-0.25)z(z-1)}{z^{2} + 0.5z - 0.25}\\ \frac{Y(z)}{z} &= \frac{(-0.25)(z-1)}{z^{2} + 0.5z - 0.25}\end{aligned}$$

The denominator on the right hand side has roots at

$$z_{1}, z_{2} = \frac{-0.5 \pm \sqrt{(0.5)^{2} - 4(1)(-0.25)}}{2} = \frac{-1 \pm \sqrt{5}}{4} = 0.309 \text{ and } -0.809$$

$$\frac{Y(z)}{z} = \frac{(-0.25)(z-1)}{(z-0.309)(z+0.809)^{2}} = z - 0.309^{+} z + 0.809$$

$$A = \frac{(-0.25)(z-1)}{(z+0.809)} = 0.155$$

$$B = \frac{(-0.25)(z-1)}{(z-0.309)} = -0.405$$

$$Y(z) = \frac{0.155z}{z-0.309} - \frac{0.405z}{z+0.809}$$

$$y(n) = 0.155 \quad (0.309)^{n} - 0.405(-0.809)^{n}, \quad n \ge 0$$

The first few values of the sequence are $y(n) = \{-0.25, 0.376, -0.25, 0.219, ...\}$ and should be compared with the iterative solution.

In the context of MATLAB, we may use *filter*(*b*, *a*, *x*) to generate the sequence y(n). The coefficients of y(.) and x(.) are numbered slightly differently as below:

$$a_1 y(n) + a_2 y(n-1) + a_3 y(n-2) + \dots = b_1 x(n) + b_2 x(n-1) + b_3 x(n-2) + \dots$$

From the difference equation

$$y(n) + \frac{1}{2}y(n-1) - \frac{1}{4}y(n-2) = 0, \qquad n \ge 0$$

we note that the input is x(n) = 0 and the coefficients of y(.) and x(.) give us the *a* and *b* vectors: a = [1, 0.5, -02.5] and b = [1]. The non-zero initial conditions y(-1) = y(-2) = 1 must first be converted to *equivalent initial conditions* for the *filter* function to work. We specify the vector yic = [y(-1), y(-2)] = [1, 1] and generate the equivalent initial conditions *eic* by the function *filtic*(*b*, *a*, *yic*). The equivalent initial conditions are then used to generate the filter output through *filter*(*b*, *a*, *xn*, *eic*). The MATLAB segment follows: %Non-zero initial conditions b = [1], a = [1, 0.5, -0.25], yic = [1, 1], n = 0:25, xn = 0 .*n % %Equivalent initial conditions eic = filtic(b, a, yic), yn = filter(b, a, xn, eic) subplot(2, 1, 1), stem(n, xn); xlabel('n'), ylabel('xn'); title('Input Sequence'); subplot(2, 1, 2), stem(n, yn); xlabel('n'), ylabel('yn'); title('Output Sequence');

The output is:

yn = [-0.25 0.375 -0.25 0.2188 -0.1719 0.1406 -0.1133 ... 0.0031 -0.0025 0.0020]





Assume the system is initially relaxed. Obtain the system function and plot its poles and zeros. **Solution** The phrase "initially relaxed" means that the initial conditions are zero, that is, y(n) = 0 for n < 0 and x(n) = 0 for n < 0. The question doesn't specify what the input x(n) is, so assume $\delta(n)$. (What will the output be if both the input *and* the initial conditions are zero?)

Steady-state and transient responses for a first order system

We consider sinusoidal inputs. Although the presentation is only for a first order system, the relationship established for the steady-state response in terms of the transfer function of the system is a general result for stable systems and sinusoidal inputs.

The system is

$$y(n) = a y(n-1) + x(n), \qquad n \ge 0$$

with the initial condition y(-1) and the input $x(n) = \cos \omega_0 n u(n)$. (We have considered the timedomain behavior of this system in Unit I). Assume |a| < 1 in order to have a stable system. The system function is obtained with zero initial conditions,

$$Y(z) = a z^{-1} Y(z) + X(z), \text{ or}$$
$$H(z) = \frac{Y(z)}{X(z)} = \frac{1}{1 - az^{-1}} = \frac{z}{z - a}$$

The solution of the difference equation is obtained by taking the *z*-transform and using the given initial condition

$$Y(z) = a [z^{-1} Y(z) + y(-1)] + X(z)$$

$$Y(z) [1 - az^{-1}] = a y(-1) + X(z)$$

$$Y(z) = a y(-1) \frac{1}{1 - az^{-1}} + X(z) \frac{1}{1 - az^{-1}} = a y(-1) H(z) + X(z) H(z)$$

Since $X(z) = \Im\{x(n)\} = \Im\{\cos\omega_0 n \ u(n)\} = \frac{z - z \cos\omega_0}{z - 2z \cos\omega_0 + 1}$, we have

$$Y(z) = a y(-1) H(z) + \frac{z - z \cos \omega_0}{z^2 - 2z \cos \omega_0 + 1} H(z) = Y_1(z) + Y_2(z)$$

$$=Y_1(z) \qquad \qquad =Y_2(z)$$

Here $Y_1(z)$ is the *zero-input response* due to the initial condition(s) a y(-1) z

$$Y_{l}(z) = a y(-1) H(z) = \frac{a y(-1)}{z - a}$$

and $Y_2(z)$ is the *forced response* due to the input x(n)

$$Y_{2}(z) = \frac{z^{2} - z\cos\omega_{0}}{z^{2} - 2z\cos\omega_{0} + 1} \quad H(z) = \frac{z\left(z^{2} - z\cos\omega_{0}\right)}{(z-a)\left(z^{2} - 2z\cos\omega_{0} + 1\right)}$$

 $Y_1(z)$ is already in a convenient form for taking the inverse, but $Y_2(z)$ must be expanded into partial fractions as below.

$$\frac{I'(z)}{z} = \frac{z(z - \cos\omega_0)}{(z - a)(z^2 - 2z\cos\omega_0 + 1)_*}$$

$$= \frac{A}{z - a} + \frac{B}{z - e^{j\omega_0}} + \frac{B}{e^{-j\omega_0}} z - z^2 \cos\omega_0 + 1$$

$$A = \frac{z(z - \cos\omega_0)}{(z^2 - 2z\cos\omega_0 + 1)_{z = a}}$$

$$= \frac{a(a - \cos\omega_0)}{a^2 - 2a\cos\omega_0^2 + 1}$$

$$= \frac{a(a - \cos\omega_0)}{a^2 - 2a\cos\omega_0^2 + 1}$$

Factoring of...

$$z^2 - 2z\cos\omega_0 + 1$$

$$= z^2 - 2z \left[\left(\frac{e^{j\omega_0} + e^{-j\omega_0}}{2} \right) \right] + 1$$

$$= z^2 - ze^{j\omega_0} - ze^{-j\omega_0} + 1$$

$$= z(z - e^{j\omega_0}) - e^{-j\omega_0}(z - e^{j\omega_0})$$

$$= (z - e^{j\omega_0})(z - e^{-j\omega_0})$$

1.

`

$$B = \frac{z (z - \cos \omega_{0})}{(z - a)(z - e^{-j\omega_{0}})}\Big|_{z = e^{j\omega_{0}}}$$

$$= \frac{e^{j\omega_{0}}(e^{-\alpha} - \cos \omega)}{(e^{j\omega_{0}} - a)(e^{j\omega_{0}} - e^{-j\omega_{0}})} = \frac{e^{j\omega_{0}}\left[e^{j\omega_{0}} - e^{-j\omega_{0}}\right]}{(e^{j\omega_{0}} - a)(e^{j\omega_{0}} - e^{-j\omega_{0}})}$$

$$= \frac{e^{j\omega_{0}}\left[\frac{e^{j\omega_{0}} - e^{-j\omega_{0}}}{2}\right]}{(e^{j\omega_{0}} - a)(e^{j\omega_{0}} - e^{-j\omega_{0}})}$$

$$= \frac{e^{j\omega_{0}}}{2(e^{j\omega_{0}} - a)} = \frac{1}{2}H(z)\Big|_{z = e^{j\omega_{0}}}$$

$$B^{*} = \frac{e^{-j\omega_{0}}}{2(e^{-j\omega_{0}} - a)}$$

So

$$Y(z) = Y_1(z) + Y_2(z) = Y_1(z) + \frac{Y_2(z)}{z}z$$

= $\frac{a \ y(-1) \ z}{z-a} + \frac{A \ z}{z-a} + \frac{B \ z}{z-e^{j\omega_0}} + \frac{B \ z}{z-e^{-j\omega_0}}$

Taking the inverse *z*-transform we get

$$y(n) = a \ y(-1) \ a^{n}u(n) + A \ a^{n}u(n) + B\left(e^{j\omega_{0}}\right)^{n}u(n) + B^{*}\left(e^{-j\omega_{0}}\right)^{n}u(n)$$

= y₁(n)

$$y(n) = [a \ y(-1) + A] \ a^{n}u(n) + B(e^{j\omega_{0}})^{n}u(n) + B^{*}(e^{-j\omega_{0}})^{n}u(n)$$

Using the fact that
$$B = \left(e^{j\omega_0} = \frac{1}{2} H(z) \right)_{z=e^{j\omega_0}} = \operatorname{Re}\left[Be^{j\omega_0 n} \right] u(n)$$

Using the fact that $B = \left(e^{j\omega_0} = \frac{1}{2} H(z) \right)_{z=e^{j\omega_0}} = \operatorname{Re}\left[Be^{j\omega_0 n} \right] u(n)$

$$2\operatorname{Re}\left[Be_{j\omega_0 n} \right] = 2\operatorname{Re}\left[\frac{1}{2} H(z) \right]_{z=e^{j\omega_0}} e^{j\omega_0 n} = \operatorname{Re}\left[H(z) \right]_{z=e^{j\omega_0}} e^{j\omega_0 n}$$
Since $H(z) = 2\operatorname{Re}\left[Be_{j\omega_0} \right] = \operatorname{Re}\left[Be_{j\omega_0} \right] = \operatorname{Re}\left[Be_{j\omega_0} \right] = \operatorname{Re}\left[Be_{j\omega_0} \right]_{j \ge H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0))} \right]_{j \ge H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0))} = \operatorname{Re}\left[H(j\omega_0) \right]_{0} e^{j(\omega_0 n+2H(j\omega_0))} = \operatorname{Re}\left[H(j\omega_0) \right]_{0} e^{j(\omega_0 n+2H(j\omega_0))} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0))} e^{j(\omega_0 n+2H(j\omega_0))} e^{j(\omega_0 n+2H(j\omega_0))} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0))} e^{j(\omega_0 n+2H(j\omega_0))} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0))} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0)})} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(j\omega_0)} e^{j(\omega_0 n+2H(\omega_0)} e^{j(\omega_0 n+2H($

Thus y(n) becomes

$$y(n) = [a \ y(-1) + A] \ a^{n}u(n) + |H(j\omega_{0})| \cos(\omega_{0}n + \angle H(j\omega_{0})u(n))$$

Transient response Steady-state response

Since |a| < 1 the transient term will eventually go to zero as $n \to \infty$. Even if the initial condition is zero, y(-1) = 0, there is still a transient response $Aa^n u(n)$ which eventually dies down.

If there is a nonzero initial condition, y(-1), but the input x(n) = 0, the solution becomes $y(n) = y_1(n) = a y(-1) a^n u(n)$ which also dies down as $n \to \infty$.

Realization of digital filters

Given H(z), the system function, or h(n), the impulse response, the difference equation may be obtained. This difference equation could be implemented by *computer program, special purpose digital circuitry*, or *special programmable integrated circuit*. This direct evaluation of the difference equation is not the only possible realization of the digital filter. Alternative realizations of the digital filter are possible by breaking up the direct realization in some form.

Direct Form realization of IIR filters An important class of linear shift invariant systems can be characterized by the following rational system function (where X(z) is the input, Y(z) the output and we have taken $a_0 = 1$ in comparison with the earlier representation):

$$H(z) = \frac{Y(z)}{X(z)} = \frac{\frac{b+b}{z} + \frac{b}{z} - \frac{b}{z} + \frac{b}{z} - \frac{b}{z} + \frac{b}{z} - \frac{b}{z} + \frac{b}{z}$$

By cross multiplying and taking the inverse *z*-transform we get the difference equation

$$y(n) = -\sum_{k=1}^{N} a_k y(n-k) + \sum_{r=0}^{M} b_r x(n-r)$$

= -a₁ y(n-1) -a₂ y(n-2) -... -a_N y(n-N)
+ b₀ x(n) + b₁ x(n-1) + ... + b_M x(n-M)

To construct a filter structure we shall need three types of block diagram elements: a delay element, a multiplier and an adder, illustrated below:

We can construct a realization of the filter called the Direct Form I by starting with y(n) and generating all the delayed versions y(n-1), y(n-2), ..., y(n-N); similarly starting with x(n) and generating all the delayed versions x(n-1), x(n-2), ..., x(n-M). We then multiply the above terms by the respective coefficients and add them up. This is shown below (next page).

This is an N^{th} order system N being the order of the difference equation. There is no restriction as to whether M should be less than or greater than or equal to N. The total number of delay elements = (N+M). It is not in canonic form because it uses more than the minimum

possible number of delay elements. It is called "Direct Form" because the multipliers are the actual filter coefficients $\{a_1, a_2, ..., a_N, b_0, b_1, b_2, ..., b_M\}$.

The difference equation of this realization (or structure) continues to be

 $y(n) = -a_1 y(n-1) - a_2 y(n-2) - \dots - a_N y(n-N)$ N multiplications

 $+ b_0 x(n) + b_1 x(n-1) + ... + b_M x(n-M)$

(M+1) multiplications

and will be referred to as the Direct Form I difference equation. The total number of multiplications can be counted and is seen to be (N+M+1). We can also count and see that there are (N+M) additions. Finally, to calculate the value y(n) we need to store N past values of y(.), and M past values of x(.), that is, a total of (N+M) storage locations (storage for the present value of x(.) is not counted).



Rearrangement of Direct Form I The above diagram of Direct Form I, or the corresponding expression for H(z), is sometimes rearranged as below. This shows visually that the transfer

function H(z) is arranged as a cascade of an *all-zero system*, $H_2(z)$, followed by an *all-pole* system, $H_1(z)$:

 $H_2(z)$ $H_1(z)$

The overall block diagram then is shown thus:



The all-zero system is $H_2(z) = \frac{W(z)}{X(z)} = \left(\sum_{k=0}^{M} b_k z\right)^{-k}$ from which, by cross-multiplying and taking the inverse z-transform, we get the difference equation below: $W(z) = H_2(z) X(z) = X(z) \left(\sum_{k=0}^{M} b_k z^{-k}\right)^{-k}$

$$\bigvee_{k=0}^{k} \int_{m=0}^{k} w(n) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M)$$

$$X(z)$$
 $W(z)$ $W(z)$

The all-pole system is $H_I(z) = \frac{Y(z)}{W(z)} = \left| \frac{N^1}{1 + \sum_{k=1}^{N} a_k z^{-k}} \right|$ from which, by cross-multiplying

and taking the inverse z-transform, we get the difference equation as below:

$$Y(z) = H_{I}(z) \ W(z) = W(z) \left| \begin{array}{c} \Box 1 & | \\ 1 \\ \downarrow \\ 1 + \sum_{k=1}^{N} a_{k} z^{-k} \end{array} \right|$$

$$y(n) = w(n) - a_{I} \ y(n-1) - a_{2} \ y(n-2) - \dots - a_{N} \ y(n-N)$$



Even though it seems that there are two equations, one for w(n) and another for y(n), there is, in effect, only one since w(n) in the second equation is simply a short hand notation for the first equation and can be eliminated from the equation for y(n).

Overall, the Direct Form I has the following alternative appearance:



Derivation of Direct Form II The transfer function H(z) can be written as the product of the two transfer functions $H_1(z)$ and $H_2(z)$ as follows where we have reversed the sequence of the two blocks:

 $H_1(z)$ $H_2(z)$



$$P(z) = H_{1}(z) X(z) = \begin{vmatrix} 1 \\ \frac{1}{\left| 1 + \sum_{k=1}^{N} a_{k} z^{-k} \right|} X(z) \quad \text{and} \quad Y(z) = H_{2}(z) P(z) = \left(\sum_{k=0}^{M} b_{k} z^{-k} \right) P(z)$$

Cross-multiplying, taking the inverse *z*-transform of the above two and rearranging, we have

$$p(n) = x(n) - \sum_{k=1}^{N} a_k p(n-k)$$
, and $y(n) = \sum_{r=0}^{M} b_r p(n-r)$

The two equations are realized as below:



 $H_l(z)$ (All-pole System)
The two branches of delay elements in the middle of the above block diagram can be replaced by just one branch containing either N or M (whichever is larger) delay elements, resulting in the Direct Form II shown below:



(*Aside*) This diagram of the filter structure is properly called a *block diagram*. There is another representation called the *signal flow graph*. See below for signal flow graphs and transposed structures.



The number of delay elements = max $\{N, M\}$ – this is the minimum possible, hence called *a* **canonic form**. The multipliers are the actual coefficients from the difference equation. Hence this is also a direct form.

The numbers of multipliers and adders are also the minimum possible, but this does not mean that it is the best realization from other considerations like immunity to round off and quantization errors.

The difference equations are:

$$p(n) = x(n) - a_1 p(n-1) - a_2 p(n-2) - \dots - a_N p(n-N) , \text{ and} y(n) = b_0 p(n) + b_1 p(n-1) + \dots + b_M p(n-M)$$

The above equations show that in order to generate y(n) we need the present value of x(.) and N (or M or whichever is larger) past values of p(.). This requires N storage locations not counting the present value of x(.). We also see that the number of multiplications = N+M+1, and number of additions = N+M.

Comparing the difference equations of Direct Forms I and II:

- To compute y(n) in DF I we need the past N outputs, the present input, and the past M inputs.
- To compute y(n) in DF II we need the N (or M) values of p(n-k) for k = 1, 2, ..., N, and the present input.

This illustrates the concept of the state of a system.

Definition The **state of a system** is the minimum information required that along with the input allows the determination of the output.

From the above discussion the N (or M) values p(n-k), k = 1, 2, ..., N, make up the state of the system.

	Direct Form I	Direct Form II				
No. of multiplications	<i>N</i> + <i>M</i> +1	<i>N</i> + <i>M</i> +1				
No. of additions	N+M	N+M				
No. of storage locations	N+M	N				
State	y(n-k), k = 1 to N	p(n-k), k = 1 to N				
	x(n-k), k = 1 to M					

Comparison of Direct Forms I and II

Signal flow graph A signal flow graph representation of a difference equation is essentially the same as its *block diagram* representation, except for a few notational differences. The signal flow graph is a network of directed *branches* that connect at *nodes*. Associated with each node is a variable or node value. Each branch has an input signal and an output signal. In a linear signal flow graph the output of a branch is a linear transformation (such as a constant or a delay shown alongside the branch) of the input to the branch. Sometimes this branch transmittance is unity and not explicitly shown. By definition, the value at each node in a graph is the sum of the outputs of all the branches entering the node.

Source nodes have no entering branches and *sink nodes* have only entering branches. We repeat below the Direct Form II block diagram developed earlier.



The signal flow graph corresponding to the above structure is shown below. With regard to the "a" and "b" coefficients we have arbitrarily taken M = N, but M and N are independent of each other. In this figure the bottom-most right and left nodes are unnecessary.



Transposition (Flow graph reversal) leaves the overall system function between input and output unchanged and leads to a set of transposed system structures that provide some useful alternatives.

Transposition of a flow graph is accomplished by reversing the directions of all branches while keeping the branch transmittances as they were and reversing the roles of the input and output so that source nodes become sink nodes and vice versa. Clearly, if a signal flow graph configuration is transposed, the number of delay branches and the number of coefficients remain the same. Thus a canonic structure remains canonic in the process of transposition.

The result of applying the transposition theorem to the signal flow graph of the direct form II structure (and reorienting the diagram so that the input and output still appear on the left and right sides, respectively) is shown below. This will be referred to as the "*direct form II transposed*".



Comparison of the direct form II with its transposed version shows that

- The direct form II structure implements the poles first and then the zeros
- The transposed direct form II structure implements the zeros first and then the poles

These differences can become important in the presence of quantization in finite precision digital implementations or in the presence of noise in discrete-time analog implementations.

When the transposition theorem is applied to cascade or parallel structures, the individual second-order systems are replaced by transposed structures.

Example 1.11.1 Develop a canonic direct form realization of the transfer function

$$H(z) = \frac{6z^5 + 8z^3 - 4}{2z^5 + 6z^4 + 10z^3 + 8z}$$

Solution Write numerator and denominator as polynomials in negative powers of z with the leading term (a_0) in the denominator equal to 1

$$H(z) = \frac{z^{5}(6+8z^{-2}-4z^{-1})}{z^{5}(2+6z^{-1}+10z^{-2}+8z^{-4})} = \frac{(6+8z^{-2}-4z^{-5})}{(2+6z^{-1}+10z^{-2}+8z^{-4})}$$
$$= \frac{(6+8z^{-2}-4z^{-5})}{2(1+3z^{-1}+5z^{-2}+4z^{-4})} = \frac{3+4z^{-2}-2z^{-5}}{1+3z^{-1}+5z^{-2}+4z^{-4}}$$

Making the following comparison with the standard notation

$$H(z) = \frac{b + b z^{-1} + b z^{-2} + b z^{-3} + b z^{-4} + b z^{-3}}{1 + a z^{-1} + a z^{-2} + a z^{-3} + a z^{-3} + a z^{-4}} = \frac{3 + 4 z^{-2} - 2 z^{-3}}{1 + 3 z^{-1} + 5 z^{-2} + 4 z^{-4}}$$

we identify the following parameters:

$$b_0 = 3, b_1 = 0, b_2 = 4, b_3 = 0, b_4 = 0, b_5 = -2$$

 $a_1 = 3, a_2 = 5, a_3 = 0, a_4 = 4$



By inspection of the above structure we can write the difference equations

p(n) = x(n) - 3 p(n-1) - 5 p(n-2) - 4 p(n-4) and y(n) = 3 p(n) + 4 p(n-2) - 2 p(n-5) Example 1.11.2 [Ludeman, 5.1] A system is specified by its transfer function as

Realize the system in the following forms: (a) Direct Form I, and (b) Direct Form II. **Solution** We need to express H(z) as a ratio of polynomials in negative powers of z with the leading term (a_0) in the denominator equal to 1. Multiplying out the factors in the numerator and denominator and rearranging



By inspection of the above structure we can write the difference equations

$$p(n) = ?$$

 $y(n) = ?$

Biquadratic section When M = N = 2 we have the biquadratic section, an important building block. Thus H(z) is a ratio of two quadratics in z^{-1} given by

$$H(z) = \frac{b+bz^{-1}+bz^{-2}}{1+az^{-1}+az^{-2}}$$

The corresponding DF II structure is shown below with the difference equations.

Biquadratic section – DF II



 $p(n) = x(n) - a_1 p(n-1) - a_2 p(n-2)$ $y(n) = b_0 p(n) + b_1 p(n-1) + b_2 p(n-2)$

The biquadratic may also be implemented in the DF I with the difference equation:

DF I:
$$y(n) = b_0 x(n) + b_1 x(n-1) + b_2 x(n-2) - a_1 y(n-1) - a_2 y(n-2)$$

Shown below is an alternative realization of the biquadratic made possible by factoring out the coefficient b_0 . This is more convenient for amplitude scaling by tuning the value of b_0 .

$$H(z) = \frac{b (1 + b'z^{-1} + b'z^{-2})}{1 + a z^{-1} + a z^{-2}}$$

Alternative Biquadratic



As an exercise label the signal out of the left adder as w(n) and write down the difference equations.

Cascade realization of IIR filters Many different realizations exist depending on how we choose to write and rearrange the given transfer function. Two very important ways of decomposing the transfer function are the *cascade* and *parallel* decompositions.

In the cascade realization H(z) is broken up into a product of transfer functions H_1 , H_2 , ..., H_k , each a rational expression in z^{-1} as follows:

$$\frac{Y(z)}{X(z)} = \hat{H(z)} = \frac{H(z)}{k} H_{k-1}(z) \dots H_{2}(z) H_{1}(z)$$

so that Y(z) can be written as

$$Y(z) = H_k(z) H_{k-1}(z) \dots H_2(z) H_1(z) X(z)$$



Although H(z) could be broken up in many different ways, the most common cascade realization is to require each of the *k* product *H*'s to be a *biquadratic section*. In many cases the design procedure yields a product of biquadratic expressions so no further work is necessary to put H(z) in the required form. The product terms $H_i(z)$ could take various forms, depending on the actual problem. Some poss_hible forms are

$$H_{i}(z) = \frac{b_{0}}{1 + a_{1}z^{-1} + a_{2}z^{-2}}, \qquad H_{i}(z) = \frac{b_{0} + b_{1}z^{-1}}{1 + a_{2}z^{-1} + a_{2}z^{-2}},$$

$$H_{i}(z) = b_{0} + b_{1}z^{-1} + b_{2}z^{-2}, \qquad H_{i}(z) = \frac{b_{0} + b_{1}z^{-1} + b_{2}z^{-2}}{1 + a_{1}z^{-1}},$$

$$H_{i}(z) = \frac{b + b_{1}z^{-1} + b_{2}z^{-2}}{1 + a_{1}z^{-1} + a_{2}z^{-2}}, \qquad \text{(Biquadratic)},$$

$$H_{i}(z) = \frac{b_{0} + b_{1}z^{-1}}{1 + a_{1}z^{-1}}, \qquad \text{(Bilinear)},$$

Each of the $H_i(z)$ could then be realized using either the direct form I or II.

Different structures are obtained by changing the ordering of the sections and by changing the pole-zero pairings. In practice due to the finite word-length effects, each such cascade realization behaves differently from the others.

Example 1.11.3 Develop two different cascade canonic realizations of the following causal IIR transfer function

$$H(z) = \frac{z(0.3z - 0.5)(2z + 3.1)}{(z^2 + 2.1z - 3)(z + 0.67)}$$

Solution Write in terms of negative powers of z:

$$H(z) = \frac{z \, z^2 \, (0.3 - 0.5 z^{-1})(2 + 3.1 z^{-1})}{z^3 \, (1 + 2.1 z^{-1} - 3 z^{-2})(1 + 0.67 z^{-1})} = \frac{(0.3 - 0.5 z^{-1})(2 + 3.1 z^{-1})}{(1 + 2.1 z^{-1} - 3 z^{-2})(1 + 0.67 z^{-1})}$$

Two (of several) different cascade arrangements, based on how the factors are paired, are shown below in block diagram form. Note that the intermediate signal $y_1(n)$ is different from the intermediate signal $y_3(n)$.

Cascade – A

$$\begin{array}{c} x(n) & 0.3 - 0.5z^{-1} & y_1(n) \\ \hline 1 + 2.1z^{-1} - 3z^{-2} & 1 + 0.67z^{-1} \end{array} \xrightarrow{y_2(n)} y_2(n) = y(n)$$

Cascade - B

$$x(n) = \frac{2 + 3.1z^{-1}}{1 + 2.1z^{-1} - 3z^{-2}} = \frac{y_3(n)}{y_3(n)} = \frac{0.3 - 0.5z^{-1}}{1 + 0.67z^{-1}} = y(n)$$

Cascade *A* is shown below using the direct form II for each block separately:



Cascade *B* is shown below using the direct form II for each block separately:



Note in this example that if $H(z) = \frac{A(z) B(z)}{C(z)}$, then depending on the pole-zero pairings and the

sequence order of the blocks in the cascade we can have 4 different implementations (structures). These are equivalent from input to output though not at the intermediate point between the blocks. Moreover the quadratic $(z^2 + 2.1z - 3)$ has real roots and so can be split into two factors each of which can be combined with the other factor (z + 0.67) in the denominator. This results in more than the 4 structures shown here.



Example 1.11.4 [Ludeman, 5.1] A system is specified by its transfer function as

$$H(z) = \frac{(z-1)(z-2)(z+1)(z-1)|| - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1 || - 1$$

Realize the system as a cascade of two biquadratic sections.

Solution In the numerator all the zeros are real-valued, so any two factors can be combined to produce a quadratic; thus there are C(4, 2) ways of doing this. But in the denominator a pair of factors must be combined so as to produce real-valued coefficients, i.e., combine two factors with complex conjugate poles.

In general, whether in the numerator or denominator, a pair of factors associated with a pair of complex conjugate roots must be combined so that the resulting quadratic expression will have real coefficients.

Parallel realization of IIR filters The transfer function H(z) is written as a sum of transfer functions $H_1(z)$, $H_2(z)$, ..., $H_k(z)$ obtained by partial fraction expansion:

$$\frac{Y(z)}{X(z)} = H(z) = H(z) + H(z) + \dots + H_{k-1}(z) + H_k(z)$$

Thus

$$Y(z) = H(z) X(z) = [H_1(z) + H_2(z) + \dots + H_{k-1}(z) + H_k(z)] X(z)$$

= $H_1(z) X(z) + H_2(z) X(z) + \dots + H_{k-1}(z) X(z) + H_k(z) X(z)$

and is shown in block diagram fashion below. Note that the outputs $y_1(n)$, $y_2(n)$, ..., $y_k(n)$ are independent of each other; they are not coupled as in the case of the cascade structure.

Based on whether H(z)/z or H(z) is the starting point for partial fractions we have parallel forms I and II (S. K. Mitra). Both of these methods are illustrated below.



Parallel Form I (This corresponds to expanding H(z)/z instead of H(z) into partial fractions). This structure arises when the $H_i(z)$ are all selected to be of the form

$$H_{i}(z) = \frac{b_{0} + b_{1}z^{-1}}{1 + a_{1}z^{-1} + a_{2}z^{-2}} = \frac{b_{0}z^{2} + b_{1}z}{z^{2} + a_{1}z + a_{2}} = \frac{z(bz + b)}{z^{2} + a_{1}z + a_{2}}$$

where the quadratic terms are used for each pair of complex conjugate poles so that the coefficients of the corresponding $H_i(z)$ will be real. Each $H_i(z)$ is then realized as either a DF I or DF II. Special cases of $H_i(z)$ are

Parallel Form II (This corresponds to expanding H(z) directly into partial fractions). This structure arises when the $H_i(z)$ are all selected to be of the form

$$H_{i}(z) = \frac{b z}{1 + a z} + b z = \frac{b z + b}{1 + a z} = \frac{b z + b}{z^{2} + a z + a}$$

Example 4.11.5 Obtain the parallel realization for

$$H(z) = \frac{8z^{3} - 4z^{2} + 11z - 2}{(z - (1/4))(z^{2} - z + (1/2))}$$

Solution For the **Parallel Form I** we expand H(z)/z. Note that in the denominator the factor $(z^2 - z + (1/2))$ represents a complex conjugate pair of poles at $((1/2) \pm j(1/2))$.

$$\frac{H(z)}{z} = \frac{8z^3 - 4z^2 + 11z - 2}{z(z - (1/4))(z^2 - z + (1/2))} = \frac{A + B}{z(z - (1/4))} + \frac{Cz + D}{(z^2 - z + (1/2))}$$

$$A = \frac{8z^3 - 4z^2 + 11z - 2}{(z - (1/4))(z^2 - z + (1/2))} = \frac{-2}{(-1/4)(1/2)} = 16$$

$$B = \frac{8z^3 - 4z^2 + 11z - 2}{z(z^2 - z + (1/2))} = \dots = \frac{-8}{z}$$

To determine C and D_2

$$\frac{8z^{3}-4z^{2}+11z-2}{z(z-(1/4))(z^{2}-z+(1/2))}$$

$$=\frac{A(z-(1/4))(z^{2}-z+(1/2))+Bz(z^{2}-z+(1/2))+(Cz+D)z(z-(1/4))}{z(z-(1/4))(z^{2}-z+(1/2))}$$

Putting A = 16 and B = 8, and equating the numerators on both sides $8z^3 - 4z^2 + 11z - 2$

$$= 16(z - (1/4))(z^{2} - z + (1/2)) + 8z(z^{2} - z + (1/2)) + (Cz + D)z(z - (1/4))$$

ients of like powers of z on both sides we have

Equating the coefficients of like powers of z on both sides we have

$$z^{3}: 8 = 16 + 8 + C \rightarrow C = -16$$

$$z^{0}: -2 = 16 (-1/4) (1/2) \text{ which is an identity} \rightarrow \text{doesn't help}$$

$$z^{1}: 11 = 16 (1/2) + 16 (-1/4) (-1) + 8 (1/2) + D (-1/4) \rightarrow D = 20$$

have

Therefore we have

$$\frac{H(z)}{z} = \frac{16}{z} + \frac{8}{(z - (1/4))} + \frac{(-16) z + 20}{(z^2 - z + (1/2))}$$

$$H(z) = 16 + \frac{8z}{(z - (1/4))} + \frac{z(20 - 16 z)}{(z^2 - z + (1/2))}$$

$$H(z) = 16 + \frac{8}{1 - 0.25z^{-1}} + \frac{-16 + 20z^{-1}}{1 - z^{-1} + 0.5z^{-2}} = H_1(z) + H_2(z) + H_3(z)$$

The corresponding parallel form I diagram is shown below.



HW Identify other signals on the diagram as needed and write down the implementation equations in full.

> $p_2(n) = x(n) + 0.25 p_2(n-1),$ $y_2(n) = 8 p_2(n)$ $p_3(n) = x(n) + 1 p_3(n-1) - 0.5 p_3(n-2),$ $y_3(n) = -16 p_3(n) + 20 p_3(n-1)$ $y(n) = y_1(n) + y_2(n) + y_3(n)$

For the **Parallel Form II** we expand H(z) $H(z) = \frac{8z^3 - 4z^2 + 11z - 2}{(z - (1/4))(z^2 - z + (1/2))} = \frac{8z^3 - 4z^2 + 11z - 2}{z^3 - 1.25z^2 + 0.75z - 0.125}$

By long division we reduce the degree of the numerator by 1

Long Division
B
Denominator
$$z^3 - 1.25z^2 + 0.75z - 0.125$$

B
 $z^3 - 4z^2 + 11z - 2$
 $8z^3 - 4z^2 + 6z - 1$
 $6z^2 + 5z - 1$
 $-$ Remainder

$$H(z) = 8 + \frac{6z^2 + 5z - 1}{z^3 - 1.25z^2 + 0.75z - 0.125} = 8 + \frac{6z^2 + 5z - 1}{(z - (1/4))(z^2 - z + (1/2))}$$

The proper fraction part can now be expanded into partial fractions. Let $67^2 \pm 57 = 1$

$$\frac{6z^{2} + 5z - 1}{(z - (1/4))(z^{2} - z + (1/2))} = \frac{A}{(z - (1/4))} \frac{Bz + C}{(z^{2} - z + (1/2))}$$
$$= \frac{6z^{2} + 5z - 1}{(z^{2} - z + (1/2))} = \dots = \frac{m}{m} = 2$$
$$\dots$$

To determine *B* and *C* $6z^2 + 5z - 1$

1
$$A(z^2 - z + (1/2)) + (Bz + C)(z - (1/4))$$

$$\overline{(z-(1/4))(z^2-z+(1/2))}^{-1} \overline{(z-(1/4))(z^2-z+(1/2))}$$

Putting $A = 2$, and equating the numerators on both sides

$$6z^{2} + 5z - 1 = 2(z^{2} - z + (1/2)) + (Bz + C)(z - (1/4))$$

Equating the coefficients of like powers of z on both sides we have

$$z^{2}: \quad 6 = 2 + B \quad \to B = 4$$

$$z^{0}: \quad -1 = 2 (1/2) + C (-1/4) \quad \to C = 8$$

we have

$$z^{0}: \quad -1 = 2 (1/2) + C (-1/4) \quad \to C = 8$$

Therefore we have

$$H(z) = 8 + \frac{2}{(z - (1/4))} + \frac{4z + 8}{(z^2 - z + (1/2))}$$
$$H(z) = 8 + \frac{2z^{-1}}{1 - 0.25z^{-1}} + \frac{4z^{-1} + 8z^{-2}}{1 - z^{-1} + 0.5z^{-2}} = H_1(z) + H_2(z) + H_3(z)$$

The corresponding parallel form II diagram is shown below.



HW Identify other signals on the diagram as needed and write down the implementation equations in full.

Example 4.11.6 [Ludeman, 5.1] A system is specified by its transfer function as

Implement the parallel realization in constant, linear and biquadratic sections. **Solution** For the Parallel Form I expand H(z)/z and for the Parallel Form II expand H(z). In the denominator keep the complex conjugate roots together.

Realization of FIR filters A causal FIR filter is characterized by its transfer function H(z) given by

$$\frac{Y(z)}{X(z)} = H(z) = {}^{M}b \ z = b + b \ z^{-1} + \dots + b - M^{Z}M$$

or, by the corresponding difference equation

$$y(n) = \sum_{r=0}^{M} b_r x(n-r) = b_0 x(n) + b_1 x(n-1) + b_2 x(n-2) \dots + b_M x(n-M)$$

Note that some use the notation below with M coefficients instead of M + 1

$$y(n) = \sum_{r=0}^{M-1} b_r x(n-r) = b_0 x(n) + b_1 x(n-1) + b_2 x(n-2) \dots + b_{M-1} x(n-M+1)$$

We see that the output y(n) is a weighted sum of the present and past input values; it does not depend on past output values such as y(n-1), etc. The block diagram is shown below. It is also called a **tapped delay line** or a **transversal filter**.



It can be seen that this is the same as the direct form I or II shown earlier for the IIR filter, except that the coefficients a_1 through a_N are zero and $a_0 = 1$; further the delay elements are arranged in a horizontal line. As in earlier diagrammatic manipulation the multipliers can all feed into the rightmost adder and the remaining adders removed.

Other simplifications are possible based on the symmetry of the coefficients $\{b_r\}$, as we shall see in FIR filter design.

Cascade realization of FIR filters The simplest form occurs when the system function is factored in terms of quadratic expressions in z^{-1} as follows:

$$H(z) = \prod_{i=1}^{K} H_i(z) = \prod_{i=1}^{K} \left(b_{0i} + b_{1i} z^{-1} + b_{2i} z^{-2} \right)$$

Selecting the quadratic terms to correspond to the complex conjugate pairs of zeros of H(z) allows a realization in terms of real coefficients b_{0i} , b_{1i} and b_{2i} . Each quadratic could then be realized using the direct form (or alternative structures) as shown below.



Parallel realization of FIR filters These are based on interpolation formulas by Lagrange, Newton, or Hermite methods. In general, these realizations require more multiplications, additions and delays than the others.

The Lattice structure – Introduction

In order to introduce the Lattice structure consider the all-pole IIR filter (b_1 through b_M are zero)

$$H(z) = \frac{b_0}{\sum_{k=0}^{N} a_k z^{-k}} = \frac{b_0}{a_0 + a_1 z} + a_2 z^{-2} + \dots + a_N z^{-N}$$

We shall take $a_0 = 1$ and $b_0 = 1$ iso that
$$H(z) = \frac{1}{1 + \sum_{k=0}^{N} a_k z^{-k}} = \frac{1}{1 + a_1 z} + a_2 z^{-k} + \dots + a_N z^{-N}$$

We shall write the subscript k on the coefficients a as an index within parentheses: thus a_k becomes a(k); further we shall also incorporate the order N of the filter as a subscript on the coefficients a: thus for an Nth order filter we shall have the set of coefficients $a_N(k)$, with k = 1 to N. With this notational refinement the system function goes through the following steps

$$H(z) = \frac{1}{1 + \sum_{k=1}^{N} a(k) z^{-k}} = 1 + a(1) z^{-1} + \frac{1}{a(2) z^{-1}} + \dots + a(N) z^{-N}$$

$$H(z) = \frac{1}{1 + \sum_{k=1}^{N} a_{N}(k) z^{-k}} = \frac{1}{1 + a_{N}(1) z^{-1} + a_{N}(2) z^{-2} + \dots + a_{N}(N) \overline{z}^{N}} = \frac{1}{A_{N}(z)}$$

where $A_N(z)$ denotes the denominator polynomial.

Using the above notation we shall consider the implementation of a first order filter, that is N = 1. We then have

$$H_{I}(z) = \frac{1}{1 + \sum_{k=1}^{1} a_{1}(k) z^{-k}} = \frac{1}{1 + a_{1}(1)z^{-1}} = \frac{1}{A_{1}(z)}$$

Note that we have also awarded a subscript to the system function H, that is, H(z) is written $H_I(z)$ just to remind ourselves that we are only dealing with a first order system. The difference equation can be written down from

$$H_{I}(z) = \frac{Y(z)}{X(z)} = \frac{1}{1 + \sum_{k=1}^{1} q(k)z^{-k}} = \frac{1}{1 + a_{1}(1)z}$$
$$y(n) = x(n) - a_{I}(1)y(n-1) \quad or \quad x(n) = y(n) + a_{I}(1)y(n-1)$$

This difference equation is implemented by the following structure (which can be visually verified). Concerning the labeling of the diagram note that $a_I(1)$ is not a signal but a multiplier that multiplies the signal coming out of the delay element and before it reaches the adder (we have previously used a triangular symbol for the multiplier).



An embellished version of the above structure is shown below and is used as a building block in the lattice structure. This being a first order filter the relationship between the direct form coefficient $a_I(1)$ and the **lattice coefficient** K_I is obvious, that is, $a_I(1) = K_I$. The implementation equations, in terms of the lattice coefficient, are

 $f_0(n) = f_1(n) - K_1 g_0(n-1)$ [which amounts to $y(n) = x(n) - a_1(1)y(n-1)$] $g_1(n) = K_1 f_0(n) + g_0(n-1)$ [which amounts to $g_1(n) = K_1 y(n) + y(n-1)$]

At this point the need for the additional symbols f(.) and g(.) and the equation for $g_1(n)$ is not obvious, but they become more useful as we increase the order of the filter and the relationship between the coefficients of the two structures becomes more involved.



Consider the second order (all-pole) filter (N = 2) whose transfer function is

$$H_{2}(z) = \frac{Y(z)}{X(z)} = \frac{1}{1 + \sum_{k=1}^{2} a_{2}(k) z^{-k}} = \frac{1}{1 + a_{2}(1)z^{-1} + a_{2}(2)z^{-2}} = \frac{1}{A_{2}(z)}$$
$$y(n) = x(n) - a_{2}(1) y(n-1) - a_{2}(2) y(n-2)$$

The corresponding lattice structure is obtained by adding a second stage at the left end of the previous first order structure:



We can write an equation for y(n) in terms of the lattice coefficients K_1 and K_2 and the signal values x(n), y(n-1) and y(n-2):

$$y(n) = x(n) - K_1 (1+K_2) y(n-1) - K_2 y(n-2)$$

Comparing this equation with the direct form equation for y(n) given above we have the relationship between the direct form and lattice coefficients

$$a_2(1) = K_1 (1 + K_2)$$
 and $a_2(2) = K_2$ $a_2(0) = 1$

Note that we have thrown in a freebie in the form of $a_2(0) = 1$ for future (actually it corresponds to the leading term in the denominator polynomial, $A_2(z)$).

Digital Signal Processing – 2

II. Discrete Fourier series

Properties of discrete Fourier series, DFS representation of periodic sequences, Discrete Fourier transforms, Properties of DFT, Linear convolution of sequences using DFT, Computation of DFT, Relation between z-transform and DFS.

Contents:

Fourier analysis – Recapitulation Discrete Fourier series Properties of discrete Fourier series The discrete Fourier transform (DFT) Properties of DFT Filtering through DFT/FFT Picket-fence effect

Fourier analysis - Recapitulation

(1) The *Fourier series* (FS) of a continuous-time periodic signal, x(t), with fundamental period T_0 , is given by the *synthesis equation*

$$x(t) = \sum_{k=-\infty}^{\infty} X_k e^{j2\pi k F_0 t}$$

The Fourier coefficients, X_k , are given by the *analysis equation*

$$X_{k} = \frac{1}{T_{0}} \int_{T_{0}} x(t) e^{-j 2\pi k F_{0} t} dt$$

The fundamental frequency, F_0 (Hz), and the period, T_0 (seconds), are related by $F_0 = 1/T_0$.

(2) The *Fourier transform* (FT) of a continuous-time aperiodic signal, x(t), is given by the *analysis equation*

$$X(F) = \int_{-\infty}^{\infty} x(t) e^{-j2\pi Ft} dt \qquad \text{or} \qquad X(\Omega) = \int_{-\infty}^{\infty} x(t) e^{-j\Omega t} dt$$

Here Ω and *F* are analog frequencies, with $\Omega = 2\pi F$. The *inverse Fourier transform* is given by the *synthesis equation*

$$x(t) = \int_{-\infty}^{\infty} X(F) e^{j2\pi Ft} dF \qquad \text{or} \qquad x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(\Omega) e^{j\Omega t} d\Omega$$

(3) The Fourier series (DTFS/DFS) for a discrete-time periodic signal (periodic sequence),

x(n), with fundamental period N is given by the synthesis equation

$$x(n) = \sum_{k=0}^{N-1} X_k e^{j 2\pi k n / N}$$
, $0 \le n \le N-1$

The Fourier coefficient X_k are given by the *analysis equation*

$$X_k = \frac{1}{N} \sum_{n=0}^{N} x(n) e^{-j 2\pi k n/N}, \quad 0 \le k \le N-1$$

This is called the **discrete-time Fourier series** (**DTFS**) or just **discrete Fourier series** (**DFS**) for short. The sequence of coefficients, X_k , also is periodic with period N.

These two equations are derived below.

(Note that if the factor (1/N) is associated with x(n) rather than with X_k the two DFS equations are identical to the two DFT equations which are derived below in their *standard form*.)

(4) The Fourier transform (DTFT) of a finite energy discrete-time aperiodic signal

(aperiodic sequence), x(n), is given by the *analysis equation* (some write $X(e^{j\omega})$ instead of $X(\omega)$)

$$X(\omega) = \sum_{n = -\infty} x(n) \ e^{-j\omega n}$$

Certain convergence conditions apply to this analysis equation concerning the type of signal x(n). We shall call this the **discrete-time Fourier transform** (**DTFT**). Physically $X(\omega)$ represents the frequency content of the signal x(n). $X(\omega)$ is periodic with period 2π .

The inverse discrete-time Fourier transform is given by the synthesis equation

$$x(n) = \frac{1}{2\pi} \int_{2\pi}^{\omega} \omega \omega d \omega$$

The basic difference between the Fourier transform of a continuous-time signal and the Fourier transform of a discrete-time signal is this: For continuous time signals the Fourier transform, and hence the spectrum of the signal, have a frequency range $(-\infty, \infty)$; in contrast, for a discrete-time signal the frequency range of the DTFT is unique over the interval of $(-\pi, \pi)$ or, equivalently, $(0, 2\pi)$.

Since $X(\omega)$ is a periodic function of the frequency variable ω , it has a Fourier series expansion; in fact, the Fourier coefficients are the x(n) values.

Discrete Fourier series

Let x(n) be a real periodic discrete-time sequence of period *N*. If x(n) can be expressed as a weighted sum of complex exponentials, the response of a linear system to x(n) is easily determined by superposition. By analogy with the Fourier series representation of a periodic *continuous-time* signal, we can expect that we can obtain a similar representation for the periodic *discrete-time* sequence x(n). That is, we seek a representation for x(n) of the form

$$x(n) = \sum_{k} X_{k} e^{j k \omega_{0} n} \quad \text{for all } n$$

Here X_k are the Fourier coefficients and $\omega_0 = 2\pi/N$ is the fundamental (digital) frequency (as $\Omega_0 = 2\pi/T_0 = 2\pi F_0$ is in the case of continuous-time Fourier series). With $k\omega_0 = \omega_k = k \begin{vmatrix} 2\pi \\ -\pi \\ N \end{vmatrix}$, the

above is also written

$$\sum_{k} X_{k} e^{jk 2\pi n/N}$$

The function $e^{jk 2\pi n/N}$ is periodic in k with a periodicity of N and there are only N distinct functions in the set $\{e^{jk 2\pi n/N}\}$ corresponding to k = 0, 1, 2, ..., N-1. Thus the representation for x(n) contains only N terms (as opposed to infinitely many terms in the continuous-time case)

$$x(n) = \sum_{k = } X_k e^{jk 2\pi n/N}$$

The summation can be done over any *N* consecutive values of *k*, indicated by the summation index $k = \langle N \rangle$. For the most part, however, we shall consider the range $0 \le k \le N-1$, and the representation for x(n) is then written as

$$x(n) = \sum_{k=0}^{N-1} X_k e^{jk 2\pi n/N} \text{ for all } n$$

This equation is the **discrete-time Fourier series (DTFS)** or just **discrete Fourier series (DFS)** of the periodic sequence x(n) with coefficients X_k .

The coefficients X_k or X(k) are given by (we skip the algebra – S&S)

$$X_{k} = \frac{1}{N} \sum_{n=0}^{N-1} x(n) e^{-j2\pi k n/N}, \quad 0 \le k \le N-1 \quad < (B)$$

Note that the sequence of Fourier coefficients $\{X_k\}$ is periodic with period = N. That is, $X_k = X_{k+N}$. The coefficients can be interpreted to be a sequence of finite length, given by Eq. (B) for k = 0, 1, 2, ..., N-1 only and zero otherwise, or as a periodic sequence defined for all k by Eq. (B). Clearly both of these interpretations are equivalent.

Because the Fourier series for discrete-time periodic signals is a finite sum defined entirely by the values of the signal over one period, the series always converges. The Fourier series provides an *exact* alternative representation of the time signal, and issues such as *convergence* or the *Gibbs phenomena* do not arise. The periodic sequence X(k) has a convenient interpretation as samples on the unit circle, equally spaced in angle, of the z-transform of *one period* of x(n). Let $x_1(n)$ represent one period of x(n). That is,

$$x_{I}(n) = x(n), \ 0 \le n \le N-1$$

0, otherwise
Then $X_{I}(z) = \sum_{n=-\infty}^{\infty} x(n) z^{-n} = \sum_{n=0}^{N-1} x(n) z^{-n}$, and $X(k) = X^{1}(z) \Big|_{z=e^{j2\pi k/N}}$. This then corresponds to

sampling the *z*-transform $X_I(z)$ at *N* points equally spaced in angle around the unit circle, with the first such sample occurring at z = 1. (Note that the periodic sequence x(n) cannot be represented by its *z*-transform since there is no value of *z* for which the *z*-transform will converge. However, $x_I(n)$ does have a *z*-transform.)

Properties of discrete Fourier series

Properties of discrete Fourier series (DFS) for periodic sequences The following notation is used:

p = periodic; e = even; o = odd $W_N = e^{-j 2\pi/N}$ Re [.] = Real part of Im [.] = Imaginary part of |.| = Magnitude of Arg (.) = Argument of

The following properties should be noted.

	Sequence	DFS		Sequence	DFS
1	$x_p(n+m)$	$W_N^{-\kappa m} X_p(\mathbf{k})$	4	Re $[x_p(n)]$	$X_{pe}(k)$
2	$x_p^{*}(n)$	$X_{p}^{*}(-k)$	5	$j \operatorname{Im} [x_p(n)]$	$X_{po}(k)$
3	$x_p^*(-n)$	$X_{p}^{*}(k)$			

Example 2.3.1 Show that DFS $\{x_p(n+m)\} = W_N^{-km} X_p(k)$. Solution We have

DFS {
$$x_p(n+m)$$
 } = $\sum_{n=0}^{N-1} x_p (n+m) W_N^{kn}$

N = 1

Set $n+m = \lambda$ so that $n = \lambda - m$ and the limits n = 0 to N-1 become $\lambda = m$ to N-1+m. Then the RHS becomes

$$= {}^{N} \sum_{\substack{\lambda = m}}^{+} {}^{m} {}_{x_{p}}(\lambda) W_{N} W_{N}^{-km}$$

Since $x_p(\lambda)$ is periodic with period *N* the summation can be done over any interval of length *N*. Thus

DFS {
$$x_p(n+m)$$
 } = $\sum_{k=0}^{N-1} x^p (\lambda) W^{k\lambda} W^{-km} = W^{-km} \sum_{\lambda=0}^{N-1} x^{-\lambda} (\lambda) W^{k\lambda}$
= $\sum_{k=0}^{-km} W^{-km} = W^{-km} QED$

Example 2.3.2 Show that DFS $\{x_p^*(n)\} = X_p^*(-k)$.

Solution We have

DFS {
$$x^{*}(n)$$
 } = $\sum_{n=0}^{N-1} x^{*}(n) W^{kn} = \left\{ \left\{ \sum_{n=0}^{N-1} x^{*}(n) W^{kn} \right\}^{*} \right\}$
 $\left\{ \left\{ \sum_{n=0}^{N-1} x^{n} W_{N} \right\}^{*} = \left\{ X(-k) \right\}^{*} = X_{p}(-k)$
 $= \left\{ \sum_{n=0}^{N-1} x^{n} (1) - x^{n} \right\}^{*}$

Based on the properties above we can show that for a real periodic sequence $x_p(n)$, the following symmetry properties of the discrete Fourier series hold:

1. Re
$$[X_p(k)] = \text{Re } [X_p(-k)]$$

2. Im $[X_p(k)] = -\text{Im } [X_p(-k)]$
3. $|X_p(k)| = |X_p(-k)|$
4. arg $X_p(k) = -\arg X_p(-k)$

The discrete Fourier transform (DFT)

(*Omit*) *The discrete Fourier transform (DFT) derived from the Fourier series* The exponential Fourier series of a continuous time periodic signal x(t) with fundamental period T₀ is given by the synthesis equation

$$x(t) = \sum_{k=-\infty}^{\infty} X_k e^{j2\pi k F_0 t}$$
 < (1)

where the Fourier coefficients X_k are given by the analysis equation

$$X_{k} = \overline{T_{0}}_{T_{0}} x(t) e^{-j2\pi k F_{0} t} dt$$
 < (2)

with the fundamental frequency F_0 and the period T_0 related by F_0 (Hz) = $1/T_0$ (sec).

To obtain finite-sum approximations for the above two equations, consider the analog periodic signal x(t) shown in Figure and its sampled version $x_s(nT)$. Using $x_s(nT)$, we can approximate the integral for X_k by the sum

$$X_{k} = \frac{1}{T_{0}} \sum_{n=0}^{N-1} x_{s}(nT) e^{-j2\pi k F_{0}nT} T, \qquad k = 0, 1, ..., N-1$$
$$= \frac{1}{N} \sum_{n=0}^{N-1} x(n) e^{-j2\pi k n/N}, \qquad k = 0, 1, ..., N-1$$

where we used the relation $F_0T = 1/N$, and approximated dt (or Δt) by T, and have used the shorthand notation $x(n) = x_s(nT)$. (This procedure is similar to that used in a typical introduction to integral calculus).

A finite series approximation for x(t) is obtained by truncating the series for x(t) in equation (1) to *N* terms and substituting t = nT and $F_0 = 1/TN$. This will necessarily give the discrete sequence x(n) instead of the continuous function x(t):



The above two equations define the discrete Fourier transform (DFT) pair. A *slight* adjustment of the (1/N) factor is needed so as to conform to **standard usage**. The adjustment consists of moving the (1/N) factor from one equation to the other. Then the direct DFT of the time series $x_0, x_1, ..., x_{N-1}$ is defined as

$$X_{k} = \sum_{n=0}^{N-1} x_{n} e^{-j2\pi k n/N} , k = 0, 1, ..., N-1$$
 < (3)

And the inverse DFT is defined as $\sum_{n=1}^{N} \frac{1}{n}$

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k e^{j2\pi k n/N}, \qquad n = 0, 1, \dots, N-1 \qquad < (4)$$

It can be shown that substituting equation (3) into equation (4) produces an identity, so that the two equations are indeed mutually inverse operations and therefore constitute a valid transform pair.

(End of Omit)

The discrete Fourier transform as a discretized (sampled) version of the DTFT A finiteduration sequence x(n) of length N (the length N may have been achieved by zero-padding a sequence of shorter length) has a Fourier transform denoted $X(\omega)$ or $X(e^{i\omega})$,

$$X(\omega) = \sum_{n=0}^{N-1} x(n) e^{-j\omega n}, \qquad 0 \le \omega < 2\pi$$

 $X(\omega)$ is a continuous function of ω and has a period of 2π . If we take N samples of $X(\omega)$ at equally spaced frequencies { $\omega_k = 2\pi k/N, k = 0, 1, ..., N-1$ }, along the interval [0, 2π), the resulting samples are (Figure)

$$X(k) \equiv X \left(\begin{array}{c} 2\pi k \\ \square \\ N \end{array} \right) = \sum_{n \in \mathbb{N}} x(n) e^{-j 2\pi k n/N},$$

Since these frequency samples are obtained by evaluating the Fourier transform $X(\omega)$ at N equally spaced discrete frequencies, the above relation is called the **discrete Fourier transform** (**DFT**) of x(n). In other words, X(k) are discrete samples of the continuous $X(\omega)$.



The corresponding **inverse discrete Fourier transform (IDFT)** is given by $x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k) e^{j2\pi k n/N}, n = 0, 1, ..., N-1$

Example 2.4.1 Find the DFT of the **unit sample** $x(n) = \{1, 0, 0, 0\}$. (Aside. What is the DTFT of $x(n) = \{1, 0, 0, 0\}$?)



Solution The number of samples is N = 4. The DFT is given by

$$X(k) = \sum_{n=0}^{\infty} x(n) e^{-jk 2\pi n/N} , \qquad 0 \le k \le N - 1$$

$$= \sum_{n=0}^{4-1} x(n) e^{-jk2\pi n/4}, \quad 0 \le k \le 3$$

$$= \sum_{n=0}^{3} x(n) e^{-jk2\pi n/4}, \quad k = 0, 1, 2, 3$$

$$\sum_{n=0}^{3} x(n) e^{-jk2\pi n/4}, \quad k = 0, 1, 2, 3$$

$$\frac{k=0}{k=0} X(0) = \sum_{n=0}^{3} x(n) e^{-j02\pi n/4} = \sum_{n=0}^{3} x(n).1 = \sum_{n=0}^{3} x(n) = x(0) + x(1) + x(2) + x(3) = 1 + 0 + 0 + 0 = 1$$

$$k = 1 X(1) = \sum_{n=0}^{3} x(n) e^{-j1.2\pi n/4} = \sum_{n=0}^{3} x(n) e^{-j\pi n/2} = x(0) e^{-j0} = 1 \cdot 1 = 1$$

$$\frac{k=2}{k=3} X(2) = \sum_{n=0}^{3} x(n) e^{-j2.2\pi n/4} = \sum_{n=0}^{3} x(n) e^{-j\pi n} = x(0) e^{-j0} = 1 \cdot 1 = 1$$

$$k = 3 X(3) = \sum_{n=0}^{3} x(n) e^{-j3.2\pi n/4} = \sum_{n=0}^{3} x(n) e^{-j3\pi n/2} = x(0) e^{-j0} = 1 \cdot 1 = 1$$

The DFT is $X(k) = \{1, 1, 1, 1\}$ and contains all (four) frequency components. In this example X(k) is real-valued.

In MATLAB use fft(x) for the DFT. The magnitude and phase plots of X(k) and the program segment follow.



k=0:3; subplot(2, 1, 1), stem(k, Mag, 'bo'); %Two rows, one column, #1 xlabel ('k'), ylabel('Magnitude'); title ('4-point DFT of $\{1, 0, 0, 0\}$ ') grid; subplot(2, 1, 2), stem(k, Phase, 'bo'); %Two rows, one column, #2 xlabel ('k'), ylabel('Phase');

Matrix formulation To facilitate computation the DFT equations may be arranged as a matrixvector multiplication. We define the twiddle factor $W_N = e^{-j2\pi/N}$, which for N = 4 becomes

 $W_4 = e^{-j2\pi/4}$. The equations are rewritten using the twiddle factor

$$X(k) = \sum_{n=0}^{\infty} x(n) e^{-jk 2\pi n/4} = \sum_{n=0}^{\infty} x(n) W_4^{kn}, \qquad k = 0, 1, 2, 3$$

There are a total of 4 values of X(.) ranging from X(0) to X(3):

X(0)=x(0)W	$u^{0} + x(1)W^{0} + x(2)$	$2W_{4}^{0} + x(3)W_{4}^{0}$
X(1) = x(0)W	$\frac{1}{2} + x(1)W^{1} + x(2)$	$2)W^{\frac{3}{2}} + x(3)W^{\frac{3}{3}}$
X(2) = x(0)W	$\frac{1}{4} + x(1)W^{2} + x(1)$	$(2)W_{4}^{4} + x(3)W_{4}^{6}$
X(3) = x(0)W	$^{\ddot{0}} + x(1)W^{\ddot{3}} + x(2)$	$2)W^{\frac{3}{6}} + x(3)W^{\frac{3}{9}}$
	4 4	4 4

These equations can be put in matrix from:

$$\begin{bmatrix} X(0) \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} x(0) \\ x(2) \\ 1 & x(2) \end{bmatrix} = \begin{bmatrix} 1 & W_4 & W_4 & W_4 \\ W_4 & W_4 & W_6 \end{bmatrix} \begin{bmatrix} x(1) \\ x(2) \\ 1 & y \end{bmatrix}$$

This last form is perhaps the most convenient to perform the actual computations by plugging in the twiddle factors W_4^m and the signal values x(.). Note that

$$W_N^{k+(N/2)} = -W_N^k$$
 and $W_N^{mN+k} = W_N^k$

and

$$W_4^1 = e^{-j 2\pi/4} = -j, W_4^2 = (-j)^2 = -1, W_4^3 = (-j)^3 = j, \text{ etc.}$$

The above matrix form then can be written,

$$\begin{vmatrix} X(0) & | & 1 & 1 & 1 & 1 \\ X(1) & | & | & 1 & -j & -1 & j \\ X(2) & | & | & 1 & -1 & 1 & -1 \\ \end{bmatrix} \begin{vmatrix} 1 & -j & -1 & j \\ 0 & | & | & 1 \\ 0 & | & | & 1 \\ 0 & | & | & 1 \\ 1 & | & 1 \\ \end{bmatrix}$$

Example 2.4.2 Find the DFT of the "dc" sequence $x(n) = \{1, 1, 1, 1\}$. (Aside. What is the DTFT of $x(n) = \{1, 1, 1, 1\}$? Give 4 samples of $X(\omega)$ at intervals of $2\pi/4$ starting at $\omega = 0$.) (Compare Proakis, 3^{rd} Ed., Ex. 5.1.2)



Solution The number of samples is N = 4. The DFT is given by

$$X(k) = \sum_{n=0}^{3} x(n) e^{-jk2\pi n/N}, \qquad 0 \le k \le N-1$$
$$= \sum_{n=0}^{3} x(n) e^{-jk2\pi n/4}, \qquad k = 0, 1, 2, 3$$

$$X(k) == \sum_{n=0}^{3} x(n) \ e^{-jk \ 2\pi n/4}, \ k = 0, \ 1, \ 2, \ 3$$

$$k = 0 \qquad X(0) = \sum_{n=0}^{3} x(n) \ e^{-j0 \ 2\pi n/4} = \sum_{n=0}^{3} x(n) \ 1 = \sum_{n=0}^{3} x(n) = x(0) + x(1) + x(2) + x(3) = 1 + 1 + 1 + 1 = 4$$

$$k = 1 \qquad X(1) = \sum_{n=0}^{3} x(n) \ e^{-j1 \ 2\pi n/4} = \sum_{n=0}^{3} x(n) \ e^{-j\pi n/2} = \sum_{n=0}^{3} 1 \ (e^{-j\pi/2})^n = \sum_{n=0}^{3} (-j)^n = 1 - j + (-j)^2 + (-j)^3 = 1 - j + j^2 - j^3 = 1 - j - 1 + j = 0$$

$$k = 2 \qquad X(2) = \sum_{n=0}^{3} x(n) \ e^{-j2 \ 2\pi n/4} = \sum_{n=0}^{3} 1 \ (e^{-j\pi})^n = 1 - 1 + 1 - 1 = 0$$

$$k = 3 \qquad X(3) = \sum_{n=0}^{3} x(n) \ e^{-j3 \ 2\pi n/4} = \sum_{n=0}^{3} 1 \ e^{-j3\pi n/2} = 1 + j - 1 - j = 0$$

The DFT is $X(k) = \{4, 0, 0, 0\}$ and contains only the "dc" component and no other. Here again X(k) is real-valued.

Example 2.4.3 Find the DFT of the sequence $x(n) = \{1, 0, 0, 1\}$





Solution The number of samples is N = 4. The DFT is given by

$$X(k) = \sum_{n=0}^{3} x(n) e^{-jk2\pi n/N}, \qquad 0 \le k \le N-1$$
$$= \sum_{n=0}^{3} x(n) e^{-jk2\pi n/4}, \qquad k = 0, 1, 2, 3$$

$$X(k) = \sum_{n=0}^{3} x(n) e^{-jk 2\pi n/4}, k = 0, 1, 2, 3$$

$$k = 0$$

$$X(0) = \sum_{n=0}^{3} x(n) e^{-j0.2\pi n/4} = \sum_{n=0}^{3} x(n).1 = \sum_{n=0}^{3} x(n)$$

$$= x(0) + x(1) + x(2) + x(3) = 1 + 0 + 0 + 1 = 2$$

$$k = 1$$

$$X(1) = \sum_{n=0}^{3} x(n) e^{-j1.2\pi n/4} = \sum_{n=0}^{3} x(n) e^{-j\pi n/2} = \sum_{n=0}^{3} x(n)(e^{-j\pi/2})^n = \sum_{n=0}^{3} x(n)(-j)^n$$

$$= 1 \cdot 1 + 0 + 0 + 1 \cdot (-j)^3 = 1 + j = \sqrt{2} e^{j\pi/4} = \sqrt{2} \angle \pi/4$$

$$k = 2$$

$$X(2) = \sum_{n=0}^{3} x(n) e^{-j2.2\pi n/4} = \sum_{n=0}^{3} x(n) e^{-j\pi n} = x(0) e^{-j0} + x(3) e^{-j\pi 3}$$

$$= 1 \cdot 1 + 1 \cdot (-1) = 0$$

$$k = 3$$

$$X(3) = \sum_{n=0}^{3} x(n) e^{-j3.2\pi n/4} = \sum_{n=0}^{3} x(n) e^{-j3\pi n/2} = x(0) e^{-j0} + x(3) e^{-j3\pi 3/2}$$

$$= 1 \cdot 1 + 1 \cdot (-j) = 1 - j = \sqrt{2} e^{-j\pi/4} = \sqrt{2} \angle -\pi/4$$

The DFT is $X(k) = \{2, \sqrt{2}e^{j\pi/4}, 0, \sqrt{2}e^{-j\pi/4}\}$. In general X(k) is complex-valued and has a magnitude and a phase. See figure.

In MATLAB use fft(x) for the DFT and ifft(X) for the IDFT. The magnitude and phase plots and the program segment follow.

xlabel ('k'), ylabel('Magnitude'); title ('4-point DFT of \{1, 0, 0, 1\}') grid; subplot(2, 1, 2), stem(k, Phase, 'bo'); %Two rows, one column, #2 xlabel ('k'), ylabel('Phase');



:

.



The frequency components {X(k), k = 0, 1, 2, 3} are "harmonics". The spacing between successive components, denoted by F_0 , is the **resolution of the DFT** and is given by $F_0 = F_s /N$: it is the sampling interval in the *frequency domain*. It is determined by the sampling frequency F_s or sampling interval T in the *time domain* ($F_s = 1/T$) and the number of samples, N. Note that N is the number of time domain samples as well as the number of frequency domain samples.

Record length, sampling time and frequency resolution of the DFT We shall use Example 3 above to illustrate. Suppose that the sampling frequency is 8000 Hz, then the sampling time is $T = (1/8000) \sec = 125 \ \mu \sec c$. The time domain samples x(n) are spaced 125 $\mu \sec c$ apart. In the frequency domain the DFT values, X(k), are spaced (8000/4) = 2000 Hz apart. The DFT spectrum is re-plotted below with these parameters.



The terms introduced in this example and their interrelationships are summarized below:

Record length (one period) = T_0 seconds = N samples Sampling interval = T seconds = $1/F_s$ Sampling interval = T sec = (T_0/N) seconds

Sampling frequency (one period) = F_s Hz Frequency resolution of the DFT = F_0 Hz = $1/T_0$ Frequency resolution of the DFT = F_0 Hz = (F_s/N) Hz The situation in the frequency domain is shown below.



Example 2.4.4 Find the inverse discrete Fourier transform of $X(k) = \{3, (2+j), 1, (2-j)\}$. Solution The number of samples is N = 4. The IDFT is given by the synthesis equation

$$\begin{aligned} x(n) &= \frac{1}{N} \sum_{k=0}^{\infty} X(k) e^{j2\pi k n/N}, & n = 0, 1, \dots, N-1 \\ &= \frac{1}{4} \sum_{k=0}^{3} X(k) e^{j2\pi k n/4}, \\ &= \frac{1}{4} \sum_{k=0}^{3} X(k) e^{jkn\pi/2}, & 0 \le n \le 3 \end{aligned}$$

The calculations for $\{x(n), n = 0 \text{ to } 3\}$ are shown in table below.

$$x(n) = \frac{1}{4} \sum_{k=0}^{3} X(k) e^{jk n\pi/2}$$

$$\begin{split} n &= 0 \\ x(0) &= \frac{1}{4} \sum_{k=0}^{3} X(k) e^{jk \, 0\pi/2} = \frac{1}{4} \sum_{k=0}^{3} X(k) \\ &= (1/4) \left\{ X(0) + X(1) + X(2) + X(3) \right\} = (1/4) \left\{ 3 + 2 + j + 1 + 2 - j \right\} = 2 \\ n &= 1 \\ x(1) &= (1/4) \sum_{k=0}^{3} X(k) e^{jk \, 1\pi/2} = (1/4) \sum_{k=0}^{3} X(k) \left(e^{j\pi/2} \right)^{k} = (1/4) \sum_{k=0}^{3} X(k) (j)^{k} \\ &= (1/4) \left\{ X(0) (j)^{0} + X(1) (j)^{1} + X(2) (j)^{2} + X(3) (j)^{3} \right\} \\ &= (1/4) \left\{ 3 \cdot 1 + (2 + j) \cdot j + 1 \cdot (-1) + (2 - j) \cdot (-j) \right\} = 0 \\ n &= 2 \\ x(2) &= (1/4) \sum_{k=0}^{3} X(k) e^{jk \, 2\pi/2} = (1/4) \sum_{k=0}^{3} X(k) e^{jk \, \pi} = (1/4) \sum_{k=0}^{3} X(k) (-1)^{k} \\ &= (1/4) \left\{ X(0) (1) + X(1) (-1) + X(2) (1) + X(3) (-1) \right\} \\ &= (1/4) \left\{ X(0) - X(1) + X(2) - X(3) \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) + 1 \cdot (-1) - (2 - j) \right\} = 0 \\ n &= 3 \\ x(3) &= (1/4) \sum_{k=0}^{3} X(k) e^{jk \, 3\pi/2} = (1/4) \sum_{k=0}^{3} X(k) \left(e^{j3\pi/2} \right)^{k} = (1/4) \sum_{k=0}^{3} X(k) (-j)^{k} \\ &= (1/4) \left\{ 3 \cdot (1 + (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-1) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-j) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3 \cdot (2 + j) \cdot (-j) + 1 \cdot (-j) + (2 - j) \cdot j \right\} \\ &= (1/4) \left\{ 3$$

Thus $x(n) = \{2, 0, 0, 1\}.$

In MATLAB use ifft(X) for the IDFT. The magnitude and phase plots of x(n) and the program segment follow.

$$\begin{split} X &= [3, (2+j), 1, (2-j)]; x = ifft(X); Mag = abs(x); Phase = angle(x); \\ n &= 0:3; \\ subplot(2, 1, 1), stem(n, Mag, 'bo'); %Two rows, one column, #1 \\ xlabel ('n'), ylabel('Magnitude'); \\ title ('4-point IDFT of \{3, (2+j), 1, (2-j)\}') \\ grid; \\ subplot(2, 1, 2), stem(n, Phase, 'bo'); %Two rows, one column, #2 \\ xlabel ('n'), ylabel('Phase'); \end{split}$$



Matrix formulation Here again to facilitate computation the IDFT equations may be arranged as a matrix veqtor multiplication.

$$\begin{aligned} x(n) &= \sum_{k=0}^{n} X(k) e^{jk2\pi n/4} = \frac{1}{4} \sum_{n=0}^{3} X(k) \left(\overset{*}{W} \right)^{k_{n}}, \qquad n = 0, 1, 2, 3 \\ \text{There are a total of 4 values of } x(.) \text{ ranging from } x(0) \text{ to } x(3): \\ \begin{bmatrix} x(0) \end{bmatrix} & \left[\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 0 \end{bmatrix} \end{bmatrix} \begin{bmatrix} x(0) \end{bmatrix} \\ \begin{vmatrix} x(1) \\ \end{vmatrix} = \frac{1 \begin{vmatrix} 1 \\ 1 \end{vmatrix} \begin{pmatrix} 1 & 0 \\ 4 \end{vmatrix} \begin{pmatrix} W \\ 4 \end{vmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{vmatrix} W \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{vmatrix} X(1) \\ 1 \end{pmatrix} \\ \begin{vmatrix} x(2) \\ 1 \end{pmatrix} \begin{vmatrix} 4 \\ 1 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \begin{vmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 2 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \begin{pmatrix} W \\ 4 \end{pmatrix} \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \end{vmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \end{vmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \end{vmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \end{vmatrix} \\ \end{vmatrix} \\ \end{vmatrix} \\ \begin{vmatrix} 2 \\ 4 \end{pmatrix} \\ \end{vmatrix} \\ \end{vmatrix} \\ \end{vmatrix} \\ \end{vmatrix}$$

.

This last form is perhaps the most convenient to perform the actual computations by plugging in the twiddle factors $(W_4^*)^m$ and the values X(.). The above matrix form then can be written,

:

$$\begin{bmatrix} x(0) \\ x(1) \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0$$


Example 2.4.6 [N not an integral power of 2] Using MATLAB find the 6-point IDFT of

$$X(k) = \{6, 0, 0, 0, 0, 0\}$$

Solution

$$\begin{split} &X = [6, 0, 0, 0, 0, 0]; \ x = ifft(X); \ Mag = abs(x); \ Phase = angle(x); \\ &n = 0:5; \\ &subplot(2, 1, 1), \ stem(n, Mag, 'bo'); \ \% \ Two \ rows, \ one \ column, \ \#1 \\ &xlabel \ ('n'), \ ylabel('Magnitude'); \\ &title \ ('6-point \ IDFT \ of \ \{6, 0, 0, 0, 0, 0, 0\}') \\ &grid; \\ &subplot(2, 1, 2), \ stem(n, \ Phase, \ 'bo'); \ \% \ Two \ rows, \ one \ column, \ \#2 \\ &xlabel \ ('n'), \ ylabel('Phase'); \end{split}$$



:

Example 2.4.7 Consider a sequence $x(n) = \{2, -1, 1, 1\}$ and the sampling time T = 0.5 sec. Compute its DFT and compare it with its DTFT.

Solution The record length of the sequence is $T_0 = 4T = 2$ sec.

The DFT is a sequence of 4 values given by

$$X(k) = \sum_{n=0}^{\infty} x(n) e^{-jk 2\pi n/4}, \qquad k = 0, 1, 2, 3$$

The periodicity of X(k) is 4. The frequency resolution of the DFT is $1/T_0 = 0.5$ Hz.

The DTFT, $X(\omega)$, is a continuous function of ω

$$X(\omega) = \sum_{n = -\infty} x(n) e^{-j\omega n} = \sum_{n = 0}^{\infty} x(n) e^{-j\omega n} = 2 e^{j0} - 1 e^{-j\omega} + 1 e^{-j2\omega} + 1 e^{-j3\omega}$$
$$= 2 - e^{-j\omega} + e^{-j2\omega} + e^{-j3\omega}$$

The periodicity of $X(\omega)$, in terms of ω , is 2π . In terms of the Hertz frequency the periodicity is the sampling frequency = $F_s = 1/T = 2$ Hz.

The DFT is a sampled version of the DTFT, sampled at 4 points along the frequency axis spaced 0.5 Hz apart.

You should evaluate completely both X(k) (a set of 4 numbers) and $X(\omega)$ (magnitude and phase). Note that $X(\omega)$ may be evaluated directly at $\omega = 0$, $\pi/2$, π , and $3\pi/2$ by plugging in the ω values into the expression given above; these are then the DFT numbers as well. The MATLAB solutions are given below.

In MATLAB: The magnitude and phase plots of the DFT can be generated by the following segment:

 $\begin{aligned} x &= [2, -1, 1, 1]; X = fft(x); Mag = abs(X); Phase = angle(X); \\ k &= 0:3; \\ subplot(2, 1, 1), stem(k, Mag, 'bo'); %Two rows, one column, #1 \\ xlabel ('k'), ylabel('Magnitude'); \\ title ('4-point DFT of \{2, -1, 1, 1\}') \\ grid; \\ subplot(2, 1, 2), stem(k, Phase, 'bo'); %Two rows, one column, #2 \\ xlabel ('k'), ylabel('Phase'); \end{aligned}$

The MATLAB solution:

X = 3, (1 + j2), 3, (1 - j2) Mag = 3, 2.2361, 3, 2.2361 Phase = 0, 1.1071, 0, -1.1071



The magnitude and phase plots of the DTFT can be generated by the following segment:

b = [2, -1, 1, 1]; % Numerator coefficientsa = [1]; % Denominator coefficientw = 0: pi/256: 2*pi; % A total of 512 points[h] = freqz(b, a, w);subplot(2, 1, 1), plot(w/pi, abs(h));xlabel('Frequency, Hz'), ylabel('Magnitude'); grid $title ('4-point DTFT of \{2, -1, 1, 1\}')$ subplot(2, 1, 2), plot(w/pi, angle(h));xlabel('Frequency, Hz'), ylabel('Phase - Radians'); grid



Example 2.4.8 Compute the discrete Fourier transform of the following finite length sequences considered to be of length *N*.

1) $x(n) = \delta(n+n_0), \quad 0 < n_0 < N$ 2) $x(n) = a^n, \quad 0 < a < 1$

Solution See Ramesh Babu 3.16.

Note that $x(n) = \delta(n+n_0)$ would be zero everywhere except at $n = -n_0$ which is not in the range [0, *N*). So make it $\delta(n-n_0)$.

Example 2.4.9 [2008] Compute the *N*-point DFT *X*(*k*) of the sequence

 $x(n) = \cos (2\pi n/N), \quad 0 \le n \le N-1$ for $0 \le k \le N-1$. **Solution** Express $\cos (2\pi n/N)$ as $(e^{j2\pi n/N} + e^{-j2\pi n/N})/2$.

Solution Express $\cos(2\pi n/n)$ as $(e^{2\pi} + e^{2\pi})/2$.

Example 2.4.10 Obtain the 7-point DFT of the sequence $x(n) = \{1, 2, 3, 4, 3, 2, 1\}$ by taking 7 samples of its DTFT uniformly spaced over the interval $0 \le \omega \le 2\pi$.

Solution The sampling interval in the frequency domain is $2\pi/7$. From Example 4 we have $X(e^{j\omega})$ or $X(\omega) = 1 + 2e^{-j\omega 1} + 3e^{-j\omega 2} + 4e^{-j\omega 3} + 3e^{-j\omega 4} + 2e^{-j\omega 5} + 1e^{-j\omega 6}$

 $= (2\cos 3\omega + 4\cos 2\omega + 6\cos \omega + 4) e^{-j3\omega}$

The DFT, X(k), is given by replacing ω with $k(2\pi/7)$ where k is an index ranging from 0 to 6: DFT = $X(\omega) = 2\pi \sum_{k/7} X(2\pi k/7)$, k = 0 to 6

This is denoted XkfromDTFT in the MATLAB segment below.

MATLAB:

w = 0: 2*pi/7: 2*pi-0.001XkfromDTFT = (4+6*cos(w)+4*cos(2*w)+2*cos(3*w)) .* exp(-j*3*w)

MATLAB solution:

 $\begin{aligned} Xk from DTFT &= [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), \\ (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)] \end{aligned}$

This is the 7-point DFT obtained by sampling the DTFT at 7 points uniformly spaced in $(0, 2\pi)$. It should be the same as the DFT directly obtained, for instance, by using the *fft* function in MATLAB:

MATLAB:

xn = [1 2 3 4 3 2 1]Xkusingfft = fft(xn)

MATLAB solution:

Xkusingfft = [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)]

It can be seen that " \downarrow kfromDTFT" = " \downarrow kusingfft".

Example 2.4.11 Obtain the 7-point inverse DTFT x(n) by finding the 7-point inverse DFT of X(k):

 $X (2\pi k / 7) = [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)]$

MATLAB:

 $\begin{aligned} Xk &= [16, (-4.5489 - 2.1906i), (0.1920 + 0.2408i), (-0.1431 - 0.6270i), (-0.1431 + 0.6270i), (0.1920 - 0.2408i), (-4.5489 + 2.1906i)] \\ xn &= ifft(Xk) \end{aligned}$

MATLAB solution:

 $xn = [1.0000 \quad 2.0000 \quad 3.0000 \quad 4.0000 \quad 3.0000 \quad 2.0000 \quad 1.0000]$

This is the original sequence we started with in Example 4.

Example 2.4.12 What will be the resulting time sequence if the DTFT of the 7-*point* sequence is sampled at 6 (or fewer) uniformly spaced points in $(0, 2\pi)$ and its inverse DFT is obtained? **Solution** The sampling interval in the frequency domain now is $2\pi/6$. From Example 4 we have

$$X(e^{j\omega}) \text{ or } X(\omega) = 1 + 2e^{-j\omega 1} + 3e^{-j\omega 2} + 4e^{-j\omega 3} + 3e^{-j\omega 4} + 2e^{-j\omega 5} + 1e^{-j\omega 6}$$
$$= (2\cos 3\omega + 4\cos 2\omega + 6\cos \omega + 4)e^{-j3\omega}$$

The DFT then is given by

DFT = $X(\omega)|_{\omega = \mathfrak{X} k/6} = X(2\pi k/6), \qquad k = 0 \text{ to } 5$

This is denoted Xk6point in the MATLAB segment below.

MATLAB:

MATLAB solution:

Xk6point = [16, (-3.0000 - 0.0000i), (1.0000 + 0.0000i), 0, (1.0000 + 0.0000i), (-3.0000 - 0.0000i)]

This is the 6-point DFT obtained by sampling the DTFT at 6 points uniformly spaced in $(0, 2\pi)$.

Example 2.4.13 Obtain the 6-point inverse DTFT x(n) by finding the 6-point inverse DFT of Xk6point:

 $X (2\pi k / 6) = [16, (-3.0000 - 0.0000i), (1.0000 + 0.0000i), 0, (1.0000 + 0.0000i), (-3.0000 - 0.0000i)]$

MATLAB:

Xk6point = [16, (-3.0000 - 0.0000i), (1.0000 + 0.0000i), 0, (1.0000 + 0.0000i), (-3.0000 - 0.0000i)] xn = ifft(Xk6point)

MATLAB solution:

 $xn = \begin{bmatrix} 2 & 2 & 3 & 4 & 3 & 2 \end{bmatrix}$

Comparing with the original 7-point sequence, $xn = \begin{bmatrix} 1 & 2 & 3 & 4 & 3 & 2 & 1 \end{bmatrix}$, we see the consequence of *under-sampling* the continuous- ω function $X(\omega)$: the corresponding time domain sequence x(n) is said to suffer *time-domain aliasing*. This is similar to the situation that occurs when a continuous-time function x(t) is under-sampled: the corresponding frequency domain function $X_s(\omega)$ contains *frequency-domain aliasing*.

Example 2.4.14 What will be the resulting time sequence if the DTFT of the 7-*point* sequence is sampled at 8 (or more) uniformly spaced points in $(0, 2\pi)$ and its inverse DFT is obtained? **Solution** The sampling interval in the frequency domain now is $2\pi/8$. From Example 4 we have

$$X(e^{j\omega}) \text{ or } X(\omega) = 1 + 2e^{-j\omega 1} + 3e^{-j\omega 2} + 4e^{-j\omega 3} + 3e^{-j\omega 4} + 2e^{-j\omega 5} + 1e^{-j\omega 6}$$
$$= (2 \cos 3\omega + 4 \cos 2\omega + 6 \cos \omega + 4) e^{-j3\omega}$$

The DFT then is given by

DFT =
$$X(\omega)|_{\omega = 2\pi k/8} = X(2\pi k/8), \qquad k = 0 \text{ to } 7$$

This is denoted Xk8point in the MATLAB segment below.

MATLAB:

w = 0: 2*pi/8: 2*pi-0.001 Xk8point = (4+6*cos(w)+4*cos(2*w)+2*cos(3*w)) .* exp(-j*3*w)

MATLAB solution:

Xk8point = [16, (-4.8284 - 4.8284i), (0.0000 - 0.0000i), (0.8284 - 0.8284i), 0, (0.8284 + 0.8284i), (0.0000 - 0.0000i), (-4.8284 + 4.8284i)]

This is the 8-point DFT obtained by sampling the DTFT at 8 points uniformly spaced in $(0, 2\pi)$.

Example 2.4.15 Obtain the 8-point inverse DTFT x(n) by finding the 8-point inverse DFT of Xk8point:

 $X (2\pi k / 8) = [16, (-4.8284 - 4.8284i), (0.0000 - 0.0000i), (0.8284 - 0.8284i), 0, (0.8284 + 0.8284i), (0.0000 - 0.0000i), (-4.8284 + 4.8284i)]$

MATLAB:

Xk8point = [16, (-4.8284 - 4.8284i), (0.0000 - 0.0000i), (0.8284 - 0.8284i), 0, (0.8284 + 0.8284i), (0.0000 - 0.0000i), (-4.8284 + 4.8284i)] xn = ifft(Xk8point)

MATLAB solution:

 $xn = [1 \ 2 \ 3 \ 4 \ 3 \ 2 \ 1 \ 0]$

We see that the original 7-point sequence has been preserved with an *appended zero*. The original sequence and the *zero-padded sequence* (with any number of zeros) have the same DTFT. This is a case of *over-sampling* the continuous- ω function $X(\omega)$: there is no *time-domain aliasing*. This is similar to the situation that occurs when a continuous-time function x(t) is over-sampled: the corresponding frequency domain function $X_s(\omega)$ is free from *frequency-domain aliasing*.

% Sketch of sequences n = 0:1:6; xn = [1, 2, 3, 4, 3, 2, 1]; subplot (3, 1, 1), stem(n, xn) xlabel('n'), ylabel('x(n)-7point'); grid % n = 0:1:5; xn = [2 2 3 4 3 2]; subplot (3, 1, 2), stem(n, xn)
xlabel('n'), ylabel('x(n)-6point'); grid

%

n = 0:1:7; xn = [1, 2, 3, 4, 3, 2, 1, 0]; subplot (3, 1, 3), stem(n, xn) xlabel('n'), ylabel('x(n)-8point'); grid



Properties of DFT

The properties of the DFT (for finite duration sequences) are essentially similar to those of the DFS for periodic sequences and result from the *implied* periodicity in the DFT representation.

(1) **Periodicity** If x(n) and X(k) are an *N*-point DFT pair, then X(k) (and x(n)) is periodic with a periodicity of *N*. That is

$$X(k+N) = X(k)$$
 for all k

This can be proved by replacing k by k+N in the defining equation for X(k).

(2) Linearity For two sequences $x_1(n)$ and $x_2(n)$ defined on [0, N-1], if $x_3(n) = a x_1(n) + b x_2(n)$ then

DFT {
$$x_3(n)$$
} = DFT { $a x_1(n) + b x_2(n)$ } = a DFT { $x_1(n)$ } + b DFT { $x_2(n)$ }

If one of the two sequences $x_1(n)$ and $x_2(n)$ is shorter than the other then the shorter one must be padded with zeros to make both sequences of the same length.

(3) Circular shift or circular translation of a sequence (The sequence wraps around). The circular shift, by an amount n_0 to the right, of the sequence x(n) defined on [0, N-1] is denoted by $x((n-n_0)_{mod N})$ or $x((n-n_0))_N$. For example, if x(n) is

$$x(n) = \{x(0), x(1), x(2), ..., x(N-3), x(N-2), x(N-1)\}$$

then $x((n-2))_N$ is given by

$$x((n-2))_N = \{x(N-2), x(N-1), x(0), x(1), x(2), \dots, x(N-3)\}$$

The operation can be thought of as wrapping the part that falls outside the region of interest around to the front of the sequence, or equivalently, just a straight (linear) translation of its *periodic extension*.

Example 2.5.1 Given $x(n) = \{1, 2, 2, 0\}$. Here N = 4. The circular shift of x(n) by one unit to the right is $x((n-1))_4$, and is given by

$$x((n-1))_4 = \{0, 1, 2, 2\}$$

where the 0 has been wrapped around to the start and the other values are shifted one unit to the right.

Alternatively, we can view this as a straight translation of the periodic extension outside the range [0, 3] of the given sequence. The periodic extension $x_p(n)$ is shown below:

$$x_p(n) = \{ \dots 2, 2, 0, 1, 2, 2, 0, 1, 2, 2, 0, 1, 2, 2, 0, 1, \dots \}$$

$$n=0$$

The periodic extension, when shifted to the right by one unit, appears as below; and the circularly shifted version $x((n-1))_4$ is the shaded part defined over $0 \le n \le 3$ only:



Example 2.5.2 Given $x(n) = \{1, 2, 2, 0\}$, sketch $x((n-k))_4$ where k is the independent variable.



DFT of circularly-shifted sequence Given that DFT $\{x(n)\} = X(k)$, then

DFT { $x((n-m))_N$ } = $W_N^{km}X(k)$

Conversely, if the X(k) is circularly shifted, the resulting inverse transform will be the multiplication of the inverse of X(k) by a complex exponential: that is, if DFT $\{x(n)\} = X(k)$, then

DFT {
$$W_N^{-ln} x(n)$$
} = $X((k-l))_N$

Note from the above property that a shift in the frequency domain values X(k) generally results in a *complex-valued* inverse sequence x(n) even though the original sequence in the time domain could have been *real-valued*.

Circular convolution The *N*-point circular convolution of two sequences $x_1(n)$ and $x_2(n)$ denoted by $x_1(n) \otimes_N x_2(n)$ is defined as follows:

$$x_{1}(n) \otimes_{N} x_{2}(n) \equiv \sum_{k=0}^{N-1} x_{1}(k) x_{2}((n-k))_{N} = \sum_{k=0}^{N-1} x_{1}((n-k))_{N} x_{2}(k)$$

where $x_I((n-k))_N$ is the reflected and circularly translated version of $x_I(n)$. Note that k is the independent variable, so that $x_I(-k)$ is the reflected version and $x_I(-(k-n))_N$ is simply the reflected version shifted right by n units; the "mod N" makes it a circular shift instead of a linear shift.

If $X_1(k)$ and $X_2(k)$ represent the N-point DFTs of $x_1(n)$ and $x_2(n)$ respectively, i.e.,

$$X_{1}(k) = DFT_{N}\{x_{1}(n)\}$$
 and $X_{2}(k) = DFT_{N}\{x_{2}(n)\}$

Then

$$IDFT_N{X_1(k) X_2(k)} = x_1(n) \otimes_N x_2(n),$$
 or $DFT_N{x_1(n) \otimes_N x_2(n)} = X_1(k) X_2(k)$

This property is used to perform circular convolution of two sequences by first obtaining their DFTs, multiplying the two DFTs, then taking the inverse DFT of the product.

Example 2.5.3 [Circular convolution] For the two sequences $x_1(n) = \{1, 2, 2, 0\}$ and $x_2(n) = \{0, 1, 2, 3\}$ find $y(n) = x_1(n) \bigcirc_4 x_2(n)$.

Solution We use the form $y(n) = x_1(n) \otimes_N x_2(n) = \sum_{k=0}^3 x_1((n-k))_4 x_2(k)$ which uses the circularly

shifted version $x_1((n-k))_4$. The values of the sequence $x_1(k) = \{1, 2, 2, 0\}$ are arranged on a circle in counterclockwise direction starting at point *A*. The sequence $x_1((-k))_4$ is then read off in the clockwise direction starting at *A*. Thus $x_1((-k))_4 = \{1, 0, 2, 2\}$. See figure below.



As an alternative we may also obtain the sequence $x_1((-k))_4$ by *periodically* extending $x_1(k)_4$, reflecting it about k = 0, and truncating it outside the range $0 \le k \le 3$.

The value

$$y(0) = \sum_{k=0}^{3} x_1((0-k))_4 x_2(k)$$

is obtained by lining up $x_1((-k))_4$ below $x_2(k)$, multiplying and adding:

$x_2(k)$	0	1	2	3
$x_1((-k))_4$	1	0	2	2

Thus y(0) = (0) (1) + (1) (0) + (2) (2) + (3) (2) = 10.

For n = 1 the value

$$y(1) = \sum_{k=0}^{\infty} x_1((1-k))_4 x_2(k)$$

is obtained as follows: the sequence $x_1((1-k))_4$ is obtained from $x_1((-k))_4$ by shifting the latter to the right by 1 with wrap around; we then line up $x_1((1-k))_4$ below $x_2(k)$, multiply and add to get y(1):

$x_2(k)$	0	1	2	3
$x_1((1-k))_4$	2	1	0	2

The result is y(1) = (0)(2) + (1)(1) + (2)(0) + (3)(2) = 7.

The procedure is continued for successive values of *n*, at each step using the circularlyshifted-by-1 version of the previous $x_l((n-k))_4$.

Circular convolution – **Matrix method** The circular convolution of the two sequences $x_I(n)$ and $x_2(n)$ is given by:

$$x_{I}(n) \otimes_{N} x_{2}(n) \equiv \sum_{k=0}^{N-1} x_{1}((n-k))_{N} x_{2}(k)$$

- Step 1. By *zero-padding* make sure the two sequences are of the same length, say, *N*.
- Step 2. Arrange the various circularly shifted versions of $x_1(.)$ as a matrix and $x_2(.)$ as a vector; then multiply to get the vector $x_3(.)$ which is the desired result.

The matrix formed by the shifted versions of $x_1(.)$ is shown below. It displays somewhat more terms than is possible to show in the complete multiplication equation shown farther down below.



The complete multiplication step is shown below:



As an example the element $x_3(0)$ is given by

$$x_3(0) = x_1(0) x_2(0) + x_1(N-1) x_2(1) + \dots + x_1(2) x_2(N-2) + x_1(1) x_2(N-1)$$

Carrying out circular convolution to obtain linear convolution If both signals $x_1(n)$ and $x_2(n)$ are of finite lengths N_1 and N_2 respectively, and defined on $[0, N_1-1]$ and $[0, N_2-1]$, respectively, as shown below, the value of N needed so that circular and linear convolution are the same on [0, N-1] can be shown to be $N \ge N_1 + N_2 - 1$.



Example 2.5.4 [Circular and linear convolution] (a) Determine the 4-point circular convolution of the sequences

$$x_1(n) = [1, 2, 3, 1]$$
 and $x_2(n) = [4, 3, 2, 1]$

(b) Evaluate the linear convolution of the above sequences.

(c) Evaluate the linear convolution of the above sequences using circular convolution. **Solution**

(a) The 4-point circular convolution is given by

$$y(n) = x_1(n) \otimes_4 x_2(n) = \{15, 16, 21, 18\}$$

(b) The linear convolution was done in Unit I:

$$y(n) = x_1(n) * x_2(n) = \{4, 11, 20, 18, 11, 5, 1\}$$

(c) Length of sequence $x_1(n) = N_1 = 4$; length of sequence $x_2(n) = N_2 = 4$. Let $N = N_1 + N_2 - 1 = 4 + 4 - 1 = 7$ be the length of each of the zero-padded sequences $x'_1(.)$ and $x'_2(.)$.

$$x_{I}(n) \otimes_{N} x_{2}(n) \equiv \sum_{k=0}^{N-1} x_{1}((n-k))_{N} x_{2}(k)$$

Example 2.5.5 [2007] Compute the circular convolution of the sequences $x_l(n) = \{1, 2, 0, 1\}$ and $x_2(n) = \{2, 2, 1, 1\}$ using the DFT approach.

Solution The sequences are of the same length, so no zero padding is needed. The length of $\{x_1(n) \otimes_N x_2(n)\}$ is 4 (= N). Use the property that if $x_1(n) \leftrightarrow X_1(k)$ and $x_2(n) \leftrightarrow X_2(k)$ and $x_3(n) = x_1(n) \otimes_N x_2(n)$, then $x_3(n) \leftrightarrow X_3(k) = X_1(k) X_2(k)$:

$$x_3(n) = x_1(n) \otimes_N x_2(n) = IDFT\{X_3(k)\} = IDFT_N\{X_1(k) \mid X_2(k)\}$$
 with $N = 4$

where $X_1(k)$ and $X_2(k)$ are the *N*-point DFTs of $x_1(n)$ and $x_2(n)$, respectively. The following steps are involved in computing $x_1(n) \otimes_N x_2(n)$:

- 1. Find $X_1(k) = DFT_4\{x_1(n)\}$ and $X_2(k) = DFT_4\{x_2(n)\}$
- 2. Compute the product $X_1(k) X_2(k)$

3. Compute $x_I(n) \otimes_N x_2(n) = IDFT\{X_I(k)X_2(k)\}$

Example 2.5.6 Compute the linear convolution of the sequences $x(n) = \{1, 2, 0, 1\}$ and $y(n) = \{2, 2, 1, 1\}$ using the DFT approach.

Solution The length of x(n)*y(n) is 7 (= 4+4–1). We zero-pad the sequences to a length of 7 each and perform circular convolution of the 7-point sequences; the result will be the same as the linear convolution of the original 4-point sequences. The following steps are involved in computing x(n)*y(n):

- 1. Augment the sequences x(.) and y(.) by zero-padding: $x_a(n) = \{1, 2, 0, 1, 0, 0, 0\}$ and $y_a(n) = \{2, 2, 1, 1, 0, 0, 0\}$
- 2. Find $X_a(k) = DFT_7\{x_a(n)\}$ and $Y_a(k) = DFT_7\{y_a(n)\}$.
- 3. Compute the product $X_a(k) Y_a(k)$
- 4. Compute $x(n)*y(n) = x(n) \otimes_7 y(n) = IDFT\{X(k) Y(k)\}$

Example 2.5.7 Compute the linear convolution of the sequences $x(n) = \{1, 2\}$ and $y(n) = \{2, 2, 1\}$ using the DFT approach.

Solution The length of x(n)*y(n) is 4 (= 2+3-1). We zero-pad the sequences to a length of 4 each and perform circular convolution of the 4-point sequences; the result will be the same as the linear convolution of the original 2- and 3-point sequences. The following steps are involved in computing x(n)*y(n):

- 1. Augment the sequences x(.) and y(.) by zero-padding: $x_a(n) = \{1, 2, 0, 0\}$ and $y_a(n) = \{2, 2, 1, 0\}$
- 2. Find $X_a(k) = DFT_4\{x_a(n)\}$ and $Y_a(k) = DFT_4\{y_a(n)\}$.
- 3. Compute the product $X_a(k) Y_a(k)$
- 4. Compute $x(n)*y(n) = x_a(n) \otimes_4 y_a(n) = \text{IDFT}\{X_a(k) Y_a(k)\}$

Convolution – **Overlap-and-add** The response, y(n), of a LTI system, h(n), can be obtained by linear convolution

$$y(n) = h(n) * x(n)$$

Let the impulse response $\{h(n), n = 0 \text{ to } M-1\}$ be of finite length *M*. The input sequence $\{x(n), n = 0 \text{ to } S-1\}$ is long but of finite length *S*. Recall that

Length
$$\{y(n)\}$$
 = Length $\{h(n)\}$ + Length $\{x(n)\} - 1 = M + S - 1$,

Further, let h(n) be defined to be zero everywhere except over the interval $[N_1, N_2]$. Similarly, let x(n) be defined to be non-zero over $[N_3, N_4]$. Then y(n) is non-zero over $[(N_1 + N_3), (N_2 + N_4)]$.



One way to perform the convolution in pseudo real time (i.e., real time with a finite delay) is by sectionalizing the input. We divide x(n) into K sections of length M each, where $K = \lceil S/M \rceil$:

$$x_{I}(n) = x(n), \qquad 0 \le n \le M - 1$$

$$0, \qquad \text{elsewhere}$$

$$\dots$$

$$x_{i}(n) = x(n), \qquad (i-1)M \le n \le iM - 1$$

$$0, \qquad \text{elsewhere}$$

$$\dots$$

$$x_{K}(n) = x(n), \qquad (K-1)M \le n \le KM - 1$$

$$0, \qquad \text{elsewhere}$$

In the K^{th} section (the last section) zeros may have to be appended. If, for instance, $x(n) = \{3, -1, 0, 1, 3, 2, 0, 1, 2, 1\}$ with S = 10 and $h(n) = \{1, 1, 1\}$ with M = 3, we have $K = \lfloor 10/3 \rfloor = 4$, with the 4th section containing two appended zeros, and the sections are

 $x_{1}(n) = \{3, -1, 0\}$ $x_{2}(n) = \{1, 3, 2\}$ $x_{3}(n) = \{0, 1, 2\}$ $x_{4}(n) = \{1, 0, 0\}$

In general, then, x(n) can be written as the sum of all the sections

(v

٦

$$x(n) = \sum_{i=1}^{n} x_i(n)$$

and the output y(n) becomes

$$y(n) = h(n) * x(n) = h(n) * \left\{ \sum_{i=1}^{n} x_i(n) \right\}$$

Using the linearity property this becomes

$$y(n) = \sum_{i=1} \{x_i(n) * h(n)\} = \sum_{i=1} y_i(n)$$

where $y_i(n) = x_i(n) * h(n)$ are the output sections. Let us examine $y_i(n)$ and $y_2(n)$. For i = 1 we have

$$y_1(n) = x_1(n) * h(n)$$

Since $x_1(n)$ and h(n) are of lengths M each and they are defined to be non-zero over [0, M-1] and [0, M-1] respectively, the result $y_1(n)$ will be non-zero over [0, 2M-2] and of length (2M-1). Similarly, for i = 2, $x_2(n)$ is defined over [M, 2M-1] while h(n) remains unchanged. The resulting $y_2(n)$ then is non-zero from (0+M) to (M-1+2M-1), i.e., over [M, 3M-2] with a length of (2M-1). Comparing $y_1(n)$ and $y_2(n)$ it is seen that they overlap in the interval $M \le n \le (2M-2)$, over a range of (2M-2) - M + 1 = M-1 points. Consequently the two must be added in this range (see figure). This amounts to adding (M-1) pairs of data.



In a similar fashion $x_3(n)$ is defined over $2M \le n \le (3M-1)$, so that $y_3(n)$ is non-zero over [2*M*, (4*M*-2)]. Comparing $y_2(n)$ and $y_3(n)$ it is seen that they overlap in the interval $2M \le n \le (3M-2)$, consequently the two must be added in this range. This amounts to adding ((3M-2) - 2M + 1) or (M-1) pairs of data.

I	I	Overlap of y_1 and y_2 (<i>M</i> -1) points	I	Overlay y_2 and $(M-1)$ p	p of y_3 oints		Over y ₃ a (<i>M</i> -1)	rlap of and y_4) points	
0	 <i>M</i> –1	<i>M</i> 2 <i>M</i> -2	2	2 <i>M</i>	3 <i>M</i> -2	2	3 <i>M</i>		n

The overlap interval of $y_1(n)$ and $y_2(n)$ is disjoint from that of $y_2(n)$ and $y_3(n)$. In general, the overlap intervals of successive pairs of $y_i(n)$ are mutually exclusive. Thus we calculate successive $y_i(n)$ for i = 0 to K and add each successive $y_i(n)$ to the previous $y_i(n)$ in

the overlap region. Hence the procedure is called the overlap-and-add method. Each convolution could be obtained by using the DFT of size (2M-1) or greater so that the resulting circular convolution would be a linear convolution. In principle rather than using the DFT we could zero-pad h(n) and each of the $x_i(n)$'s to a length of (2M-1) and perform circular convolution to generate the $y_i(n)$'s which are then overlapped and added.

Input divided into sections of length L In the above development we divided x(n) into K sections of length M each, where $K = \lceil S/M \rceil$. This need not be the case. We could divide x(n) into (some number of) sections of length L each. The situation now looks as below and the overlap occurs over a range of (M-1) points – the same as before.



Once again each new section $x_i(n)$ and h(n) are zero-padded to a length of (L+M-1) and circular convolution performed to generate the new $y_i(n)$'s which are overlapped and added to generate y(n).

Note A little reflection shows that the input sequence x(n) need not be of finite length. As the stream of input samples arrives we could sectionalize it into blocks of size *L* and proceed as discussed above to generate the stream of blocks of y(n) as a continuous process.

Example 2.5.7 [Ramesh Babu's Example 3.14] Find the output y(n) of a filter with impulse response $h(n) = \{1, 1, 1\}$ and input $x(n) = \{3, -1, 0, 1, 3, 2, 0, 1, 2, 1\}$.

Symmetry properties of the DFT Notation: $R_N(n) = 1$ in [0, N-1], and 0 elsewhere. Thus $x((n+m))_N R_N(n)$ means the circularly shifted version of the finite length sequence x(n) defined over [0, N-1]. Sometimes the $R_N(n)$ is omitted.

The following properties should be noted.

	Sequence	DFT		Sequence	DFT
1	$x((n+m))_N R_N(n)$	$W_N^{-\kappa m} X(k)$	4	Re $[x(n)]$	$X_{ep}(k)$
2	$x^{*}(n)$	$X^*((-k))_N R_N(k)$	5	$j \operatorname{Im} [x(n)]$	$X_{op}(k)$
3	$x^*((-n))_N R_N(n)$	$X^{*}(k)$			

Example 2.5.8 Show that DFT $\{x((n+m))_N\} = W_N^{-km}X(k)$. Solution By definition DFT $\{x((n+m))_N\} = \sum_{k=1}^{n-1} x((n+m)) W$

FT {
$$x((n+m))_N$$
} = $\sum_{n=0}^{-1} x((n+m)) W_{N} W_{N}$

Set $n+m = \lambda$ so that $n = \lambda - m$ and the limits n = 0 to N-1 become $\lambda = m$ to N-1+m. Then the RHS becomes

$$= \sum_{\lambda=m}^{N-1+m} x((\lambda))_N W_N^{k\lambda} W_N^{-km}$$

The limits on λ will be changed to 0 to N–1 resulting in

$$=\sum_{\lambda=0}^{k-1} x((\lambda))_{N} W_{N}^{k\lambda} W_{N}^{-km} = W_{N}^{-km} X(k)$$

Based on the properties above, we can show that for a real sequence the following symmetry properties of the DFT hold:

1. Re
$$[X(k)] = \text{Re} [X((-k))_N] R_N(k)$$

2. Im $[X(k)] = -\text{Im} [X((-k))_N] R_N(k)$
4. arg $[X(k)] = -\text{arg} [X((-k))_N] R_N(k)$

Example 2.5.9 [2009] Given that the real-valued sequence x(n) defined over $0 \le n \le N-1$ has the DFT $X(k) = X_R(k) + jX_I(k)$, $0 \le k \le N-1$ show that $X_R(k)$ is an even function and $X_I(k)$ is an odd function of k.

Solution By definition we have N^{-1}

$$X(k) = \sum_{n=0}^{N-1} x(n) e^{-jk 2\pi n/N}, \quad 0 \le k \le N-1$$

= $\sum_{n=0}^{N-1} x(n) \{ \cos(2\pi kn/N) - j \sin(2\pi kn/N) \}, \quad 0 \le k \le N-1$
= $\sum_{n=0}^{N-1} x(n) \cos(2\pi kn/N) - j \sum_{n=0}^{N-1} x(n) \sin(2\pi kn/N), \quad 0 \le k \le N-1$
= $\sum_{n=0}^{N-1} x(n) \cos(2\pi kn/N)$ and $X_{I}(k) = -\sum_{n=0}^{N-1} x(n) \sin(2\pi kn/N)$ the above DFT may

be

written

With $X_R(k)$

$$X(k) = X_R(k) + jX_I(k), \qquad 0 \le k \le N-1$$

Since $\cos(2\pi kn/N)$ is an even function of k, that is, $\cos(2\pi (-k)n/N) = \cos(2\pi kn/N)$ for all k, it follows that $X_R(k)$ is an even function of k, that is, $X_R(-k) = X_R(k)$ for all k.

Similarly, since $\sin(2\pi kn/N)$ is an odd function of k, it follows that $X_I(k)$ is an odd function of k.

Do the above results depend on whether x(n) is real-valued or not?

Filtering through DFT/FFT

Filtering of a sequence x(n) may be done in the discrete time domain using the difference equation. Alternatively, we may work in the frequency domain: Given x(n) we first find its DFT, X(k) and then set selected components of X(k) = 0 (this is done so as to preserve the symmetry

properties of X(k) about the mid point of the sequence k = N/2). The resulting DFT is denoted $X_{f}(k)$. We then find the IDFT of $X_{f}(k)$ which we shall denote as $x_{f}(n)$, which should be a filtered version of x(n). Thus x(n) has been filtered entirely in the discrete frequency domain.

We illustrate with a signal consisting of two frequency components, a 2 Hz and a 4 Hz. Given the signals $x_l(t) = \cos 2\pi 2t$ and $x_2(t) = \cos 2\pi 4t$, the signal $x(t) = x_l(t) + x_2(t)$ is sampled at 16 Hz. We construct below the discrete-time sequence x(n) and find its 8-point DFT, i.e., the X(k) values. We then filter the signal in the frequency domain, i.e., work on the DFT values instead of on the x(n) values and denote the "filtered" DFT values by $X_f(k)$. We then take the IDFT of $X_f(k)$ resulting in the sequence $x_f(n)$.

$$\begin{aligned} x(t) &= \cos 2\pi 2t + \cos 2\pi 4t \\ x(nT) &= x(n) = \cos 2\pi 2nT + \cos 2\pi 4nT = \cos 2\pi 2n(1/16) + \cos 2\pi 4n(1/16) \\ &= \cos (\pi n/4) + \cos (\pi n/2) \end{aligned}$$

The cos $2\pi 2t$ component has an analog frequency of 2 Hz. When sampled at 16 Hz its digital frequency is 1/8 cycles/sample. Similarly, the 4 Hz component, cos $2\pi 4t$, sampled at 16 sps, has a digital frequency of 1/4 cycles/sample.

Example 2.6.1 (a) Find the frequency and period of (i) $x_1(n) = \cos (\pi n/4)$ and (ii) $x_2(n) = \cos (\pi n/2)$. Sketch the sequences $x_1(n)$, $x_2(n)$, and $x(n) = x_1(n) + x_2(n)$ for $0 \le n \le 7$. **Solution** Arrange $\cos (\pi n/4)$ in the format $\cos (2\pi f n)$. Thus, $\cos (\pi n/4) = \cos (2\pi (1/8)n)$, from which the digital frequency is identified as f = 1/8 cycle/sample or $\omega = \pi/4$ rad/sample. The sequence values are:

$$x_{I}(n) = \begin{cases} 1 & 1 & 1 & 1 & 1 \\ 1, \frac{1}{\sqrt{2}}, - & \frac{1}{\sqrt{2}}, - & \frac{1}{\sqrt{2}}, - & \frac{1}{\sqrt{2}} \end{cases}$$
$$x_{2}(n) = \{1, 0, -1, 0, 1, 0, -1, 0\}$$
$$x(n) = x_{I}(n) + x_{2}(n) = \begin{cases} 2, \frac{1}{\sqrt{2}}, -1, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \end{cases}$$

The sequences $x_1(n)$, $x_2(n)$ and x(n) are sketched below.



In MATLAB the following segment plots the three functions $x_1(t)$, $x_2(t)$ and x(t).

t = 0:1/160:0.5; x1t = cos(2*pi*2*t); x2t = cos(2*pi*4*t); xt = x1t + x2t; % subplot(3,1,1), plot(t,x1t); xlabel('t'), ylabel('x1(t)'); title('x1(t) = cos 2 pi2t'; % subplot(3,1,2), plot(t,x2t); xlabel('t'), ylabel('x2(t)'); title('x2(t) = cos 2 pi4t'); % subplot(3,1,3), plot(t,xt); xlabel('t'), ylabel('x(t)'); title('x(t) = cos 2 pi2t + cos 2 pi4t');



In MATLAB the following segment plots the two functions $x_1(n)$, $x_2(n)$ and x(n).

$$\label{eq:states} \begin{split} \%t &= 0:1/160:0.5; \ x1t = \cos(2^*pi^*2^*t); \ x2t = \cos(2^*pi^*4^*t); \\ \%xt &= x1t + x2t; \\ \% \\ n &= 0:8; \ x1n = \cos(pi^*n/4); \ x2n = \cos(pi^*n/2); \\ xn &= x1n + x2n, \\ \% \\ subplot(3,1,1), \ stem(n,x1n); \ xlabel('n'), \ ylabel('x1(n)'); \\ title('x1(n) &= \cos(\langle pi \ n/4 \rangle'); \\ \% \end{split}$$

subplot(3,1,2), stem(n,x2n); xlabel('n'), ylabel('x2(n)'); title('x2(n) = cos(\pi n/2)'); % subplot(3,1,3), stem(n,xn); xlabel('n'), ylabel('x(n)'); title('x(n) = cos(\pi n/4)+ cos(\pi n/2)');

The sequence is $x(n) = \{2, 0.707, -1, -0.707, 0, -0.707, -1, 0.707\}$



Example 2.6.1 (b) Find the 8-point DFT (8-point FFT using either DIT or DIF, or use the direct calculation) of $x_l(n) = \cos(\pi n/4)$, for $0 \le n \le 7$. Sketch the sequence $X_l(k)$.

In MATLAB the following segment plots the magnitude and phase angle of $X_l(k)$:

n = 0:7; x1n = cos(pi*n/4), X1k = fft(x1n), %
% MX1k = magnitude of X1k, AX1k = phase angle of X1k MX1k = abs(X1k); AX1k = angle(X1k); %
subplot(2,1,1), stem(n, abs(X1k)); xlabel('k'), ylabel('|X1(k)|'); title ('Magnitude of X1(k)'); %
subplot(2,1,2), stem(n, angle(X1k)); xlabel('k'), ylabel('<X1(k)'); title ('Phase of X1(k)');

The sequence and DFT values are:

 $x_{l}(n) = \{1 \quad 0.707 \quad 0 - 0.707 - 1 - 0.707 - 0 \quad 0.707 \}$ $X_{l}(k) = \{-0 \ (4 - 0i) \ 0 \ (-0 - 0i) \ 0 \ (-0 + 0i) \ 0 \ (4 + 0i) \}$ $|X_1(k)| = \{0 \ 4 \ 0 \ 0 \ 0 \ 0 \ 0 \ 4\}$ $\angle X_1(k) = \{3.1416 \ -0.0000 \ 0 \ -2.9873 \ 0 \ 2.9873 \ 0 \ 0.0000\}$

!!! Note that all phase angle values should be zeros, that is, $\angle X_1(k) = 0$ for all k.

The MATLAB plots for magnitude and phase angle are shown below. The 2 Hz component is indicated by $X_l(1) = 4$ (and, from symmetry considerations, $X_l(7) = 4$).

With regard to the phase angle, strictly speaking, the phase of $X_l(k) = 0$ for all k since $X_l(k)$ is real valued for all k. *Check out why!!!*: When the input is listed explicitly as

x1n = [1, 1/sqrt(2), 0, -1/sqrt(2), -1, -1/sqrt(2), 0, 1/sqrt(2)]

the program gives the correct phase angle calculations, but not when it is specified implicitly as



$$n = 0:7; x1n = cos(pi*n/4)$$

Example 2.6.1 (c) Find the 8-point DFT of $x(n) = \cos(\pi n/4) + \cos(\pi n/2)$, for $0 \le n \le 7$. Sketch the sequence X(k).

Solution Consider the 8-point sequence x(n) for n = 0 to 7 obtained by sampling x(t) at 16 Hz. The length of the sequence is N = 8. The sequence values are

$$x(n) = \begin{cases} 2, \frac{1}{\sqrt{2}}, -1, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \end{cases}$$

The corresponding DFT is given by $\sum_{N=1}^{N-1}$

$$X(k) = \sum_{n=0}^{\infty} x(n) \ e^{-j 2\pi k \ n \ / N} , \qquad k = 0 \ \text{to} \ (N-1)$$

and may be obtained using either the DIT or the DIF form of FFT:

The DFT is sketched below: there is a component at 2 Hz. (k = 1) and another at 4 Hz. (k = 2) as would be expected.



Note that the DFT sequence shows a component at 2 Hz and another at 4 Hz, corresponding to k = 1 and 2 respectively.

In MATLAB the following segment plots the magnitude and phase angle of X(k). Note that x(n) is specified by an explicit list.

%n = 0:7; xn = cos(pi*n/4) + cos(pi*n/2), xn = [2, 1/sqrt(2), -1, -1/sqrt(2), 0, -1/sqrt(2), -1, 1/sqrt(2)], Xk = fft(xn), % % MXk = magnitude of Xk, AXk = phase angle of Xk MXk = abs(Xk), AXk = angle(Xk), % k = 0:7; subplot(2,1,1), stem(k, abs(Xk)); xlabel('k'), ylabel('|X(k)|'); title ('Magnitude of X(k)'); % subplot(2,1,2), stem(k, angle(Xk)); xlabel('k'), ylabel('<X(k)'); title ('Phase of X(k)');

The sequence and DFT values are:

 $x(n) = \{2, 0.707, -1, -0.707, 0, -0.707, -1, 0.707\}$ $X(k) = \{0 \quad 4 \quad 4 \quad 0 \quad 0 \quad 0 \quad 4 \quad 4\}$ $| \quad X \ \langle k \rangle = \{0 \quad 4 \quad 4 \quad 0 \quad 0 \quad 0 \quad 4 \quad 4\}$ $\angle X \ \langle k \rangle = \{0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0\}$ The MATLAB plots for magnitude and phase angle are shown below. The 2 Hz component is indicated by X(1) = 4 (and, from symmetry considerations, X(7) = 4). Similarly, the 4 Hz component is indicated by X(2) = X(6) = 4.



Example 2.6.1 (d) [Filtering] Next we would like to filter the sequence $x(n) = \cos(\pi n/4) + \cos(\pi n/2)$ so that the 4 Hz component is removed.

Solution In terms of the DFT sequence X(k) this means setting X(2) and X(6) to zero. In order to preserve the symmetry properties of the DFT we should set *both* X(2) and X(6) = 0, not just X(2). The resulting DFT sequence is denoted $X_f(k)$ and is given by

$$X_{f}(k) = \{0, 4, 0, 0, 0, 0, 0, 4\}$$

$$\uparrow \qquad k=0$$

We next find the inverse DFT of the above $X_f(k)$ by using either the DIT or the DIF form of the FFT. The result is denoted by $x_f(n)$. Using the IDFT formula we have

$$X_{f}(n) = \frac{1}{N} \sum_{\substack{k=0\\k=0}}^{N-1} X_{f}(k) e^{j2\pi k n/N}, \qquad n = 0, 1, \dots, N-1$$
$$= \frac{1}{8} \sum_{\substack{k=0\\k=0}}^{N-1} X_{f}(k) e^{j2\pi k n/8}, \qquad n = 0 \text{ to } 7$$

It will be seen that $x_f(n)$ is equal to the original $x_I(n)$ component (= cos $\pi n/4$, from the 2 Hz. component; see plot of $x_I(n)$ earlier); that is

$$x_{f}(n) = x_{I}(n) = \begin{cases} 1 & 1 & 1 & 1 \\ 1, \frac{1}{\sqrt{2}}, - & \frac{-1}{\sqrt{2}}, - & \frac{1}{\sqrt{2}}, & \frac{1}{\sqrt{2}} \end{cases}$$

This is low pass filtering where we have selectively removed the 4 Hz component. In MATLAB the following segment finds the inverse DFT, $x_f(n)$, from the given $X_f(k)$.

> Xfk = [0, 4, 0, 0, 0, 0, 0, 4], n = 0:7; xfn = ifft(Xfk), stem(n, xfn); xlabel('n'), ylabel('xf(n)'); title ('Filtered signal xf(n)');

The filtered version is $x_f(n) = \{1 \ 0.707 \ 0 \ -0.707 \ -1 \ -0.707 \ 0 \ 0.707\}$



To remove all frequency components above 2 Hz (in this example the 4, 6 and 8 Hz components), we set X(2) = X(6) = 0 for the 4 Hz, X(3) = X(5) = 0 for the 6 Hz, and X(4) = 0 for the 8 Hz, once again preserving symmetry. In this example, of course, there are no 6 or 8 Hz components.

Similarly high pass filtering is done by deleting X(1) and X(7) – set them to zero – preserving symmetry once again. In this case $X_{fHP}(k) = \{0, 0, 4, 0, 0, 0, 4, 0\}$.

Picket-fence effect

Example 2.7.1 The signal $x(t) = \cos 2\pi 2t$ is sampled at 16 Hz. (a) What frequency components do you expect to see in its DFT? (b) Take 8 samples and calculate the 8-point DFT. **Solution** (b) Using $x(n) = \cos 2\pi 2n(1/16) = \cos (\pi n/4)$, the sequence values are $x(n) = \{1, 1/\text{sqrt}(2), 0, -1/\text{sqrt}(2), -1, -1/\text{sqrt}(2), 0, 1/\text{sqrt}(2)\}$

Note that the average value of the sequence (the dc component) is zero. The MATLAB program follows:

xn = [1, 1/sqrt(2), 0, -1/sqrt(2), -1, -1/sqrt(2), 0, 1/sqrt(2)]; Xk = fft(xn), %
% MXk = magnitude of Xk, AXk = phase angle of Xk MXk = abs(Xk); AXk = angle(Xk); %
k = 0:7; subplot(2,1,1), stem(k, abs(Xk)); xlabel('k'), ylabel('|X(k)|'); title ('Magnitude of X(k)'); %
subplot(2,1,2), stem(k, angle(Xk)); xlabel('k'), ylabel('<X(k)'); title ('Phase of X(k)');

The DFT values are:

The frequency resolution is $F_s/N = 16/8 = 2$ Hz. The table below shows that the 8-point DFT contains a component at 2 Hz, corresponding to k = 1. The DFT values are all real numbers symmetrically disposed about k = 4, the center of symmetry.

$$k=0$$

$$\$$$

$$X(k) = \{0, 4, 0, 0, 0, 0, 0, 4\}$$

$$Hz < 0 2 4 6 8 10 12 14$$

$$(F_s/2)$$



Example 2.7.2 [Zero-padding] The first 8 sample values of a 2-Hz cosine, $x(t) = \cos 2\pi 2t$, obtained at a sampling rate of 16 samples/second are given below:

$$x(n) = \{1, 1/\text{sqrt}(2), 0, -1/\text{sqrt}(2), -1, -1/\text{sqrt}(2), 0, 1/\text{sqrt}(2)\}$$

Find the 16-point DFT using zero padding. **Solution** The zero-padded sequence is given by

Note that when zero padded the average value of the sequence is no longer zero. The frequency resolution of the DFT is

Frequency resolution =
$$\frac{Sampling \ Frequency}{Number \ of \ point \ s} = \frac{F_s}{N} = \frac{16 \ Hz}{16} = 1 \ Hz$$

The 2 Hz component corresponds to X(k) with k = 2.

In MATLAB:

xn = [1, 1/sqrt(2), 0, -1/sqrt(2), -1, -1/sqrt(2), 0, 1/sqrt(2)], Xk = fft(xn, 16), %x(n) is zero-padded to a length of 16 % % MXk = magnitude of Xk, AXk = phase angle of Xk MXk = abs(Xk), AXk = angle(Xk), % k = 0:15; subplot(3,1,1), stem(k, abs(Xk)); xlabel('k'), ylabel('|X(k)|'); title ('Magnitude of X(k)'); % subplot(3,1,2), stem(k, angle(Xk)); xlabel('k'), ylabel('<X(k)'); title ('Phase of X(k)'); % xn16 = ifft(Xk), n = 0:15; subplot(3,1,3), stem(n, xn16); xlabel('n'), ylabel('x(n)'); title ('Zero-padded sequence xn16');

The DFT (reproduced from the MATLAB output) is

$$X(k) = \{0, (1+1.7654i), 4, (1-2.8478i), 0, (1-0.8478i), 0, (1-0.2346i), 0, \{1+0.2346i), 0, (1+0.8478i), 0, (1+2.8478i), 4, (1-1.7654i)\}$$





Frequency resolution =
$$\frac{Sampling \ Frequency}{Number of \ point \ s} = \frac{F_s}{N} = \frac{16 \ Hz}{16} = 1 \ Hz.$$

The 2 Hz component corresponds to X(k) with k = 2.

t = 0: 1/16: 15/16; xn = cos (2*pi*2*t), Xk = fft(xn), %
MXk = magnitude of Xk, AXk = phase angle of Xk MXk = abs(Xk), AXk = angle(Xk), %
k = 0:15; subplot(3,1,1), stem(k, abs(Xk)); xlabel('k'), ylabel('|X(k)|'); title ('Magnitude of X(k)');
%
subplot(3,1,2), stem(k, angle(Xk)); xlabel('k'), ylabel('<X(k)');
title ('Phase of X(k)');
%
n = 0:15;
subplot(3,1,3), stem(n, xn); xlabel('n'), ylabel('x(n)');
title ('Sequence x(n)');</pre>

The DFT (reproduced from the MATLAB output) is

 $X(k) = \{0, (-0+0i), (8-0i), (0-0i), (0-0i), (0-0i), (0-0i), (0+0i), (-0, (0-0i), (0+0i), (0+0i), (0+0i), (0+0i), (0+0i), (0+0i), (0+0i), (-0-0i)\}$



Specifying the sequence values explicitly as below produces correct phase values (rather than implicitly by n = 0: 1 : 15; x = cos (pi*n/4), X = fft(x))

%Listing the sequence values explicitly xn1 = [1, 1/sqrt(2), 0, -1/sqrt(2), -1, -1/sqrt(2), 0, 1/sqrt(2)], xn2 = [1, 1/sqrt(2), 0, -1/sqrt(2), -1, -1/sqrt(2), 0, 1/sqrt(2)], xn = [xn1,xn2], Xk = fft(xn), % % MXk = magnitude of Xk, AXk = phase angle of Xk MXk = abs(Xk), AXk = angle(Xk), % k = 0:15 subplot(3,1,1), stem(k, abs(Xk)); xlabel('k'), ylabel('|X(k)|'); title ('Magnitude of X(k)'); % subplot(3,1,2), stem(k, angle(Xk)); xlabel('k'), ylabel('<X(k)'); title ('Phase of X(k)'); % n = 0:15; subplot(3,1,3), stem(n, xn); xlabel('n'), ylabel('x(n)'); title ('16-point sequence xn');



k<	0		2						8							15
X(k)	{0	0	8	0	0	0	0	0	0	0	0	0	0	0	8	0}
Hz <			2						8							

Example 2.7.4 The signal $x(t) = \cos 2\pi 2t$ is sampled at 12 Hz (still satisfies the sampling theorem). Calculate (1) the 6-point DFT and (2) the 8-point DFT. Compare the results. **Solution** The resulting sequence is $x(n) = \cos 2\pi 2n(1/12) = \cos (\pi n/3)$. (1) 6-point DFT.

$$n = 0: 1:5; x = cos (pi*n/3), X = fft(x)$$

The samples are

$$x(n) = \left\{ \frac{1}{2}, \frac{1}{2}, -\frac{1}{2}, -1, -\frac{1}{2}, \frac{1}{2} \right\}$$

The average value of the sequence (the dc component) is zero. The frequency resolution is $F_s/N = 12/6 = 2$ Hz. The DFT is

$$X(k) = \{0, 3, 0, 0, 0, 3\}$$

which does show a component at 2 Hz.

(2) 8-point DFT.

 $n = 0: 1: 7; x = \cos(pi*n/3), X = fft(x)$

The samples are

re

$$x(n) = \left\{ \begin{array}{c} 1 & 1 \\ 1, -, -\frac{1}{2}, -1, -\frac{1}{2}, \frac{1}{2}, 1, \frac{1}{2} \end{array} \right\}$$

Note that the average value of the sequence (the dc component) is *not* zero. The frequency resolution is $F_s/N = 12/8 = 1.5$ Hz. The DFT is

$$X(k) = \{1.5, (2.5607 + j2.5607), -j1.5, (0.4393 - j0.4393), 0.5, (0.4393 + j0.4393), j1.5, (2.5607 - j2.5607)\}$$

It is not possible to know that there is a component at 2 Hz.

Example 2.7.5 For the signal $x(t) = \cos 2\pi 2t + \cos 2\pi 4t$ choose a sampling frequency of 12Hz (still satisfies sampling theorem). Find the *N*-point DFT for (a) N = 6, (b) N = 8, and (c) N = 12.

(a) For N = 6

 $x(n) = \cos((2\pi 2n/12)) + \cos((2\pi 4n/12)) = \cos((\pi n/3)) + \cos((2\pi n/3))$

 $n = 0: 1: 5; x = \cos(pi*n/3) + \cos(2*pi*n/3), X = fft(x)$

The samples are

$$x(n) = \{2, 0, -1, 0, -1, 0\}$$

The dc component is zero. The frequency resolution is $F_s/N = 12/6 = 2$ Hz. The DFT is

$$X(k) = \{0, 3, 3, 0, 3, 3\}$$

Both the 2 Hz and the 4 Hz components show up.

k٢	0			3		5
X(k)	{0	3	3	0	3	3
Hz <	0	2	4	6	8	10

(b) For N = 8

$$x(n) = \cos((2\pi 2n/12)) + \cos((2\pi 4n/12)) = \cos((\pi n/3)) + \cos((2\pi n/3))$$

$$n = 0: 1: 7; x = \cos(pi*n/3) + \cos(2*pi*n/3), X = fft(x)$$

The samples are

$$x(n) = \{2, 0, -1, 0, -1, 0, 2, 0\}$$

This particular set has a *nonzero* dc component. The frequency resolution is $F_s/N = 12/8 = 1.5$ Hz. The DFT is

$$X(k) = \{2, (3 + j3), 0, (3 - j3), 2, (3 + j3), 0, (3 - j3)\}$$

Neither the 2 Hz nor the 4 Hz component can show up.

k<	0				4			7
X(k)	{2	3+j3	0	3—ј3	2	3+j3	0	3—j3
Hz <	0	1.5	3	4.5	6	7.5	9	10.5



Example 2.7.6 For the signal $x(t) = \cos 2\pi 2t + \cos 2\pi 3t + \cos 2\pi 4t$ choose a sampling frequency of 12Hz (still satisfies sampling theorem). Find the *N*-point DFT for (a) N = 6, (b) N = 8.

(a) For N = 6

 $x(n) = \cos (2\pi 2n/12) + \cos (2\pi 3n/12) + \cos (2\pi 4n/12)$ = cos (\pi n/3) + cos (\pi n/2) + cos (2\pi n/3)

 $n = 0: 1: 5; x = \cos(pi*n/3) + \cos(\pi n/2) + \cos(2*pi*n/3), X = fft(x)$

The samples are

$$x(n) = \{3, 0, -2, 0, 0, 0\}$$

The dc component is not zero. The frequency resolution is $F_s/N = 12/6 = 2$ Hz. The DFT is

$$X(k) = \{1, (4 + j1.732), (4 - j1.732), 1, (4 + j1.732), (4 - j1.732)\}$$

Both the 2 Hz and the 4 Hz components show up, but the 3 Hz component is missing.

k<	0			3		5
X(k)	{1	4 + j1.732	4 - j1.732	1	4 + j1.732	4 - j1.732
Hz <	0	2	4	6	8	10
Digital Signal Processing – 2(B)

II(B). Fast Fourier transform

Fast Fourier Transform (FFT) – Radix-2 decimation in time and decimation in frequency FFT Algorithms, Inverse FFT, and FFT for composite N.

Contents:

Introduction Radix-2 decimation-in-time FFT (Cooley-Tukey) Radix-2 decimation-in-frequency FFT (Sande-Tukey) Inverse DFT using the FFT algorithm *Decimation-in-time algorithm for N = 4 (Cooley-Tukey formulation) *Decimation-in-frequency algorithm for N = 4 (Sande-Tukey formulation) FFT with general radix

Introduction

For a finite-duration sequence x(n) of length N, the DFT sum may be written as

$$X(k) = \sum_{n=0}^{\infty} x(n) W_N^{kn}, \ k = 0, \ 1, \ \dots, \ N-1$$

where $W_N = e^{-j2\pi/N}$. There are a total of *N* values of *X*(.) ranging from *X*(0) to *X*(*N*-1). The calculation of *X*(0) involves no multiplications at all since every product term involves $W_N^0 = e^{-j0} = 1$. Further, the first term in the sum always involves W_N^0 or $e^{-j0} = 1$ and therefore does not require a multiplication. Each *X*(.) calculation other than *X*(0) thus involves (*N*-1) complex multiplications. And each *X*(.) involves (*N*-1) complex additions. Since there are *N* values of *X*(.) the overall DFT requires (*N*-1)² complex multiplications and *N*(*N*-1) complex additions. For large *N* we may round these off to N^2 complex multiplications and the same number of complex additions.

Each complex multiplication is of the form

$$(A + jB) (C + jD) = (AC - BD) + j(BC + AD)$$

and therefore requires four real multiplications and two real additions. Each complex addition is of the form

$$(A + jB) + (C + jD) = (A + C) + j(B + D)$$

and requires two real additions. Thus the computation of all N values of the DFT requires $4N^2$ real multiplications and $4N^2(=2N^2+2N^2)$ real additions.

Efficient algorithms which reduce the number of multiply-and-add operations are known by the name of **fast Fourier transform** (FFT). The Cooley-Tukey and Sande-Tukey FFT algorithms exploit the following properties of the **twiddle factor** (phase factor), $W_N = e^{-j2\pi/N}$ (the factor $e^{j2\pi/N}$ is called the N^{th} principal root of 1):

1. Symmetry property
$$W_N^2 = -W_N$$

1. Symmetry property $W_N^2 = -W_N$ 2. Periodicity property $W_N^{k+N} = W_N^k$

To illustrate, for the case of N = 8, these properties result in the following relations:

$$W_{8}^{0} = -W_{8}^{4} = 1 \qquad W_{8}^{1} = -W_{8}^{5} = \frac{1-j}{\sqrt{2}}$$
$$W_{8}^{2} = -W_{8}^{6} = -j \qquad W_{8}^{3} = -W_{8}^{7} = -\frac{\sqrt{2}}{\sqrt{2}}$$

The use of these properties reduces the number of complex multiplications from N^2 to $\frac{N}{2}\log_2 N$ (actually the number of multiplications is less than this because several of the multiplications by W_N^r are really multiplications by ± 1 or $\pm j$ and don't count); and the number of complex additions are reduced from N^2 to $N \log_2 N$. Thus, with each complex multiplication requiring four real multiplications and two real additions and each complex addition requiring two real additions, the computation of all N values of the DFT requires

Number of real multiplications = $4 \begin{pmatrix} N \\ -2 \log_2 N \end{pmatrix} = 2N \log_2 N$ Number of real additions = $2N \log_2 N + 2 \left| \begin{pmatrix} N \\ -2 \log_2 N \end{pmatrix} \right| = 3N \log_2 N$

We can get a rough comparison of the speed advantage of an FFT over a DFT by computing the number of multiplications for each since these are usually more time consuming than additions. For instance, for N = 8 the DFT, using the above formula, would need $8^2 = 64$ complex multiplications, but the radix-2 FFT requires only 12 (= $\log 8 = 4 \times 3$).

No. of points	No. of complex multiplicationsDFTFFT		No. of real multiplication	
N			DFT	FFT
32	1024	80	4096	320
128	16384	448	65536	1792
1024	1048576	5120	4194304	20480

Number of multiplications: DFT vs. FFT

We consider first the case where the length N of the sequence is an integral power of 2, that is, $N = 2^{v}$ where v is an integer. These are called **radix-2 algorithms** of which the **decimation-in-time (DIT)** version is also known as the **Cooley-Tukey algorithm** and the **decimation-in-frequency (DIF)** version is also known as the **Sande-Tukey algorithm**. We show first how the algorithms work; their derivation is given later.

For a radix of (r = 2), the **elementary computation** (*EC*) known as the **butterfly** consists of a single complex multiplication and two complex additions.

If the number of points, N, can be expressed as $N = r^m$, and if the computation algorithm is carried out by means of a succession of *r*-point transforms, the resultant FFT is called a **radix***r* **algorithm**. In a radix-*r* FFT, an elementary computation consists of an *r*-point DFT followed by the multiplication of the *r* results by the appropriate twiddle factor. The number of *EC*s required is

$$C_r = \frac{N}{\log_r N}$$

which decreases as r^{r} increases. Of course, the complexity of an *EC* increases with increasing *r*. For r = 4, the *EC* requires three complex multiplications and several complex additions.

Suppose that we desire an *N*-point DFT where *N* is a composite number that can be factored into the product of integers

$$N=N_1 N_2 \ldots N_m$$

If, for instance, N = 64 and m = 3, we might factor N into the product $64 = 4 \times 4 \times 4$, and the 64-point transform can be viewed as a three-dimensional $4 \times 4 \times 4$ transform.

If *N* is a prime number so that factorization of *N* is not possible, the original signal can be *zero-padded* and the resulting new composite number of points can be factored.

Radix-2 decimation-in-time FFT (Cooley-Tukey)

Procedure and important points

- 1. The number of input samples is $N = 2^{v}$ where v is an integer.
- 2. The input sequence is shuffled through bit-reversal. The index n of the sequence x(n) is expressed in binary and then reversed.
- 3. The number of stages in the flow graph is given by $v = \log_2 N$.
- 4. Each stage consists of N/2 butterflies.
- 5. Inputs/outputs for each butterfly are separated as follows:

Separation = 2^{m-1} samples where m = stage index, stages being numbered from left to right (that is, m = 1 for stage 1, m = 2 for stage 2 etc.).

This amounts to separation increasing from left to right in the order 1, 2, 4, ..., N/2.



- 6. The number of complex additions = $N \log_2 N$ and the number of complex multiplications N is $\log N$.
- 7. The elementary computation block in the flow graph, called the butterfly, is shown here. This is an **in-place calculation** in that the outputs $(A + BW_N^k)$ and $(A BW_N^k)$ can be

computed and stored in the same locations as A and B.



Example 2.2.1 Radix-2, 8-point, decimation-in-time FFT for the sequence

$$n < 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7$$

$$x(n) = \{1, \ 2 \quad 3 \quad 4 \quad -4 \quad -3 \quad -2 \quad -1\}$$

1

Solution The twiddle factors are

$$W_{8}^{0} = 1 \qquad \qquad W_{8}^{2} = \begin{pmatrix} e^{-j2\pi/8} \\ e \end{pmatrix}^{2} = e^{-j\pi/2} = -j \qquad \qquad W_{8}^{1} = e^{-j2\pi/8} = e^{-j\pi/4} = \frac{1}{\sqrt{2}} j \frac{1}{\sqrt{2}} \frac{1}{\sqrt$$

One of the elementary computations is shown below:



The signal flow graph follows:



8-point FFT using DIT							
	Results of the first stage						
Input	Stage 1	Stage 2	Stage 3 (Output)				
1	$1 + (-4) \cdot 1 = -3$						
- 4	$1 - (-4) \cdot 1 = 5$						
3	$3 + (-2) \cdot 1 = 1$						
- 2	$3 - (-2) \cdot 1 = 5$						
2	$2 + (-3) \cdot 1 = -1$						
- 3	$2 - (-3) \cdot 1 = 5$						
4	$4 + (-1) \cdot 1 = 3$						
- 1	$4 - (-1) \cdot 1 = 5$						

Results of the second stage						
Input	Stage 1	Stage 2	Stage 3 (Output)			
1	-3	$-3 + 1 \cdot 1 = -2$				
-4	5	$5+5.(-j)=5$ $2e^{-j\pi/4}$				
3	1	$-3 - 1 \cdot 1 = -4$				
- 2	5	$5-5.(-j)=5$ $2e^{j\pi/4}$				
2	-1	$-1 + 3 \cdot 1 = 2$				
- 3	5	$5+5.(-j)=5$ $2e^{-j\pi/4}$				
4	3	$-1 - 3 \cdot 1 = -4$				
- 1	5	$5-5.(-j)=5$ $2e^{j\pi/4}$				

	Results of the third stage				
Input	Stage 1	Stage 1 Stage 2 Stage 3 (Output)			
1	-3	-2	-2+2.1=0		
-4	5	$5 \ 2e^{-j\pi/4}$	$52e^{-j\pi/4}$ + $52e^{-j\pi/4}$. $e^{-j\pi/4}$ = $5-j12.07$		
3	1	-4	$-4 + (-4) \cdot (-j) = -4 + j4 = 4 2e^{j3\pi/4}$		
- 2	5	ο jπ/4	$5 2e^{j\pi/4} + 5 2e^{j\pi/4}$. $e^{-j3\pi/4} = 5 - j 2.07$		
2	-1	2	-2-2.1 = -4		
- 3	5	$5 \; 2e^{-j\pi/4}$	$5 2e^{-j\pi/4} - 5 2e^{-j\pi/4}$. $e^{-j\pi/4} = 5 + j2.07$		
4	3	-4	$-4 - (-4) \cdot (-j) = -4 - j4 = 4 2e^{-j3\pi/4}$		
- 1	5	e j \pi / 4	$5 \ 2e^{j\pi/4} - 5 \ 2e^{j\pi/4} \cdot e^{-j3\pi/4} = 5 + j12.07$		

The DFT is $X(k) = \{0, (5 - j12.07), (-4 + j4), (5 - j2.07), -4, (5 + j2.07), (-4 - j4), (5 + j12.07)\}$

The MATLAB program:

x = [1, 2, 3, 4, -4, -3, -2, -1], X = fft(x)

Radix-2 decimation-in-frequency FFT (Sande-Tukey)

Procedure and important points

- 1. The number of input samples is $N = 2^{v}$ where v is an integer.
- 2. The input sequence is in natural order; the output is in bit-reversed order.
- 3. The number of stages in the flow graph is given by $v = \log_2 N$.
- 4. Each stage consists of N/2 butterflies.
- 5. Inputs/outputs for each butterfly are separated in the reverse order from that of the DIT. The separation decreases *from left to right* in the order N/2, ..., 4, 2, 1.
- 6. The number of complex additions = $N \log_2 N$ and the number of complex multiplications is $\sum_{n=1}^{N} \log N$.

2, 2^{2}

7. The basic computation block in the flow graph of the DIF FFT is the butterfly shown here. This is an **in-place calculation** in that the two outputs (A + B) and $(A - B) W_N^k$ can be computed and stored in the same locations as A and B.



Example 2.3.1 Radix-2, 8-point, decimation-in-frequency FFT for the sequence

 $n < 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7$ $x(n) = \{1, 2 \quad 3 \quad 4 \quad -4 \quad -3 \quad -2 \quad -1\}$

Solution The twiddle factors are the same as in the DIT FFT done earlier (both being 8-point DFTs):

$$\begin{array}{c} 0 \\ W_8 = 1 \\ W_8 = \left(e^{-j2\pi/8}\right)^2 = e^{-j\pi/2} = -j \\ W_8 = \left(e^{-j2\pi/8}\right)^2 = e^{-j\pi/8} = -j \\ W_8 = \left(e^{-j2\pi/8}\right)^2 = e^{-j\pi/8} = -j \\ W_8 = \left(e^{-j\pi/8}\right)^2 = -j \\$$

One of the elementary computations is shown below:



The signal flow graph follows:



8-point FFT using DIF							
	Results of the first stage						
Input	Stage 1	Stage 2	Stage 3 (Output)				
1	1 + (-4) = -3						
2	2 + (-3) = -1						
3	3 + (-2) = 1						
4	4 + (-1) = 3						
-4	(1 - (-4)) 1 = 5						
- 3	$(2 - (-3)) e^{-j\pi/4} = 5 e^{-j\pi/4}$						
-2	(3 - (-2))(-j) = -j5						
- 1	$(4 - (-1)) e^{-j3\pi/4} = 5 e^{-j3\pi/4}$						

	Results of the second stage						
Input	Stage 1	Stage 2	Stage 3 (Output)				
1	-3	-3 + 1 = -2					
2	-1	-1 + 3 = 2					
3	1	(-3 - 1) 1 = -4					
4	3	(-1 - 3) (-j) = j4					
- 4	5	$5 + (-j5) = 5\sqrt{2}e^{-j\pi/4}$					
- 3	$5 e^{-j\pi/4}$	$5 e^{-j\pi/4} + 5 e^{-j3\pi/4} = -j5 2$					
-2	-j5	$(5 - (-j5)) 1 = 5 \sqrt{2}e^{j\pi/4}$					
- 1	$5 e^{-j3\pi/4}$	$(5 e^{-j\pi/4} - 5 e^{-j3\pi/4}) (-j) = -j5 2$					

	Results of the third stage					
Input	Stage 1	Stage 2	Stage 3 (Output)			
1	-3	-2	-2 + 2 = 0			
2	-1	2	(-2 - 2) 1 = -4			
3	1	-4	$-4 + j4 = -4 + j4 = 4 - 2e^{j3\pi/4}$			
4	3	j4	$(-4, -j4)$ 1 = $-4 - j4 = 4$ $2e^{-j3\pi/4}$			
- 4	5	$5 \frac{1}{2} e^{-j\pi/4}$	$5\sqrt{2}e^{-j\pi/4} + (-j5\sqrt{2}) = 5 - j12.07$			
- 3	$5 e^{-j\pi/4}$	_j5 𝔄	$(5 \mathcal{Q}e^{-j\pi/4} - (-j5 2)) 1 = 5 + j2.07$			
- 2	—j5	$5 2e^{ j\pi/4}$	5 $2e^{j\pi/4}$ + (-j5 2) = 5 - j2.07			
- 1	$5 e^{-j3\pi/4}$	-j5 2	$(5 \ 2e^{j\pi/4} - (-j5 \ 2)) \ 1 = 5 + j12.07$			

The DFT is $X(k) = \{0, (5 - j12.07), (-4 + j4), (5 - j2.07), -4, (5 + j2.07), (-4 - j4), (5 + j12.07)\}$

The MATLAB progarm is the same as shown in Example 1.

(DIT Template) The elementary computation (Butterfly):



The signal flow graph:



(DIF Template) The elementary computation (Butterfly):



The signal flow graph:





Inverse DFT using the FFT algorithm

The inverse DFT of an *N*-point sequence $\{X(k), k = 1, 2, ..., (N-1)\}$ is defined as

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k) W_{N}^{-kn}, \qquad n = 0, 1, \dots, N-1$$

where $W_N = e^{-j 2\pi/N}$. Take the complex conjugate of x(n) and multiply by N to get

$$N x^{*}(n) = \sum_{k=0}^{N-1} X^{*}(k) W_{N}^{kn}$$

The right hand side of the above equation is simply the DFT of the sequence $X^*(k)$ and can be computed by using any FFT algorithm. The desired output sequence is then found by taking the conjugate of the result and dividing by $N_{\lambda *}$

$$x(n) = \frac{1}{N} \left(\sum_{k=0}^{N-1} X^{*}(k) W^{kn} \right)^{*}$$

Example 2.4.1 Given the DFT sequence $X(k) = \{0, (-1-j), j, (2+j), 0, (2-j), -j, (-1+j)\}$ obtain the IDFT x(n) using the DIF FFT algorithm.

Solution This is an 8-point IDFT. The 8-point twiddle factors are, as calculated earlier,

$$\begin{array}{c} 0 \\ W_8 = 1 \\ W_8 = \left(e^{-j2\pi/8}\right)^2 = e^{-j\pi/2} = -j \\ \end{array} \qquad \begin{array}{c} 1 \\ W_8 = e \\ = -j \\ W_8 = \left(e^{-j2\pi/8}\right)^3 = e^{-j\pi/2} = -j \\ W_8 = \left(e^{-j\pi/8}\right)^3 = e^{-j\pi/8} = -j \\$$

The elementary computation (Butterfly) is shown below:



The signal flow graph follows:



	8-point IDFT using DIF FFT						
	Results of the first stage						
Input	Stage 1	Stage 2	Stage 3 (Output)				
$X^{*}(k)$							
0	0 + 0 = 0						
<u>-1+j</u>	-1+j+2+j=1+j2						
—j	-j + j = 0						
2—ј	2-j + (-1-j) = 1-j2						
0	(0-0) 1 = 0						
2+j	$(-1+j-(2+j)) e^{-j\pi/4} = -3 e^{-j\pi/4}$						
j	(-j - j) (-j) = -2						
—1—j	$(2-j-(-1-j)) e^{-j3\pi/4} = 3 e^{-j3\pi/4}$						

	Results of the second stage						
Input	Stage 1	Stage 2	Stage 3 (Output)				
0	0	0 + 0 = 0					
-1+j	1+j2	1+j2 + 1-j2 = 2					
-j	0	(0-0) 1 = 0					
2—ј	1-j2	(1+j2-(1-j2))(-j) = 4					
0	0	0 + (-2) = -2					
2+j	$-3 e^{-j\pi/4}$	$-3 e^{-j\pi/4} + 3 e^{-j3\pi/4} = -3 2$					
j	-2	(0 - (-2)) 1 = 2					
—1—j	$3 e^{-j3\pi/4}$	$(-3 e^{-j\pi/4} - 3 e^{-j3\pi/4}) (-j) = 3 2$					

	Results of the third stage				
Input	Stage 1	Stage 2	Stage 3 (Output)		
0	0	0	0 + 2 = 2	$\leftarrow 8 x^{*}(0)$	
-1+j	1+j2	2	(0-2) 1 = -2	$\leftarrow 8 x^{*}(4)$	
—j	0	0	0 + 4 = 4	$\leftarrow 8 x^{*}(2)$	
2-ј	1-ј2	4	(0-4) 1 = -4	$\leftarrow 8 x^*(6)$	
0	0	-2	$-2 + (-3\sqrt{2}) = -6.24$	$\leftarrow 8 x^*(1)$	
2+j	$-3 e^{-j\pi/4}$	-3 ⋧	$(-2 - (-3\sqrt{2})) 1 = 2.24$	$\leftarrow 8 x^*(5)$	
j	-2	2	2 + 3 2 = 6.24	$\leftarrow 8 x^*(3)$	
-1-j	$3 e^{-j3\pi/4}$	32	(2-32)1=-2.24	$\leftarrow 8 x^{*}(7)$	

The output at stage 3 gives us the values $\{8 x^*(n)\}$ in bit-reversed order:

$$\left\{8x^{*}(n)\right\}_{bit\ rev\ order} = \{2, -2, 4, -4, -6.24, 2.24, 6.24, -2.24\}$$

The IDFT is given by arranging the data in normal order, taking the complex conjugate of the sequence and dividing by 8:

$$\left\{8\overset{*}{x}\right\}_{normalorder} = \left\{2, -6.24, 4, 6.24, -2, 2.24, -4, -2.24\right\}$$

Note Because of the conjugate symmetry of $\{X(k)\}$, we should expect the sequence $\{x(n)\}$ to be *real-valued*.

The MATLAB program:

$$X = [0, (-1-j), j, (2+j), 0, (2-j), -j, (-1+j)], x = ifft(X)$$

Example 2.4.2 Given the DFT sequence $X(k) = \{0, (1-j), j, (2+j), 0, (2-j), (-1+j), -j\}$ obtain the IDFT x(n) using the DIF FFT algorithm. **Solution** There is no conjugate symmetry in $\{X(k)\}$. Using MATLAB

$$X = [0, 1-1j, 1j, 2+1j, 0, 2-1j, -1+1j, -1j]$$

x = ifft(X)

The IDFT is

$$x(n) = \{0.5, (-0.44 + 0.037i), (0.375 - 0.125i), (0.088 + 0.14i), (-0.75 + 0.5i), (0.44 + 0.21i), (-0.125 - 0.375i), (-0.088 - 0.39i)\}$$

FFT with general radix

As mentioned in the introduction, if the number of points, N, can be expressed as $N = r^m$, and if the computation algorithm is carried out by means of a succession of *r*-point transforms, the resultant FFT is called a **radix-***r* **algorithm**. In a radix-*r* FFT, an **elementary computation** (*EC*) consists of an *r*-point DFT followed by the multiplication of the *r* results by the appropriate twiddle factor. The number of *EC*s required is

$$C_r = \frac{N}{r} \log_r N$$

which decreases as r increases.

Of course, the complexity of an *EC* increases with increasing *r*. For r = 2, the *EC* (the butterfly) consists of a single complex multiplication and two complex additions; for r = 4, the *EC* requires three complex multiplications and several complex additions.

Suppose that we desire an *N*-point DFT where *N* is a composite number that can be factored into the product of integers

$$N = N_1 N_2 \ldots N_m$$

If, for instance, N = 64 and m = 3, we might factor N into the product $64 = 4 \times 4 \times 4$, and the 64-point transform can be viewed as a three-dimensional $4 \times 4 \times 4$ transform.

If *N* is a prime number so that factorization of *N* is not possible, the original signal can be *zero-padded* and the resulting new composite number of points can be factored.

We illustrate in the table below the situation for N = 64. Since $64 = 2^6$, we can have a radix-2 FFT; alternatively, since $64 = 4^3$, we can also have a radix-4 FFT.

$N = 64 = 2^6 = 4^3 = 8^2$					
Radix-2 Radix-4 Radix-8					
No. of stages	$\log_2 64 = 6$	$\log_4 64 = 3$	$\log_8 64 = 2$		
No. of <i>EC</i> s per stage	64/2 = 32	64/4 = 16	64/8 = 8		

Digital Signal Processing – 3

III. IIR Digital Filters

Analog filter approximations – Butterworth and Chebyshev, Design of IIR digital filters from analog filters, Bilinear transformation method, Step and Impulse invariance techniques, Spectral transformations, Design examples: Analog-Digital transformations

Contents:

Introduction The normalized analog, low pass, Butterworth filter Time domain invariance Bilinear transformation Nonlinear relationship of frequencies in bilinear transformation Digital filter design – The Butterworth filter Analog design using digital filters Frequency transformation The Chebyshev filter The Elliptic filter

Introduction

Nomenclature With $a_0 = 1$ in the linear constant coefficient difference equation,

$$a_0 y(n) + a_1 y(n-1) + \dots + a_N y(n-N) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M), \qquad a_0 \neq 0$$

we have,

$$H(z) = \frac{\sum_{i} b^{i} z^{-i}}{1 + \sum_{i=1}^{N} a_{i} z^{-i}}$$

М

This represents an IIR filter if at least one of a_1 through a_N is nonzero, and all the roots of the denominator are not canceled exactly by the roots of the numerator. In general, there are M finite zeros and N finite poles. There is no restriction that M should be less than or greater than or equal to N. In most cases, especially digital filters derived from analog designs, $M \le N$. Systems of this type are called N^{th} order systems. This is the case with IIR filter design in this Unit.

When M > N, the order of the system is no longer unambiguous. In this case, H(z) may be taken to be an N^{th} order system in cascade with an FIR filter of order (M - N).

When N = 0, as in the case of an FIR filter, according to our convention the *order* is 0. However, it is more meaningful in such a case to focus on *M* and call the filter an FIR filter of *M* stages or (M+1) coefficients.

Example The system $H(z) = (1 - z^{-8}) / (1 - z^{-1})$ is not an IIR filter. Why (verify)?

IIR filter design An analog filter specified by the Laplace transfer function, $H_a(s)$, may be designed to either frequency domain or time domain specifications. Similarly, a digital filter, H(z), may be required to have either (1) a given frequency response, or (2) a specific time domain response to an impulse, step, or ramp, etc.

Analog design using digital filters, $\omega_i = \Omega_i T$ Another possibility is that a digital filter may be required to simulate a continuous-time (analog) system. To simulate an analog filter the discrete-time filter is used in the A/D - H(z) - D/A structure shown below. The A/D converter can be thought of roughly as a sampler and coder, while the D/A converter, in many cases, represents a decoder and holder followed by a low pass filter (*smoothing filter*). The A/D converter may be preceded by a low pass filter, also called an *anti-aliasing filter* or *pre-filter*.



We will usually be given a set of analog requirements with critical frequencies $\Omega_1, \Omega_2, ..., \Omega_N$ in radians/sec., and the corresponding frequency response magnitudes $K_1, K_2, ..., K_N$ in dB. The sampling rate 1/T of the A/D converter will be specified or can be determined from the input signals under consideration. The general approach for the design is to first convert the analog requirements to digital requirements and then design the digital filter using the bilinear transformation. The conversion of the analog specifications to digital specifications is through the formula $\omega_i = \Omega_i T$. To show that this is true, suppose that the input to the equivalent analog filter is $x_a(t) = \sin \Omega_i t$. The output of the *A/D* converter with sampling rate 1/T becomes

 $x(n) = x_a(nT) = \sin\Omega_i nT = \sin(\Omega_i T)n = \sin\omega_i n$

Thus, the magnitude of the discrete-time sinusoidal signal is the same as the continuous time sinusoid, while the digital frequency ω_i is given in terms of the analog frequency Ω_i by $\omega_i = \Omega_i T$.

Thus, the specifications for the digital filter become $\omega_1, \omega_2, ..., \omega_N$ with the corresponding frequency response magnitudes $K_1, K_2, ..., K_N$. The digital frequency, ω , is in units of radians. The procedure is conceptually shown in figure below.



There are various techniques for designing H(z):

- 1. Numerical approximation (numerical solution) to the derivative operation or the integration operation (this latter results in the **bilinear transformation** aka bilinear *z*-transformation BZT).
- 2. Time domain invariance, e.g., impulse invariance and step-invariance methods.

The focus is on the low pass analog filter because once designed it can be transformed into an equivalent quality high pass, band pass or band stop filter by frequency transformation. The Butterworth, Chebyshev and elliptic filters are used as a starting point in designing digital filters. We approximate the magnitude part of the frequency response, not the phase. Butterworth and Chebyshev filters are actually special cases of the more difficult elliptic filter.

Because a constant divided by an N^{th} order polynomial in Ω falls off as Ω^N it will be an approximate low pass function as Ω varies from 0 to ∞ . Therefore, an all-pole analog filter H(s) = 1/D(s) is a good and simple choice for a low pass filter form and is used in both the Butterworth and the type I Chebyshev filters. Moreover, for a given denominator order, having the numerator constant (order zero) gives (for a given number of filter coefficients) the maximum attenuation as $\Omega \to \infty$.

The normalized analog, low pass, Butterworth filter

As a lead-in to digital filter design we look at a simple analog low pass filter – an RC filter, and its frequency response.

Example 3.2.1 Find the transfer function, $H_a(s)$, impulse response, $h_a(t)$, and frequency response, $H_a(j\Omega)$, of the following system.

$$R = 10k\Omega$$
+
$$+ \qquad +$$
Input = $x(t)$

$$C = 20\mu F$$
Output = $y(t)$

$$- \qquad -$$

Solution This is a voltage divider. The transfer function is given by

$$\frac{H(s) = \underline{Y(s)}}{x(s)} = \frac{(1 sC)}{R + (1 sC)} = \frac{(1/RC)}{s + (1 rC)} = \frac{5}{s + 5}$$

Taking the inverse Laplace transform gives the impulse response,

$$h_a(t) = 5 \ e^{-5t} u(t)$$

The frequency response is

$$H_a(j\Omega) = H_a(s)|_{s=j\Omega} = \frac{5}{j\Omega+5} = \frac{1}{\sqrt{1+(\Omega/5)^2}} e^{-j\tan^{-1}(\Omega/5)}$$

The cut-off frequency is $\Omega_c = 5$ rad/sec. The gain at $\Omega = 0$ is 1.



The MATLAB plots of frequency response of $H_a(j\Omega) = 5 (j\Omega+5)$ are shown below. We use the function *fplot*. The analog frequency, Omega (Ω), extends from 0 to ∞ ; however, the plots cover the range 0 to 6π rad/sec.

subplot(2, 1, 1), fplot('abs(5/(5+j*Omega))', [-6*pi, 6*pi], 'k'); xlabel ('Omega, rad/sec'), ylabel('|H(Omega)|'); grid; title ('Magnitude') % subplot(2, 1, 2), fplot('angle(5/(5+j*Omega))', [-6*pi, 6*pi], 'k'); xlabel ('Omega, rad/sec'), ylabel('Phase of H(Omega)'); grid; title ('Phase')



If we adjust the values of the components R and C so that 1/RC = 1, we would have $H_a(s)$

 $= \frac{1}{s+1}$ which is a *normalized filter* with cut-off frequency $\Omega_c = 1$ rad/sec and gain of 1 at $\Omega = 0$.

Such a normalized LP filter could be transformed to another LP filter with a different cut-off frequency of, say, $\Omega_c = 10$ rad/sec by the *low pass to low pass transformation* $s \rightarrow (s/10)$. The transfer function then becomes

$$H_a(s) = \frac{1}{s+1} = \frac{10}{s+10}$$

which still has a dc gain of 1. The gain of this filter could be *scaled* by a multiplier, say, *K*, so that

$$H_a(s) = K \frac{10}{s+10}$$

which has a dc gain of *K* and a cut-off frequency of $\Omega_c = 10$ rad/sec.

The frequency response of the normalized filter $H_a(s) = 1/(s+1)$ is $H_a(j\Omega) = 1 \langle j\Omega+1 \rangle$. The corresponding MATLAB plots are shown below using the function *plot*. Omega is a *vector*, consequently we use "./" instead of "/" etc.

Omega = -6*pi: pi/256: 6*pi; H = 1./(1.+ j .*Omega); subplot(2, 1, 1), plot(Omega, abs(H), 'k'); xlabel ('Omega, rad/sec'), ylabel('|H(Omega)|'); grid; title ('Magnitude') subplot(2, 1, 2), plot(Omega, angle(H), 'k'); xlabel ('Omega, rad/sec'), ylabel('Phase of H(Omega)'); grid; title ('Phase')



Butterworth filter The filter $H_a(s) = 5/(s + 5)$ is a first order Butterworth filter with cut-off frequency $\Omega_c = 5$ rad/sec. Its magnitude response is given by

 $|H(j\Omega)| = \frac{1}{\sqrt{1 + (\Omega/5)^2}}$

Omega = 0: pi/256: 15; H1 = 1./sqrt(1.+ (Omega/5) .^2); plot(Omega, H1, 'k'); legend ('1st Order, Cut-off = 5 rad/sec'); xlabel ('Omega, rad/sec'), ylabel('|H(Omega)|'); grid; title ('Magnitude')



Frequency response analysis The frequency response analysis in the analog frequency (Ω) domain is given by the following equations which may be used to illustrate, qualitatively, the effect of LP, HP or BP analog filters on a signal.

$$Y(s) = H(s) X(s)$$

$$Y(j\Omega) = H(j\Omega) X(j\Omega) = |H(j\Omega) e^{j \angle H(j\Omega)} X(j\Omega) e^{j \angle X(j\Omega)} = H(j\Omega) X(|j\Omega) e^{j(\angle H(j\Omega) + \angle X(j\Omega)}$$

$$|Y(\Omega) = |H(\Omega)| |X(\Omega) \text{ and } \angle Y(\Omega) = \angle H(\Omega) + \angle X(\Omega)$$

$$X(s) = H(s) = H(s)$$

Similarly, if we have a digital filter H(z) the frequency response analysis in the digital frequency (ω) domain is given by the following equations which may be used to illustrate, qualitatively, the effect of LP, HP or BP digital filters on a signal.

$$Y(z) = H(z) X(z)$$

$$Y(j\omega) = H(j\omega) X(j\omega) = |H(j\omega)| e^{j \angle H(j\omega)} X(j\omega) e^{j \angle X(j\omega)} = H(j\omega) X(|j\omega) e^{j (\angle H(j\omega) + \angle X(j\omega))}$$

$$|Y(\omega)| = |H(\omega)|X(\omega) \quad \text{and} \quad \angle Y(\omega) = \angle H(\omega) + \angle X(\omega)$$



The N^{th} order Butterworth filter In general the magnitude response of the N^{th} order Butterworth filter with cut-off frequency Ω_c is given by

$$|H(j\Omega)| = \frac{1}{\sqrt{1 + \left(\Omega / \Omega_c \gamma^2 N\right)}}$$

With the normalized frequency variable defined as $r = \Omega / \Omega_c$, the MATLAB segment below plots the magnitude response for 1st, 3rd, and 10th order filters, that is, N = 1, 3, and 10, respectively. Note that as the filter order increases the response becomes flatter on either side of the cut-off; and the transition (cut-off) becomes sharper.

r = 0: 0.1: 3; H1 = $1./sqrt(1.+r.^2)$; H3 = $1./sqrt(1.+r.^6)$; H10 = $1./sqrt(1.+r.^20)$; plot (r, H1, r, H3, 'r', r, H10, 'k'); legend ('1st Order', '3rd Order', '10th Order'); xlabel ('Normalized frequency, r'), ylabel('|H(r)|'); grid; title ('Magnitude')



Writing down the N^{th} order filter transfer function H(s) from the pole locations Let us look at some analog Butterworth filter theory.

- (1) $|H(j\Omega)| = \frac{1}{\sqrt{1 + (\Omega/\Omega_c)^{2N}}}$ decreases monotonically with Ω . No ripples.
- (2) Poles lie on the unit circle in the *s*-plane (for the Chebyshev filter, in contrast, they lie on an ellipse).
- (3) The transition band is wider (than in the case of the Chebyshev filter).
- (4) For the same specifications the Butterworth filter has more poles (or, is of higher order) than the Chebyshev filter. This means that the Butterworth filter needs more components to build.



The normalized analog Butterworth filter has a gain of $|H(j\Omega)| = 1$ at $\Omega = 0$, and a cut-off frequency of $\Omega_c = 1$ rad/sec. Given the order, N, of the filter we want to be able to write down its transfer function from the pole locations on the Butterworth circle.

Example 3.2.2 Given the order *N* of the filter, divide the unit circle into 2*N* equal parts and place poles on the unit circle at $(360^0/2N)$ apart. The H(s) will be made up of the *N* poles in the left half plane only. Remember complex valued poles must occur as complex conjugate pairs. There will



be no poles on the imaginary axis. Since the N poles must lie on the left half semicircle, when N is odd the odd pole must be at s = -1. Thus, for N = 1, there is one pole, $s_1 = -1$, and H(s) is given by

$$H(s) = \frac{1}{s - s_1} = \frac{1}{s - (-1)} = \frac{1}{s + 1}$$

Example 3.2.3 Filter order N = 2 so that 2N = 4 and $360^{0}/4 = 90^{0}$. The pole plot is shown above. The poles are at



so that

$$H(s) = \frac{1}{(s-s_1)(s-s_2)} = \frac{1}{\left[\begin{vmatrix} s+k\\ \sqrt{2} \end{vmatrix}} - \frac{1}{\left(\sqrt{2}\right)} \begin{vmatrix} 1\\ \sqrt{2} \end{vmatrix}} \begin{vmatrix} s+k\\ \sqrt{2} \end{vmatrix} - \frac{1}{\left(\sqrt{2}\right)} \begin{vmatrix} 1\\ \sqrt{2} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} \begin{vmatrix} 1\\ \sqrt{2} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} \begin{vmatrix} 1\\ \sqrt{2} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} \begin{vmatrix} 1\\ \sqrt{2} \end{vmatrix} + \frac{1}{\left(\sqrt{2}\right)} + \frac{1}{\left(\sqrt{$$

Denominator is

Dr. =
$$\begin{pmatrix} -1 \\ s + \frac{1}{\sqrt{2}} \\ - \frac{$$

So

$$H(s) = \frac{1}{s^2 + \sqrt{2}s + 1}$$

Example 3.2.4 Filter order N = 3, so that 2N = 6 and $360^{\circ}/6 = 60^{\circ}$. Poles are at $\begin{vmatrix} 1 & \frac{3}{\sqrt{2}} \end{vmatrix}$ $s_{1, 2, 3} = -1$, $\begin{vmatrix} -2 + j & 2 \\ 2 & 2 & 2 \end{vmatrix}$, and $\begin{vmatrix} -2 - j & 2 \\ 2 & 2 & 2 \end{vmatrix}$ $s^2 + \sqrt{2} + s + 1$



$$H(s) = \frac{1}{(s+1)\left(s+\frac{1}{2}-j\frac{\sqrt{3}}{2}\right)\left(s+\frac{1}{2}+j\frac{\sqrt{3}}{2}\right)}$$

Denominator is $\begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$$Dr. = (s+1) | | s+|^{2} | s^{2} | = (s+1)^{2} + s+1$$

$$| | (s) = 1$$

So

$$H(s) = \frac{1}{(s+1)(s^2+s+1)}$$

Example 3.2.5 Filter order N = 4, so that 2N = 8 and $360^{0}/8 = 45^{0}$. Poles are at

$$s = (-\cos 22.5^{\circ} \pm j \sin 22.5^{\circ}) = (-\cos \alpha \pm j \sin \alpha) (-\cos 67.5^{\circ} \pm j \sin 67.5^{\circ}) = (-\cos \beta \pm j \sin \beta)$$



Determining the order and transfer function from the specifications A typical *magnitude* response specification is sketched below. The magnitudes at the critical frequencies Ω_1 and Ω_2 are *A* and *B*, respectively. Typically Ω_1 is in the pass band or is the edge of the pass band and Ω_2 is in the stop band or is the edge of the stop band. For illustrative purposes we have arbitrarily



taken A = 0.707 (thus Ω_I is the cut-off frequency, but this need not be the case) and B = 0.25. The *log-magnitude* specification is diagrammed below. Note that $(20 \log A) = K_I dB$ and $(20 \log B) = K_2 dB$. Thus the analog filter specifications are

> $0 \ge 20\log_{10}|H(j\Omega)| \ge K_1 \text{ for all } \Omega \le \Omega_1$ $20\log_{10}|H(j\Omega)| \le K_2 \text{ for all } \Omega \ge \Omega_2$



With the magnitude $|H(j\Omega)|$ given by the Butterworth function,

$$|H(j\Omega)| = \frac{1}{\sqrt{1 + (\Omega/\Omega_c)^{2N}}}$$

and using the equality condition at the critical frequencies in the above specifications the order, N, of the filter is given by

$$N = \begin{bmatrix} \left[10^{-K_{1}/10} - 1 \right] \\ 10g_{10} \\ 10^{-K_{2}/10} - 1 \\ 10^{-K_{2}/10} \\ 10^{-K_$$

The result is rounded to the next larger integer. For example, if N = 3.2 by the above calculation then it is rounded up to 4, and the order of the required filter is N = 4. In such a case the resulting filter would exceed the specification at both Ω_1 and Ω_2 . The cut-off frequency Ω_c is determined from one of the two equations below.

$$\Omega_c = \frac{\Omega_1}{\sqrt[2N]{10^{-K_1/10} - 1}} \quad \text{or} \quad \Omega_c = -\frac{\Omega_2}{\sqrt[N]{-1}} \rightarrow (3.17 \text{ Ludeman})$$

The equation on the left will result in the specification being met exactly at Ω_1 while the specification is exceeded at Ω_2 . The equation on the right results in the specification being met exactly at Ω_2 and exceeded at Ω_1 .

Note that the design equation for N may be written in the alternative form

$$N = \left| \frac{10^{-K_{2}/10} - 1}{10^{-K_{1}/10} - 1} \right|$$

Example 3.2.6 What is the order and transfer function of the analog Butterworth filter that satisfies the following specification?

$$\Omega_1 = 200 \text{ rad/sec}$$
 $K_1 = -1 \text{ dB}$
 $\Omega_2 = 600 \text{ rad/sec}$ $K_2 = -30 \text{ dB}$

Solution The order
$$N$$
 is given by
 $|\log_{10}| \left(\frac{10^{-K_{1}/10} - 1}{10^{-K_{2}/10} - 1} \right)| = \left| \frac{\log_{10}\left(\frac{10^{-(-1)/10} - 1}{10^{-(-30)/10} - 1} \right)}{2\log_{10}\left(\frac{200}{10} \right)|} \right| = \left| \frac{\log_{10}\left(\frac{10^{0.1} - 1}{10^{3} - 1} \right)}{2\log_{10}\left(\frac{200}{2\log_{10}} \right)|} \right|$
 $N = \left| \frac{2\log\left(\Omega\right)}{2\log\left(\Omega\right)} \right|^{\frac{1}{2}} \left| \frac{1}{2} \right| \frac{1}{2} \left| \frac{1}{2} \right| \frac{1}{2} \left| \frac{1}{2} \right| \frac{1}{2} \left| \frac{1}{2} \right| \frac{1}{2} \left| \frac{1}{2} \right| \frac{1}{2} \left| \frac{1}{2} \left| \frac{1}{2} \left| \frac{1}{2} \left| \frac{1}{2} \left| \frac{1}{2} \left| \frac{1}{2} \right| \frac{1}{2} \left| \frac{1}{$

Now, as in an earlier example, locate 8 poles uniformly on the unit circle, making sure to satisfy all the requirements ... and write down the transfer function, H(s), of the normalized Butterworth filter (with a cut-off frequency of 1 rad/sec),

$$H(s) = \frac{1}{(s^2 + 1.848 s + 1)(s^2 + 0.765 s + 1)}$$

Next, determine the cutoff frequency Ω_c that corresponds to the given specifications and the order N = 4 determined above

$$\Omega_{c} = \frac{\Omega_{1}}{\sqrt[2^{N}]{10^{-K_{1}/10} - 1}} = \frac{200}{\sqrt[2^{(4)}]{10^{-(-1)/10} - 1}} = \frac{200}{\sqrt[8]{1.26 - 1}} = \frac{200}{0.8446} = 236.8 \text{ rad/sec}$$

Finally, we make the substitution $s \rightarrow (s/236.8)$ in H(s) and thereby move the cutoff frequency from 1 rad/sec to 236.8 rad/sec resulting in the transfer function $H_a(s)$

$$=\frac{1}{((s/236.8)^{2}+1.848(s/236.8)+1)((s/236.8)+1)((s/236.8)^{2}+0.765(s/236.8)+1)((s/236.8)^{2}+0.765(s/236.8)+1)((s/236.8)^{4}-(236.8)^{4})(s^{2}+1.848(236.8)+236.8^{2})(s^{2}+0.765(2326.8)+236.8^{2})=\dots$$

(Aside) The more general N^{th} order Butterworth filter has the magnitude response given by

$$\left|H\left(j\Omega\right)\right| = \frac{1}{\sqrt{1 + \varepsilon^{2}\left(\Omega/\Omega_{1}\right)^{N}}}$$

The parameter \langle has to do with pass band attenuation and Ω_l is the pass band edge frequency (not necessarily the same as the 3 dB cut-off frequency Ω_c). MATLAB takes $\langle = 1$ in which case $\Omega_l = \Omega_c$. See DSP-HW. [See Cavicchi, Ramesh Babu]. (*End of Aside*)

Time domain invariance

Given an analog filter's response to a specific input we require that the response of the digital filter (to be designed) to the digital version of the analog input should be the same as the analog response at sampling instants. If the input is an impulse function the corresponding design is called an impulse invariant design, if the input is a step function the corresponding design is called a step invariant design.

Impulse-invariant design If $h_a(t)$ represents the response of an analog filter $H_a(s)$ to a unit impulse $\delta(t)$, then the unit sample response of a discrete-time filter used in an A/D - H(z) - D/A structure is selected to be the sampled version of $h_a(t)$. That is, we are preserving the response to an impulse. Therefore the discrete-time filter is characterized by the system function, H(z), given

$$H(z) = \mathbf{z}\{h(n)\} = \mathbf{z}\left[\left\{h_a(t)\right\}\Big|_{t=nT}\right]$$

If we are given an analog filter with system function $H_a(s)$ the corresponding impulse-invariant digital filter, H(z), is seen from above to be

$$H(z) = \Im \left[\left(L^{-1} \{ H_a(s) \} \right) \right]_{a=nT}, \text{ where } L^{-1} \text{ means Laplace inverse}$$

Note that at this point we have not specified how $H_a(s)$ was obtained, but rather we have shown how to obtain the digital filter H(z) from any given $H_a(s)$ using impulse invariance.

Example 5.3.1 [Low pass filter] For the analog filter $H_a(s) = \frac{A}{s + \alpha}$ find the H(z) corresponding

to the impulse invariant design using a sample rate of 1/T samples/sec. Solution The analog system's impulse response is $h(t) = \mathcal{L}^{-1} | = Ae^{-\alpha t}u(t)$. The

 $\int s + \alpha$

corresponding h(n) is then given by

$$h(n) = h_{a}(t)|_{t=nT} = Ae^{-\alpha nT}u(nT) = A(e^{-\alpha T})^{n}u(n) = Aa^{n}u(n)$$

where, as previously, we have set $e^{-\alpha T} = a$. The discrete-time filter, then, is given by the *z*-transform of h(n)

$$H(z) = \mathbf{z} \{h(n)\} = \mathbf{z} \{h(e^{-\alpha T})^n u(n)\} = \frac{Az}{z - e^{-\alpha T}} = \frac{Az}{z - a}$$
$$= \frac{A}{1 - e^{-\alpha T} z^{-1}} = \frac{A}{1 - a z^{-1}}$$

which has a pole at $z = e^{-\alpha T} = a$. In effect, the pole at $s = -\alpha$ in the *s*-plane is mapped to a pole at $z = e^{-\alpha T} = a$ in the *z*-plane. (**HW** What is the difference equation?)

$$\frac{Y(z)}{X(z)} = \frac{A}{1 - az^{-1}} \longrightarrow Y(z)(1 - az^{-1}) = A X (z)$$

$$\rightarrow \qquad y(n) - ay(n-1) = Ax(n) \longrightarrow \qquad y(n) = Ax(n) + ay(n-1)$$

Relationship between the *s***-plane and the** *z***-plane** ($z = e^{sT}$) We can extend the above procedure to the case where $H_a(s)$ is given as a sum of *N* terms with distinct poles as

$$H_a(s) = \sum_{k=1}^{\infty} \frac{\Box_k}{s+\alpha}$$

For this case the impulse invariant design H(z) is given by

$$H(z) = \mathbf{z} \Big| \begin{bmatrix} -1 & A & | & | \\ -1 & \sum_{k=1}^{k} a_{k} \\ | & L & \sum_{k=1}^{k} a_{k} \\ | & k \\ | & k$$

where L^{-1} means Laplace inverse. We observe that a pole at $s = -\alpha_k$ in the *s*-plane gives rise to a pole at $z = e^{-\alpha_k T}$ in the *z*-plane and the coefficients in the partial fraction expansion of $H_a(s)$ and H(z) are equal. If the analog filter is stable, corresponding to $-\alpha_k$ being in the left half plane, then the magnitude of $e^{-\alpha_k T}$ will be less than unity, so that the corresponding pole of the digital filter is inside the unit circle, and as a result the digital filter also is stable.

While the poles in the *s*-plane "map" to poles in the *z*-plane according to the relationship $z = e^{sT}$, it is important to recognize that the impulse invariance design procedure does not correspond to a *mapping (transformation)* of the *s*-plane to the *z*-plane by that relationship or in fact by any relationship. (An example of a transformation is where we actually make a $2(1-z^{-1})$

substitution, say, $s = \overline{T} \left| \frac{1}{1+z^{-1}} \right|$ which, of course, is the bilinear transformation). For example,

the zeros of $H_a(s)$ do not map to zeros of H(z) according to this relation. See also matched *z*-transform later.

We can explore the relationship $z = e^{s^T}$ keeping in mind that it only applies to poles and that it is not a transformation. Set $s = \zeta + j\Omega$ and $z = r e^{j\omega}$ in $z = e^{s^T}$ to get $r e^{j\omega} = e^{(\sigma+j\Omega)^T} = e^{\sigma T} e^{j\Omega T}$ so that $r = e^{\sigma T}$ and $\omega = \Omega T$.

The above relations can be used to show that poles in the left half of the primary strip in the *s*-plane map into poles within the unit circle in the *z*-plane as shown in the figure for $s = s_1$.

Mapping of poles, $z = e$					
s-plane pole $s = \zeta + iQ$	<i>z</i> -plane pole $\int_{a}^{sT} dx = \int_{a}^{j\omega}$	r	ω		
$3-\zeta+Jsz$	$z = e = r e^s$				
0	1	1	0		
jΩ _s /2	-1	1	π		
$-\infty + j\Omega_s/2$	-0	0	π		
$-\infty - j\Omega_s/2$	-0	0	π		
$-j\Omega_{s}/2$	-1	1	π		

Mapping of poles, $z = e^{sT}$

For $s_1 = \zeta_1 + j\Omega_1$ we have $r = e^{\sigma_1 T}$ and $\omega = \Omega_1 T$. However, poles at s_2 and s_3 (which are a distance Ω_s from s_1) also will be mapped to the same pole that s_1 is mapped to. In fact, an infinite number of *s*-plane poles will be mapped to the same *z*-plane pole in a many-to-one relationship. These frequencies differ by $\Omega_s = 2\pi F_s = 2\pi/T$ (F_s is the sampling frequency in Hertz). This is called **aliasing (of the poles)** and is a drawback of the impulse-invariant design. The analog system poles will not be aliased in this manner if, in the first place, they are confined to the "primary strip" of width $\Omega_s = 2\pi F_s = 2\pi/T$ in the *s*-plane.

In a similar fashion poles located in the right half of the primary strip in the *s*-plane will be mapped to the outside of the unit circle in the *z*-plane. Here again the mapping of the *s*-plane poles to the *z*-plane poles is many-to-one.



Owing to the aliasing, the impulse invariant design is suitable for the design of low pass and band pass filters but not for high pass filters.

(*Omit*) *Matched z-transform* In this method we apply the mapping $z = e^{sT}$ not only to the poles but also to the zeros of $H_a(s)$. As a result the observations made above are valid for the matched ransform method of filter design.

Frequency response of the equivalent analog filter Going back to the impulse invariant design of the first order filter, how does the frequency response of the A/D - H(z) - D/A structure using this H(z) compare to the frequency response of the original system specified by $H_a(s)$?



Note
$$H_a(s) = \frac{A}{s+\alpha}$$
 and $H(z) = \frac{Az}{z_1 - e^{-\alpha T}} = \frac{Az}{z_n - a}$ with $a = e^{-\alpha T}$. For the analog filter we have
 $H_a(j\Omega) = \frac{A}{s+\alpha} \Big|_{s=j\Omega} = \frac{Az}{j\Omega + \alpha} = \frac{A}{\sqrt{\alpha^2 + \Omega^2}} e^{-j\tan^{-1}(\Omega/\alpha)}$
 $|H_a(j\Omega)| = \frac{A}{\sqrt{\alpha^2 + \Omega^2}}, \quad -\infty < \Omega < \infty$

For future reference note that $|H_a(j0)| = A / \alpha$.

To obtain the equivalent frequency response of the A/D - H(z) - D/A structure one must first find the frequency response of the discrete-time filter specified by H(z). This is given by $Ae^{j\omega}$

$$H(e^{j\omega}) = H(z)]_{z=e^{j\omega}} = \frac{\pi}{e^{j\omega} - e^{-\alpha T}}, \quad -\pi < \omega < \pi \text{ (because periodic)}$$

The analog frequency response of the *equivalent* analog filter is then determined by replacing ω by ΩT . Note, however, that since the digital frequency, ω , is restricted to $(-\pi, \pi)$, the analog frequency, Ω , is correspondingly is restricted to $(-\pi/T, \pi/T)$. We get

$$H_{eq}(j\Omega) = H(e^{j\omega}) \Big]_{e^{-\Omega T}} = \frac{Ae^{j\Omega T}}{e^{j\Omega T} - e^{-\alpha T}} = \frac{A}{1 - e^{-\alpha T}e^{-j\Omega T}}, \qquad \Omega T < \pi \text{ or } \Omega < \pi/T$$

 π/T

Denominator = $1 - e^{-\alpha T} e^{-j\Omega T} = 1 - e^{-\alpha T} (\cos \Omega T - j \sin \Omega T) = (1 - e^{-\alpha T} \cos \Omega T) + j e^{-\alpha T} \sin \Omega T$

So

$$H_{eq}(j\Omega) = \frac{A}{(1 - e^{-\alpha T} \cos \Omega T) + je^{-\alpha T} \sin \Omega T}, \qquad \Omega <$$

$$|H_{eq}(j\Omega)| = \frac{A}{\sqrt{1 + e^{-2\alpha T} - 2e^{-\alpha T} \cos \Omega T}}, \qquad \Omega < \pi/T$$

$$(j0)| = \frac{A}{1 - e^{-\alpha T}} = \frac{A}{1 - a}$$

Note that $/H_{e}(j0)/=$

We can plot $|H_{eq}(j\Omega)|$ and $|H_a(j\Omega)|$ for, say, $\alpha = 1$ and different values of T say T = 0.1and T = 1, remembering that $|H_{eq}(j\Omega)|$ is periodic, the basic period going from $-\pi/T < \Omega < \pi/T$. Ideally the two plots should be very close (in shape, over the range of frequencies of interest) but it will be found that the smaller the value of T, the closer the two plots are. Thus T = 0.1 will result in a closer match than T = 1. Therefore, using the impulse invariant design, good results are obtained provided the time between samples (T) is selected small enough. What is small enough may be difficult to assess when the $H_a(s)$ has several poles; and when it is found, it may be so small that implementation may be costly. In general, other transformational methods such as the bilinear allow designs with sample rates that are less than those required by the impulse invariant method and also allow flexibility with respect to selection of sample rate size.



Example 3.3.2 [Impulse invariant design of 2nd order Butterworth filter] Obtain the impulse

invariant digital filter corresponding to the 2nd order Butterworth filter $H_a(s) = \frac{4}{s^2 + 2\sqrt{2s} + 4}$

with sampling time T = 1 sec.

Solution Note that we are using T = 1 sec. simply for the purpose of comparing with the bilinear design done later with T = 1 sec. It is important to remember that in bilinear design calculations the value of T is immaterial since it gets cancelled in the *design process*; but in impulse invariant design there is no such cancellation, so the value of T_4 is critical (the smaller, the better). H(s) = -

$$= \frac{3}{\sqrt{2}} = \frac{3}{\sqrt{2}} = 2\sqrt{2} + 4 = \frac{3}{\sqrt{2}} = 2\sqrt{2} + (\sqrt{2})^{2} + (\sqrt{2})^{2} = 2\sqrt{2} + (\sqrt{2})^{2} = 2\sqrt{2} + (\sqrt{2})^{2} = 2\sqrt{2} + (\sqrt{2})^{2} = 2\sqrt{2} + (\sqrt{2})^{2} + (\sqrt{2})^{2} + (\sqrt{2})^{2} = 2\sqrt{2} + (\sqrt{2})^{2} + (\sqrt{2})^{2} + (\sqrt{2})^{2} = 2\sqrt{2} + (\sqrt{2})^{2} +$$

The expression in braces is in familiar form and can be converted to its impulse invariant digital filter equivalent. See 3(c) in HW.

Step invariant design Here the response of the digital filter to the unit step sequence, u(n), is chosen to be samples of the analog step response. In this way, if the analog filter has good step response characteristics, such as small rise-time and low peak over-shoot, these characteristics would be preserved in the digital filter. Clearly this idea of waveform invariance can be extended to the preservation of the output wave shape for a variety of inputs.



(*Omit*) Problem Given the analog system $H_a(s)$, let $h_a(t)$ be its impulse response and let $p_a(t)$ be its step response. The system $H_a(s)$ is given to be continuous-time linear time-invariant. Let also
h(n) be the unit sample response, p(n) be the step response, and, H(z) be the system function,

of a discrete-time linear shift-invariant filter. Then,

(a) If
$$h(n) = h_a(nT)$$
, does $p(n) = \sum_{k=-\infty}^{n} h_a(kT)$?
(b) If $p(n) = p_a(nT)$, does $h(n) = h_a(nT)$?
Solution (a) If $h(n) = h_a(nT)$, does $p(n) = \sum_{k=-\infty}^{n} h_a(kT)$? We know that

$$u(n) = \sum_{k=-\infty}^{\infty} \delta(k)$$

This is seen to be true by writing it out in full as

$$u(n) = 6(-\infty) + 6(-\infty + 1) + \dots + 6(-1) + 6(0) + 6(1) + \dots + 6(n)$$

where *n* is implicitly some positive integer. Take, for instance, n = 3; then, from the above equation u(3) = 6(0) = 1, all the other terms being zero. In other words u(n) is a linear combination of unit sample functions. And, since the response to 6(k) is h(k), therefore, the response to u(n) is a linear combination of the unit sample responses h(k). That is,

$$p(n) = \sum_{k = -\infty}^{n} h(k) = \sum_{k = -\infty}^{n} h_a(kT)$$

Therefore, the answer to the above question is, Yes.

(b) If $p(n) = p_a(nT)$, does $h(n) = h_a(nT)$? Since $\delta(n) = u(n) - u(n-1)$, the response of the digital system, H(z), to the input $\delta(n)$ is

$$h(n) = p(n) - p(n-1) = p_a(nT) - p_a(nT-T) \neq h_a(nT)$$

Therefore, the answer to the above question is, No.

Example 3.3.3 [LP filter] [Step Invariance] Consider the continuous-

time system $H_a(s) = \frac{A}{s+\alpha}$ with unit step response $p_a(t)$. Determine the system function, H(z), i.e., the *z*-transform of the unit sample response h(n) of a discrete-time system designed from this

system on the basis of step-invariance, such that $p(n) = p_a(nT)$, where

$$p(n) = \sum_{k=-\infty}^{\infty} h(k)$$
 and $p_a(t) = \oint_{-\infty}^{1} a(\tau) d\tau$

Solution Since $p_a(t) = \oint_{-\infty} (\tau) d\tau$ we have

$$\mathcal{L}[p_a(t)] = P_a(s) = \frac{H_a(s)}{s} = \frac{A}{s(s+\alpha)} = \frac{K_1}{s} + \frac{K_2}{s+\alpha}$$

where K_1 and K_2 are the coefficients of the partial fraction expansion, given by

$$K_I = \frac{A}{s+\alpha}\Big|_{s=0} = \frac{A}{\alpha}$$
 and $K_2 = \frac{A}{s_{s=-\alpha}} = -\frac{A}{\alpha}$

Therefore,

from which we write $p(n) = p_a(nT)$ and hence P(z) etc. Equivalently, we may reason as follows. The correspondence between *s*-plane poles and *z*-plane poles is

However, what we need is H(z). Since $\delta(n) = u(n) - u(n-1)$, and we know the response to u(n), therefore, the response to u(n) - u(n-1) is given by h(n) = p(n) - p(n-1), and taking the ztransform of this last equation.

$$\begin{array}{l} \sum_{z = P(z) - z} \sum_{z} P(z) = (1 - z^{-1})P(z) = \frac{|z - 1|}{P(z)} \\ = \left| \frac{(z - 1)A}{(z)\alpha} \left(1 - e^{-\alpha T} \right) \frac{z}{(z - 1)(z - e^{-\alpha T})} \right. \\ = \frac{A}{\alpha} \left(1 - e^{-\alpha T} \right) \frac{1}{(z - e^{-\alpha T})} \end{array}$$

Alternatively, we may also obtain the transfer function as the ratio of output and input transforms, H(z) = P(z)/U(z).

Frequency-response analysis As we did in the case of the impulse-invariant design, here also we can compare $|H_a(j\Omega)|$ with $|H_{eq}(j\Omega)|$. For the system $H_a(s) = \frac{1}{s+\alpha}$ the frequency response is already evaluated as $|H_a(j\Omega)| = \frac{1}{\sqrt{\alpha^2 + \Omega^2}}$. We need the frequency response, $|H_{eq}(j\Omega)|$, of the equivalent analog filter when the above H(z) is used in a A/D - H(z) - D/A structure. Start with the above H(z) and set $z = e^{j\omega} t_0$ get $A = -\alpha_T = 1_{\alpha}$ $= A = -\alpha_{\overline{T}} 1_{\alpha}$ $\sum_{z=e}^{r} \left(1-e^{-1}\right)_{z=e^{i\omega}} \alpha\left(1-e^{-1}\right)_{z=e^{i\omega}} \alpha\left(1-e^{-1}\right)$ For the equivalent analog filter we get $H_{eq}(j\Omega)$ by setting $\omega = \Omega T$ in $H(e^{-j\omega})$.

Bilinear transformation

One approach to the *numerical solution* of an ordinary linear constant-coefficient differential equation is based on the application of the trapezoidal rule to the first order approximation of an integral (or integration). Consider the following equivalent pair of equations

$$\frac{dy}{dt} = x(t) \qquad \Rightarrow \qquad \int dy = \int x(t)dt$$

Here $\int dy = \text{area shown shaded and is given by}$ $y(n) - y(n-1) = \left(\frac{x(n) + x(n-1)}{2}\right)T$

where we have used the trapezoidal rule to compute the area under a curve.



Taking the *z* transform of the above we get

$$Y(z) - z^{-1} Y(z) = \frac{1}{2} X(z)(1 + z^{-1})$$

Rearranging terms gives
$$\frac{Y(z)}{X(z)} = \frac{T}{2} \left(\begin{vmatrix} 1 + z^{-1} \\ - z \end{vmatrix}^{-1} \right) = \frac{1}{2} \left(\begin{vmatrix} 1 - z^{-1} \\ - z \end{vmatrix}^{-1} \right) = \frac{1}{2} \left(\begin{vmatrix} 1 - z^{-1} \\ - z \end{vmatrix}^{-1} \right)$$

Thus the continuous time system $y(t) = \int x(t)dt$ represented by the following block diagrams



is replaced in the discrete-time domain by the following block.



In other words, given the Laplace transfer function $H_a(s)$, the corresponding digital filter is given by replacing s with $\frac{2(1-z^{-1})}{T(1+z^{-1})}$, or $H(z) = H(s) \Big]_{a}^{2(1-z^{-1})}$

This is called *bilinear transformation* (both numerator and denominator are first order polynomials), also known as *bilinear z-transformation* (BZT).

Here again, note that at this point we have not specified how $H_a(s)$ was obtained, but rather we are showing how to obtain the digital filter H(z) from any given $H_a(s)$ using bilinear transformation.

Example 3.4.1 [LP filter] [Bilinear] Design a digital filter based on the analog system $H_a(s)$ = A, using the bilinear transformation. Give the difference equation. Use T = 2 sec. $s + \alpha$

Solution

$$\frac{H(z) = H_{\alpha}(s) \int \underline{2(1-z^{-1})}_{s=T(1+z^{-1})} = \frac{A}{\underline{2(1-z^{-1})} + \alpha} = \frac{A}{(\alpha+1) + (\alpha-1)z^{-1}}$$
$$T(1+z^{-1})$$

Example 2 [Ludeman, p. 178] Apply bilinear transformation to the 2nd order Butterworth filter $H_a(s) = \frac{4}{s^2 + 2\sqrt{2s + 4}}$ with T = 1 sec. Obtain (1) H(z), (2) the difference equation and

(3) $H(e^{j\omega})$.

Solution
1. With
$$T = 1$$
 second, $s = \frac{2}{T} \left(\frac{1 - z^{-1}}{2} \right) = \frac{2(1 - z^{-1})}{1 + z^{-1}}$, and
 $T \left(\frac{1 - z^{-1}}{2} \right) = \frac{2(1 - z^{-1})}{1 + z^{-1}} = \frac{4}{1 + z^{-1}}$, and
 $H(z) = \frac{H}{a}(s) = \frac{2(1 - z^{-1})}{s = \frac{2(1 - z^{-1})}{1 + z^{-1}}} = \frac{4}{1 + 2(1 - z^{-1})} = \frac{4}{\sqrt{2}(1 + z^{-1})} = \frac{4}{\sqrt{2}(1 + z^{-1})} = \frac{1 + 2z^{-1} + z^{-2}}{1 + 2z^{-1} + z^{-2}} = -\frac{3.4142135 + 0.5857865 z}{1 + 2z^{-1} + z^{-2}}$

2. The difference equation is obtained from

$$H(z) = \frac{Y(z)}{X(z)} = \frac{1 + 2z^{-1} + z^{-2}}{3.4142135 + 0.5857865 z^{-2}}$$

Cross-multiply and take inverse z-transform to get

$$3.4142135 y(n) + 0.5857865 y(n-2) = x(n) + 2x(n-1) + x(n-2)$$

By rearranging and scaling y(n) can be realized as

$$y(n) = 0.2928932 [x(n) + 2 x(n-1) + (n-2)] - 0.1715729 y(n-2)$$

- 3. Frequency response
 - quency response $H(z)]_{z=e^{j\omega}} = \frac{1 + 2e^{-j\omega} + e^{-j2\omega}}{3.4142135 + 0.5857865e^{-j2\omega}} = \frac{N(\omega)}{D(\omega)}$ Numerator = $N(\omega) = 1 + 2e^{-j\omega} + e^{-j2\omega} = e^{-j\omega}(e^{-j\omega} + 2 + e^{-j\omega}) = e^{-j\omega}(2 + 2\cos\omega)$ $= 2(1 + \cos \omega) e^{-j\omega}$ Denominator = $D(\omega)$ = 3.4142135 + 0.5857865 $e^{-j2\omega}$ = A + B $e^{-j2\omega}$ $= (A + B \cos 2\omega) - i B \sin 2\omega$

$$H(e^{j\omega}) = \underbrace{2(1+\cos\omega)}_{i \text{ (a -1)}} \underbrace{2(1+\cos\omega)}_{i \text{ (b -1)}} \underbrace{2(1+\cos\omega)}_{i \text{ (b -1)}} \underbrace{2(1+\cos\omega)}_{i \text{ (c -1)}} \underbrace{2($$

 $\sqrt{(A+B\cos 2\omega)^2+(B\sin 2\omega)^2}$

The magnitude of the frequency response is

$$|H(e^{j\omega})| = \frac{2(1+\cos\omega)}{\sqrt{(A+B\cos 2\omega)^2 + (B\sin 2\omega)^2}}$$

Plot 20 $\log_{10}|H(e^{j\omega})|$ vs ω for $\omega = 0$ to π .

Relationship between the s- and z-planes The bilinear transformation is

$$s = \frac{2(1-z)}{\overline{T}(1+z^{-1})}$$
 or $z = \frac{1+(sT/2)}{1-(sT/2)}$

We can map a couple of points on the $j\Omega$ axis by setting $s = j\Omega$, so that

$$z = \frac{1 + (j\Omega T/2)}{1 - (j\Omega T/2)}$$

Thus $\Omega = 0$ maps to z = 1 and $\Omega = 2/T$ maps to $z = j1 = 1 e^{j\pi/2}$ as shown in figure below. The transformation has the following properties:

- 1. The entire $j\Omega$ axis of the *s*-plane goes on to the unit circle of the *z*-plane.
- 2. The entire left half of the *s*-plane is mapped into the inside of the unit circle of the *z*-plane. (In contrast, in impulse invariant design the *poles* in the left halves of infinitely many strips of width Ω_s in the *s*-plane are mapped into the inside of the unit circle in the *z*-plane.)

So a stable analog filter, with all of its poles in the left half plane, would be transformed into a stable digital filter with all of its poles in the unit circle. The frequency response is evaluated on the $j\Omega$ axis in the *s*-plane and on the unit circle in the *z*-plane. While the frequency responses of the analog filter and digital filter have the same amplitudes there is a nonlinear relationship between corresponding digital and analog frequencies.



Nonlinear relationship of frequencies in bilinear transformation

In the bilinear transformation the analog and digital frequencies are non-linearly related. Setting 1 - 7

$$s = {}^{j\Omega} \text{ and } z = e^{j\omega} \text{ in } s = \frac{2}{||} ||_{-\frac{1}{2}} ||_{+\infty} \text{ we get}$$

$$j\Omega = \frac{2}{||} \left(\frac{T + e^{-j\omega}}{1 + e^{-j\omega}} \right)^{-1} = 2 \frac{e^{-j\omega/2} (e^{j\omega/2} - e^{-j\omega/2})}{(e^{j\omega/2} - e^{-j\omega/2})} = \frac{2((e^{j\omega/2} - e^{-j\omega/2})/2)}{T - 2((e^{j\omega/2} + e^{-j\omega/2})/2)}$$
or
$$\Omega = \frac{T + e^{-j\omega}}{\sin \omega/2} = 2T e^{-i\omega/2} (e^{j\omega/2} + e^{-j\omega/2}) = \frac{T}{T} \frac{2((e^{j\omega/2} + e^{-j\omega/2})/2)}{2((e^{j\omega/2} + e^{-j\omega/2})/2)}$$

$$= \frac{T + e^{-j\omega}}{\sin \omega/2} = 0$$
or
$$\omega = 2 \tan |\frac{\Omega T}{2}|$$

$$T \cos \omega/2 = T + (2)$$

First we sketch Ω (the analog frequency) as a function of ω (the digital frequency) as given by $\Omega = \frac{2}{T} \tan \left(\frac{1}{2} \right)$ to show qualitatively the distortion of the frequency scale that occurs due to the nonlinear nature of the relationship.



Equally spaced pass bands *A* are pushed together or warped on the higher frequency end of the digital frequency scale. This effect is normally compensated for by pre-warping the analog filter before applying bilinear transformation.

Because of warping the relationship between Ω_1 and Ω_2 on the one hand and ω_1 and ω_2 on the other is not linear. The digital frequencies ω_1 and ω_2 are pushed in towards the origin ($\omega = 0$). In this process $\Omega = \infty$ is transformed to $\omega = \pi$.



If the bilinear transformation is applied to the system $H_a(s)$ with critical frequency Ω_c , the digital filter will have a critical frequency $\omega_c = 2 \tan^{-1}(\Omega_c T/2)$. If the resulting H(z) is used in an A/D–H(z)–D/A structure, the equivalent critical frequency (of the equivalent analog filter) is obtained by replacing ω_c with $\Omega_{ceq}T$:

$$\omega_c \Rightarrow \Omega_{ceq} T = 2 \tan \left| \left(\frac{\Omega_c T}{2} \right) \right| \qquad \text{or} \qquad \Omega_{ceq} = \frac{2}{T} \tan^{-1} \left(\frac{\Omega_c T}{2} \right)$$

If $(\Omega_c T/2)$ is so small that $\tan^{-1}(\Omega T/2) \approx \Omega_c T/2$, then we have

$$\Omega_{ceq} = \frac{2 \left(\frac{\Omega T}{\Box} \right)}{T \left(\frac{\Box}{2} \right)} = \Omega$$

If this condition is not satisfied, then the warping of the critical frequency (in the bilinear design) is compensated for by *pre-warping*.

Digital filter design – The Butterworth filter

Before taking up design we reproduce below some of the material relating filter specifications to the filter order and the cut-off frequency.

A typical *magnitude* response specification is sketched below. The magnitudes at the critical frequencies Ω_1 and Ω_2 are A and B, respectively. For illustrative purposes we have



arbitrarily taken A = 0.707 (thus Ω_I is the cut-off frequency, but this need not be the case) and B = 0.25.

The *log-magnitude* specification is diagrammed below. Note that $(20 \log A) = K_1$ and $(20 \log B) = K_2$. Thus the analog filter specifications are

$$0 \ge 20\log_{10}|H(j\Omega)| \ge K_1 \text{ for all } \Omega \le \Omega_1$$

$$20\log_{10}|H(j\Omega)| \le K_2 \text{ for all } \Omega \ge \Omega_2$$

 $dB (= 20 \log_{10} |H(j\Omega)|)$



With the magnitude $|H(j\Omega)|$ given by the Butterworth function,

$$|H(j\Omega)| = \frac{1}{\sqrt{1 + \left(\Omega / \Omega_c \gamma^2 N\right)}}$$

and using the equality condition at the critical frequencies in the above specifications the order, N, of the filter is given by

$$N = \begin{bmatrix} \left[10^{-K_{1}/10} - 1 \right] \\ \left[10g_{10} \right] \\ 10g_{10} \\ 10g$$

The result is rounded to the next larger integer. For example, if N = 3.2 by the above calculation then it is rounded up to 4, and the order of the required filter is N = 4. In such a case the resulting filter would exceed the specification at both Ω_1 and Ω_2 . The cut-off frequency Ω_c is determined from one of the two equations below.

$$\Omega_c = \frac{\Omega_1}{\sqrt[2^N]{10^{-K_1/10} - 1}} \quad \text{or} \quad \Omega_c = -\frac{\Omega_2}{\sqrt[N]{-1}} \rightarrow (3.17 \text{ Ludeman})$$

The equation on the left will result in the specification being met exactly at Ω_1 while the specification is exceeded at Ω_2 . The equation on the right results in the specification being met exactly at Ω_2 and exceeded at Ω_1 .

Example 3.6.1 Design and realize a digital low pass filter using the *bilinear transformation* method to satisfy the following characteristics:

- (a) A monotonic pass band and stop band
- (b) -3.01 dB cut off frequency of 0.5π rad.
- (c) Magnitude down at least 15 dB at 0.75π rad.



Note that the given frequencies are digital frequencies. The required frequency response is shown. We use bilinear transformation on an analog prototype.

Step 1 Pre-warp the critical digital frequencies $\omega_l = 0.5\pi$ and $\omega_2 = 0.75\pi$ using T = 1 sec. That is, we find the analog frequencies Ω'_1 and Ω'_2 that correspond to ω_l and ω_2 :

$$\Omega'_{1} = \frac{2}{T} \tan\left(\frac{\omega_{1}}{2}\right) = 2 \tan\frac{0.5\pi}{2} = 2.0 \text{ rad / sec}$$

$$\Omega'_{2} = \frac{2}{T} \tan\left(\frac{\omega_{2}}{2}\right) = 2 \tan\frac{0.75\pi}{2} = 4.8284 \text{ rad / sec}$$

Step 2 Design LP analog filter with critical frequencies Ω'_1 and Ω'_2 that satisfy

$$0 \ge 20 \log |H_a(j \Omega'_1)| \ge -3.01 \text{ dB} = K_I$$
, and
 $20 \log |H_a(j \Omega'_2)| \le -15 \text{ dB} = K_2$

The Butterworth filter satisfies the monotonic property and has an order N and critical frequency Ω_c determined by Eq. 3.16 and 3.17 of Ludeman

$$N = \frac{\begin{vmatrix} \log_{10} | & 10^{-K_{2}/10} - 1 \\ | & 10^{-K_{2}/10} - 1 \end{vmatrix}}{\begin{vmatrix} \log_{10} | & \frac{\Omega_{1}}{1} \\ | & \Omega_{2} \end{vmatrix}} \text{ and } \Omega_{c} = \frac{\Omega_{1}}{\sqrt{10^{-K_{1}/10} - 1}}$$

Plugging in numerical values,

<u>,</u> ¬

Note in this case that $\Omega_c = \Omega_1$.

Therefore, the required *pre-warped*, *normalized*, *unit bandwidth*, analog filter of order 2 using the Butterworth Table 3.1b (or the Butterworth circle) is

$$H_a(s) = \frac{1}{s^2 + \sqrt{2} s + 1}$$
 (with a cut off frequency = 1 rad / sec)

Since we need a cut-off frequency of $\Omega_c = 2$ rad/sec, we next use the low pass to low pass transformation $s \rightarrow s/2$ in order to move the cut-off frequency from 1 to 2 rad/sec.

$$\frac{H_{a}(s)}{s^{2} + \sqrt{2} s + 1} = \frac{1}{(s/2)^{2} + \sqrt{2}(s/2) + 1}$$
$$= \frac{4}{s^{2} + 2\sqrt{2} s + 4} \quad \text{(with a cut-off frequency} = 2 \text{ rad/sec.)}$$

Step 3 Applying the bilinear transformation to $H_a(s)$ with T = 1 will transform the pre-warped analog filter into a digital filter with system function H(z) that will satisfy the given digital requirements:

$$H(z) = H_{a}(s) \Big]_{s \to \overline{1+z_{-1}}}^{2(1-z^{-1})} = \frac{4}{\left[2(1-z^{-1})^{2}\right]_{+}} \left[2(1-z^{-1})\right]_{+}^{2} \sqrt{2} \frac{1+2z^{-1}+z^{-2}}{\left[1+z^{-1}\right]_{+}}^{2} \sqrt{2} \frac{1+2z^{-1}+z^{-2}}{\left[1+z^{-1}\right]_{+}}^{2} + 4$$

$$= -3.414 + 0.585 z$$

HW: Obtain the difference equation. Plot $|H(e^{j\omega})|$ and $\angle H(e^{j\omega})$ vs ω .

Bilinear transformation: Cancellation of sampling time in warping and pre-warping The digital specifications are the set of critical frequencies $\{\omega_1, \omega_2, ..., \omega_N\}$ and the corresponding set of magnitude requirements $\{K_1, K_2, ..., K_N\}$. When an analog filter is used as the prototype for

the bilinear transformation method the relationship between digital and analog frequencies is nonlinear and governed by (a) $(\mathbf{O} \mathbf{T})$

$$\Omega = \frac{2}{tan} \frac{|\Omega|}{|\Omega|} \qquad \text{and} \qquad \omega = 2 \tan \left| \frac{|\Omega|}{|\Omega|} \right|$$

Therefore, to get the proper digital frequency, we must design an analog filter with analog critical frequencies Q_i : i = 1, 2, ..., N given by

$$\Omega_i = \frac{2}{T} \tan \left| \begin{array}{c} \omega \\ \frac{i}{i} \\ 1 \end{array} \right|, \qquad i = 1, 2, \dots, N$$

This operation will be referred to as pre-warping. The corresponding analog magnitude requirements are not changed and remain the same as the corresponding digital requirements. An analog filter $H_a(s)$ is then designed to satisfy the pre-warped specifications given by $\Omega_1, \Omega_2, \dots$ Ω_N and $K_1, K_2, ..., K_N$. The bilinear transformation is then applied to $H_a(s)$, i.e., $H(z) = H(s) \Big|_{2(1-z^{-1})}$

$$H_{a}(s) \int_{s \to \frac{2(1-z^{-1})}{T(1+z^{-1})}} \frac{1}{T(1+z^{-1})}$$

As the T in the Ω_i equation and the T in the bilinear transform cancel in the procedure described above for low pass filter design, it is convenient to just use T = 1 in both places. This is easily seen since if the Ω_i comes from an analog-to-analog transformation of an $H_a(s)$ with a unit radian $2(1-z^{-1})$ is

cut-off frequency, we have $s \rightarrow (s/\Omega_i)$, and when the bilinear transformation $s \rightarrow \frac{1}{T(1+z^{-1})}$

used the cascade of transformations is given by

$$s \to \frac{2(1-z^{-1})}{T(1+z^{-1})\Omega_{i}} = \frac{2(1 \neq 2^{-1}) \omega}{T(1+z^{-1}) \tan^{i}} = \frac{(1-z^{-1}) \omega}{(1-z^{-1}) \tan^{i}} + \frac{1}{2}$$

This does not contain a T. Thus it is immaterial what value of T is used as long as it is the same in both steps (which it is).

The procedure for the design of a digital filter using the bilinear transformation consists of:

Step 1: Pre-warping the digital specifications

Step 2: Designing an analog filter to meet the pre-warped specs

Step 3: Applying the bilinear transformation

In the process T is arbitrarily set to 1, but it can be set equal to any value (e.g., T = 2), since it cancels in the design. The design process is shown by the figure below.



Example 3.6.2 Design a digital low pass filter with pass band magnitude characteristic that is constant to within 0.75 dB for frequencies below $\omega = 0.2613\pi$ and stop band attenuation of at least 20 dB for frequencies between $\omega = 0.41\pi$ and π .



Use *bilinear transformation*. Determine the transfer function H(z) for the lowest order Butterworth design which meets these specifications. Draw the cascade form realization.

Step 1: Pre-warp $\omega_1 = 0.2613 \pi$, $\omega_2 = 0.41 \pi$ with T = 1 sec.

$$\Omega_{l} = \frac{2}{I} \tan \frac{\Omega_{1}}{2} = 2 \tan \frac{0.2613\pi}{0.47\pi} = 0.8703 \text{ rad / sec}$$

$$\Omega_{2} = \frac{1}{T} \tan^{2} = 2 \tan \frac{1}{2} = 1.501 \text{ rad / sec}$$
Step 2: Design $H_{a}(s)$

$$\int_{a}^{1} \frac{10^{+0.75/10} - 1}{10^{+20/10} - 1} \Big|_{1} \int_{a}^{1} \log |^{(1.19-1)}|_{1} \Big|_{1} = \int_{a}^{10} \frac{109}{2(-0.2367)} \Big|_{1} = \int_{a}^{1} \frac{2(-0.2367)}{2(-0.2367)} \Big|_{1} = \int_{a}^{1} \frac{1}{2(-0.2367)} \Big|_{1} = \int_{a$$

$$1+z^{-1}$$

$$= \frac{1}{(s-s_{1})(s-s_{1}^{*})}\Big|_{s \to \frac{2(1-z^{-1})}{1+z^{1}}} \frac{1}{(s-s_{2})(s-s_{2}^{*})}\Big|_{s \to \frac{2(1-z^{-1})}{1+z^{1}}} \frac{1}{(s-s_{2})(s-s_{3}^{*})}\Big|_{s \to \frac{2(1-z^{-1})}{1+z^{1}}}$$

Example 3.6.3 Determine H(z) for a Butterworth filter satisfying the following constraints. Use the *impulse invariance* technique.



Solution The critical frequencies are $\omega_1 = \pi/2$, $\omega_2 = 3 \pi/4$. Use $\omega = \Omega T$ to determine the analog frequencies Ω_1 and Ω_2 . Note *T* is not given. Take T = 1, so that $\omega = \Omega \cdot 1 = \Omega$. (This corresponds to the pre-warping step of the bilinear transformation with $\Omega = (2/T) \tan(\omega/2)$).

We have $\omega = \Omega$.1, so that the critical frequencies are

 $K_{I} = -3.01 \text{ dB}$ $\Omega_{I} = \omega_{I} = \pi/2 \text{ rad/sec.} (= \omega_{c})$ $K_{2} = -13.98 \text{ dB}$ $\Omega_{2} = \omega_{2} = 3\pi/4 \text{ rad/sec.}$

The order of the filter is given by

$$N = \begin{bmatrix} 10^{-K_{1}/10} - 1\\ 10^{2} - 1 \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{10^{-K_{1}/10} - 1}{10^{2} - 1} \right) \\ 10^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} \\ 2 \log_{10} \left[\frac{\pi/2}{10^{2} - 1} \right] \\ 2 \log_{10} \left[\frac{\pi/2}{10^{2} - 1} \right] \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 100^{10} \left(\frac{\pi/2}{10^{2} - 1} \right) \\ \frac{\pi/2}{10^{2} - 1} \end{bmatrix} = \begin{bmatrix} 1$$

 Ω_c is already known to be $\pi/2$ rad/sec.

The 4th order normalized Butterworth filter (with unit bandwidth) is

$$H_a(s) \text{ (normalized)} = \frac{1}{s^4 + 2.613 s^3 + 3.414 s^2 + 2.613 s + 1}$$

Using the low pass to low pass analog transformation $s \rightarrow (s/\Omega_c)$ or $s \rightarrow (s/1.57)$ we get the Butterworth filter satisfying the required specs:

$$H_{a}(s) = \frac{1}{\left(\begin{array}{c} s \end{array}\right)^{4}} \frac{1}{\left(\begin{array}{c} s \end{array}\right)^{3}} \frac{1}{\left(\begin{array}{c} s \end{array}\right)^{2}} \frac{1}{\left(\begin{array}{$$

For the impulse invariant design we need the poles of $H_a(s)$, so the above form is not much help. We need to use the factored form:

$$H_a(s) \text{ (normalized)} = \frac{1}{(s^2 + 0.76536 \ s + 1)(s^2 + 1.84776 \ s + 1)}$$

And with $s \rightarrow (s/1.57)$, we have

$$H_{a}(s) = \frac{1}{\left|\left(\frac{s}{s}\right)^{2} + 0.76536\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right|} \\ \left|\left(1.57\right)^{2} + 0.76536\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(1.57\right)^{2} + 0.76536\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 0.76536\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 0.76536\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 0.76536\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right|\right| \\ \left|\left(\frac{s}{s}\right)^{2} + 1.84776\left(\frac{s}{s}\right) + 1\right| \\ \left|\left(\frac{s}{s}\right)^{2} +$$

Put this into partial fraction form. We need the poles individually since we need to use the relation

$$(s = s_I) \rightarrow (z = e^{s_I T})$$
, with $T = 1$ of course

to get H(z) from the $H_a(s)$. Once we get the $H_a(s)$ we can then combine complex conjugate pole pairs to biquadratic form and then implement as a parallel form with two biquadratics in it. Alternatively, memorize relations 3(c) and 3(d). This latter gives the biquadratics directly.

This same problem will next be solved using the *bilinear transformation* to show the difference.

Step 1: Pre-warping according to $\Omega = \frac{2}{T} \tan \frac{\omega}{2}$ with T = 1 gives $\Omega_1 = 2 \tan \frac{\omega_1}{2} = 2 \tan \frac{(\pi/2)}{2} = 2 \operatorname{rad} / \sec$ $\Omega_2 = 2 \tan \frac{2}{2} = 2 \tan \frac{(3\pi/4)}{2} = 4.828 \operatorname{rad} / \sec$

Step 2: Design *H_a(s)*

$$N = \frac{\left| \left[\log_{10} \right] \left(\left(\underline{10}^{3,9,9,9,0} - \underline{1} \right) \right] \right|}{\left| 2\log_{10} \left(\underline{2}^{-} \right) \right|} \right| = ?$$

$$Q_{c} = \frac{\Omega_{1}}{\left| 4.828 \right|}$$

$$Q_{c} = \frac{\Omega_{1}}{\left| 2(1-z^{-1}) \right|} = ?$$
Step 3: $H(z) = H_{a}(s) \int_{s \to \frac{1+z^{-1}}{1+z^{-1}}}^{2\sqrt{10^{-K_{1}/10}}} - 1$

Example 3.6.4 Design a low pass digital filter by applying *impulse invariance* to an appropriate Butterworth continuous-time filter. The digital filter specs are:

 $-1 \text{ dB} \le 20 \text{ Log } |H(\omega)| \le 0, \qquad 0 \le |\omega| \le 0.2\pi$ $20 \text{ Log } |H(\omega)| \le -15 \text{ dB}, \qquad 0.3\pi \le |\omega| \le \pi$

Solution First convert the digital frequencies ω to analog frequencies Ω . The mapping between ω and Ω is linear in the absence of aliasing. We shall use $\omega = \Omega T$ with T = 1. Thus the specs become

$-1 \text{ dB} \le 20 \text{ Log } /H_a(\Omega) / \le 0,$	$0 \leq \Omega \leq 0.2\pi$
$20 \operatorname{Log} / H_a(\Omega) / \leq -15 \mathrm{dB},$	$0.3\pi \leq \Omega \leq \pi$



$$= \left| \frac{\log_{10} \left(\frac{30.2528}{20.0228} \right)}{2 \left(-0.1761 \right)} \right| = \left| \frac{\log_{10} \left(0.0085 \right)}{2 \left(-0.1761 \right)} \right| = \left| \frac{-2.0729}{2 \left(-0.3827 \right)} \right| = \left| 5.885 \right| = 6$$

$$\Omega_c = \frac{\Omega_1}{2\sqrt{10^{-K_1/10} - 1}} = \frac{0.2\pi}{\sqrt[12]{12} \cdot 2.589 - 1} = \frac{0.2\pi}{\sqrt[12]{12} \cdot 2.589} = \frac{0.2\pi}{0.8935} = 0.703 \text{ rad/sec}$$

Let the left-half plane poles be denoted s_1, s_1, s_2, s_2, s_3 , and s_3 . $H_a(s) = \begin{pmatrix} 1 & * & 1 \\ & * & 1 \end{pmatrix} = \begin{pmatrix} s_1, s_2, s_2, s_3, s_3 \\ & * & 1 \end{pmatrix}$

$$(s-s)(s-s)(s-s)(s-s)(s-s)(s-s)(s-s)$$

Since $\Omega_c = 0.703$ the LP to LP transformation s \rightarrow (s/0.703) results in

$$H_{a}(s) = \frac{1}{(s-s_{1})(s-s_{1}^{*})(s-s_{2})(s-s_{2}^{*})(s-s_{3})(s-s_{3}^{*})} \bigg|_{s \to \frac{s}{0.703}}$$

To be completed

Example 3.6.5 [2003] [The Butterworth circle and the bilinear transformation] Refer to Oppenheim & Schafer, Sec. 5.1.3 and Sec. 5.2.1. The bilinear transformation is given by

$$s = \frac{2}{T} \frac{(1+z^{-1})}{(1+z^{-1})}$$
 or $z = \frac{1+(sT/2)}{1-(sT/2)}$

This last equation is used to map the poles on the Butterworth circle in the s-plane into poles on the Butterworth circle in the z-plane. For the normalized Butterworth filter with a cut-off frequency of 1 rad/sec., the Butterworth circle in the s-plane has unit radius. If the cut-off frequency is Ω_c instead of 1, then the circle has a radius of Ω_c . This is the case in the example on pp. 212-214 of Oppenheim & Schafer where the order N of the filter is 3, the radius of the Butterworth circle in the s-plane is Ω_c , and $\Omega_c T = \frac{1}{2}$ which corresponds to a sampling frequency of twice the cut-off frequency (Figure 5.14).

For the two poles at $s = -\Omega_c$ and $s = \Omega_c$ and $\Omega_c T = \frac{1}{2}$ we get

$$s = -\Omega c: \qquad z = \frac{1 - (\Omega_c T/2)}{1 + (\Omega_c T/2)} = \frac{1 - (1/4)}{1 + (1/4)} = 3/5$$

$$s = \Omega c: \qquad z = \frac{1 + (\Omega_c T/2)}{1 - (\Omega_c T/2)} = \frac{1 + (1/4)}{1 - (1/4)} = 5/3$$

Both of these *z*-plane poles are on the real axis as shown in figure below.



The other *s*-plane poles are similarly mapped to *z*-plane poles, though the algebra involved is a little more. Note that the three poles in the left-half of the *s*-plane are mapped into the inside of the unit circle in the *z*-plane.

5.7 Analog design using digital filters

When we are required to simulate an analog filter using the A/D - H(z) - D/A structure, the specifications consist of the analog frequencies { $\Omega_1, \Omega_2, ..., \Omega_N$ }, the corresponding magnitudes { $K_1, K_2, ..., K_N$ } and the sampling time *T*. We convert to digital specs using the relation $\omega_i = \Omega_i T$.

Example 5.7.1 [Bilinear] [4.2, p. 180, Ludeman] Design a digital filter H(z) that when used in an A/D-H(z)-D/A structure gives an equivalent low-pass analog filter with (a) -3.01 dB cut-off frequency of 500Hz, (b) monotonic stop and pass bands, (c) magnitude of frequency response down at least 15 dB at 750 Hz, and (d) sample rate of 2000 samples/sec. **Solution**

Step 0 First we convert the analog Ω 's to digital ω 's using $\omega_i = \Omega_i T$.

$$\Omega_1 = 2\pi F_1 = 2\pi 500 = 1000 \,\pi \text{rad/sec.}, \qquad K_1 = -3.01 \text{ dB}$$

$$\Omega_2 = 2\pi F_2 = 2\pi 750 = 1500 \,\pi \text{rad/sec.}, \qquad K_2 = -15 \text{ dB}$$

Thus

$$\omega_1 = \Omega_1 T = 1000 \,\pi \, \frac{1}{2000} = 0.5 \,\pi \text{rad}, \qquad K_1 = -3.01 \,\text{dB}$$

 $\omega_2 = \Omega_2 T = 1500 \,\pi \, \frac{1}{2000} = 0.75 \,\pi \text{rad}, \qquad K_2 = -15 \,\text{dB}$

Step 1 Pre-warping. Use T = 1. (In Steps 1 and 3 we could have used T = 1/2000 but the two occurrences of T would cancel out).

$$\Omega_{1}^{'} = \frac{2}{T} \tan \left(\frac{\omega_{1}}{2}\right) = 2.0 \text{ rad/sec.} \text{ and } \Omega_{2}^{'} = \frac{2}{T} \tan \left(\frac{\omega_{2}}{2}\right) = 4.828 \text{ rad/sec.}$$
Step 2 Design $H_{a}(s)$, i.e., determine the low pass Butterworth filter (see earlier example).

$$\left[10^{+3.01/10} - 1 \right] \left[\log \left(\frac{2-1}{1}\right) \right] = \left[\frac{\log_{10}\left| \frac{10^{+15/10} - 1}{10^{+15/10} - 1} \right| \right]}{2\log_{10}\left| \frac{2}{2\log_{10}\left| \frac{10}{2} \right|} \right| = \left[\frac{10(31.62 - 1)}{2(-0.3827)} \right] = \left[\frac{1-1.486}{2(-0.3827)} \right] = \left[1.941 \right] = 2$$

$$\Omega_{c} = \frac{1}{\sqrt{10^{+3.01/10} - 1}} = \frac{1}{\sqrt{2} - 1} = 2 \text{ rad / sec}$$

Do analog low pass to low pass transformation $s \rightarrow (s/\Omega_c)$, i.e., $s \rightarrow (s/2)$ in order to move the cutoff frequency from 1 to 2 rad/sec. This gives the $H_a(s)$ with pre-warped specs and $\Omega_c = 2$ rad/sec.

$$H_{a}(s) = \frac{1}{s^{2} + \sqrt{2} s + 1} \int_{s \to s/2}^{s} = \frac{1}{(s/2)^{2} + \sqrt{2}(s/2) + 1}$$
$$= \frac{4}{s^{2} + 2\sqrt{2} s + 4}$$
(with a cut-off frequency = 2 rad/sec.)

Step 3 Applying the bilinear transformation $s \rightarrow \frac{2}{T} \frac{(1 - z_{-1})}{(1 + z_{-1})}$ to $H_a(s)$ with T = 1 will transform

the pre-warped analog filter into a digital filter with system function H(z) that will satisfy the given requirements:

$$H(z) = H_{a}(s) \Big|_{s \to \Box^{=}}^{2(1-z^{-1})} \frac{4}{\left[\frac{2}{2(1-z^{-1})}\right] + 2\sqrt{2}\left[\frac{2(1-z^{-1})}{1+z^{-1}}\right] + 4} = \frac{1+2z^{-1}+z^{-2}}{3.414+0.585z^{-2}}$$

Example 3.7.2 [2002] Derive the Butterworth digital filter having the following specs:

Pass band: 0 to 4411 rad/sec. Maximum ripple in pass band: 1 dB Stop band: beyond 25975 rad/sec. Minimum attenuation in stop band: 60 dB Sampling frequency: 20 kHz



Solution Even though pass band *ripple* may suggest a Chebyshev filter we shall comply with the request for a Butterworth filter. We use *bilinear transformation*.

Analog frequency specs are given. Convert them to digital by using the relation $\omega = \Omega T$ and T = (1/20000) sec.

$$\Omega_1 = 4411 \text{ rad/sec. becomes } \omega_1 = \Omega_1 T = 4411/20000 = 0.221 \text{ rad}$$

 $\Omega_2 = 25975 \text{ rad/sec. becomes } \omega_2 = \Omega_2 T = 25975/20000 = 1.299 \text{ rad}$

Now, starting from these digital specs the design proceeds in 3 steps as usual.

Step 1 Pre-warp the critical digital frequencies ω_I and ω_2 using T = 1 sec., to get $\Omega'_1 = \frac{1}{T} \tan\left(\frac{1}{2}\right) = 2 \tan\left(\frac{1}{2}\right) =$ Step 2 Design an analog low pass filter with critical frequencies Ω'_1 and Ω'_2 to satisfy

 $0 \ge 20 \log |H_a(j \Omega_1)| \ge -1 dB = K_1$, and 20 $\log |H_a(j \Omega_2)| \le -60 dB = K_2$

The Butterworth filter of order *N* and cut-off frequency Ω_c is given by equations (3.16) and (3.17) of Ludeman:

$$N = \begin{bmatrix} 10^{-K_{1}/10} - 1 \\ 0 \\ 10^{-2} & -1 \end{bmatrix} \begin{bmatrix} 10^{2} \\ 0 \\ 10^{-2} & -1 \end{bmatrix} \begin{bmatrix} 10^{2} \\ 0 \\ 10^{-2} & -1 \end{bmatrix} \begin{bmatrix} 10^{2} \\ 0 \\ 10^{-1} \end{bmatrix} \begin{bmatrix} 10^{2} \\ 0 \\ 0 \\ 10^{-1} \end{bmatrix} \begin{bmatrix} 10^{2} \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 10^{2} \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 10^$$

Therefore the required pre-warped Butterworth (analog) filter using Table 3.1b (Ludeman) and the analog low-pass to low pass transformation from Table 3.2, $s \rightarrow (s/\Omega_c)$, that is, $s \rightarrow (s/0.262)$, is

$$H_a(s) = \frac{1}{s^4 + 2.613 \ s^3 + 3.414 \ s^2 + 2.613 \ s + 1} \Big|_{s \to (s/0.262)}^1$$

Example 3.7.3 Design a digital LPF using *bilinear transformation* with the following specifications, and a Butterworth approximation:

2 dB at 5 rad/sec.,
23 dB at 10 rad/sec.,
Sampling frequency = 1000 per sec.



Solution Convert the analog specs to digital using $\omega = \Omega T$, with T = 1/1000 sec. The critical frequencies are

 $\omega_1 = \Omega_1 T = 5/1000 = 0.005 \text{ rad};$ $K_1 = -2 \text{ dB}$ $\omega_2 = \Omega_2 T = 10/1000 = 0.01 \text{ rad};$ $K_2 = -23 \text{ dB}$

Now apply Steps 1, 2 and 3 of the bilinear transformation design process.

Example 5.7.4 [2003] Design a digital filter that will pass a 1 Hz signal with attenuation less than 2 dB and suppress 4 Hz signal down to at least 42 dB from the magnitude of the 1 Hz signal.



Solution All the specs are given in the analog domain. The sampling period *T* is not specified. Since 1 Hz is in the pass band and 4 Hz in the stop band we shall use some multiple of 4 Hz, say, 20 Hz as the sampling frequency. Thus T = 1/20. We shall employ the *impulse invariance* method.

Step 0 Convert the analog specs to digital by using $\omega = \Omega T = 2\pi FT$. Thus

 $\Omega_1 = 2\pi . 1 = 2\pi \text{ rad/sec.}$, and $\Omega_2 = 2\pi . 4 = 8\pi \text{ rad/sec.}$

so that $\omega_1 = 2\pi T = 2\pi (1/20) = 0.1 \pi$ rad., and $\omega_2 = 8\pi T = 8\pi (1/20) = 0.4 \pi$ rad.

Step 1 Convert the digital frequencies ω_1 and ω_2 back to analog frequencies. Since we are using impulse invariance this involves using the same formula $\omega = \Omega T$ and we get the same analog frequencies as before, viz., $\Omega_1 = 2\pi$ rad/sec., and $\Omega_2 = 8\pi$ rad/sec. (Note that the value of *T* is irrelevant up to this point. We could have used a value of T = 1 in Steps 0 and 1, resulting in awkward values for the ω 's, like 2π and 8π when we expect values between 0 and π ; but this is not a problem for the design).

Step 2 Determine the order of the analog Butterworth filter.

$$N = \begin{vmatrix} 10^{-K_{1}/10} - 1 \\ |\log_{10}| \begin{pmatrix} 10^{-K_{1}/10} - 1 \\ |\Omega_{10}|^{2} & -1 \end{pmatrix} \\ |\Omega_{10}| \begin{pmatrix} 10^{+2/10} - 1 \\ |\Omega_{10}|^{2} & -1 \end{pmatrix} \\ |\Omega_{10}| & |\Omega_{10}|^{2} \\ |\Omega_{10}| & |\Omega_{$$

Cut–off frequency Ω_c is determined next:

$$\Omega_{c} = \frac{\Omega_{1}}{\sqrt{10^{-K_{1}/10} - 1}} = \frac{2\pi}{\sqrt[8]{10^{+2/10} - 1}} = \frac{2\pi}{\sqrt[8]{10^{+2/10} - 1}} = \frac{2\pi}{\sqrt[8]{10^{-K_{1}/10} - 1}} = 6.719$$

Step 3 The normalized filter $H_a(s)$ of order 4 is

$$H_a(s) = \frac{1}{(s^2 + 1.848s + 1)(s^2 + 0.765s + 1)}$$

We may break it down into partial fractions now (before making the transformation $s \rightarrow (s/\Omega_c)$).

$$H_a(s) = \frac{As+B}{s^2 + 1.848s + 1} + \frac{Cs+D}{s^2 + 0.765s + 1}$$

Determine A, B, C, and D and then substitute $s \rightarrow (s/\Omega_c)$.

Alternatively, we can get the individual pole locations from the Butterworth circle.

$$H_{eff}(s) = \frac{1}{(s+0.924+j0.383)(s+0.924-j0.383)(s+0.383+j0.924)(s+0.383-j0.924)}$$
$$= \frac{A}{s-s_1} + \frac{A}{s-s_2} + \frac{C}{s-s_2} + \frac{C}{s-s_2}$$

Determine A, A^* , C and C^* . Next find $H_a(s)|_{s \to (s/\Omega_c)}$ which is the analog prototype. From this we

can find the H(z) by mapping the *s*-plane poles to *z*-plane poles by the relation: $(s = s_I) \rightarrow (z = e^{s_I T})$. Here at last we must specify *T*; we could use T = 1/20. In general, the smaller the value of *T* the better.

$$H_{a}(z) = A \frac{z_{sT}}{z - e^{1}} + A^{*} \frac{z}{z - e^{sT}} + C \frac{z}{z - e^{sT}} + C^{*} \frac{z}{z - e^{sT}} + C^{*} \frac{z}{z - e^{2}}$$

If we were to use the *bilinear transformation* use some value of *T* like 1/20 (justifying it on the basis of the sampling theorem) in $\omega = \Omega T$ in Step 0.

Example 5.7.5 [Low pass filter] Design a digital low pass filter to approximate the following transfer function:

$$H_a(s) = \frac{1}{s^2 + \sqrt{2} s + 1}$$

Using the **BZT** (*Bilinear z-transform*) method obtain the transfer function, H(z), of the digital filter, assuming a 3 dB cut-off frequency of 150 Hz and a sampling frequency of 1.28 kHz. Solution This problem is a slight variation from the pattern we have followed so far in that it specifies *one* critical frequency (cut-off frequency) and the filter order instead of *two* critical frequencies with the filter order unknown.

Step 0 Convert the analog specs to digital

$$\Omega_c = 2\pi F_c = 2\pi 150 = 300 \pi \text{ rad/sec.}$$

The corresponding digital frequency is

Step

$$\omega_{c} = \Omega_{c} T = 300 \ \pi \frac{1}{1280} = 0.2343 \ \pi \text{ rad.}$$

1 Pre-warp (with $T = 1$)
 $\Omega_{c}' = \frac{1}{T} \tan \left(\frac{1}{2} \right) = 0.7715$

Step 2 $H_a(s)$ is given to be of order 2. We only need to do the analog low pass to low pass transformation.

$$H_{a}(s) \xrightarrow{s}{}^{s \to \frac{s}{0.7715}} = \frac{1}{|(\square| + \sqrt{2|\square| + 1})| + 1} \sqrt{2|\square| + 1} \sqrt{2|\square| + 1} \sqrt{0.7715}$$

Step 3 Applying the BZT (with T = 1)

$$H(z) = H_{a}(s) \Big|_{s \to \frac{2(1-z^{-1})}{T(1+z^{-1})}} = \frac{1}{\left[\frac{2(1-z^{-1})}{\lfloor (1+z^{-1}) \ 0.7715}\right]^{2} + \sqrt{2}\left[\frac{2(1-z^{-1})}{\lfloor (1+z^{-1}) \ 0.7715}\right] + 1}$$
$$= \frac{0.0878 (1+2z^{-1}+z^{-2})}{1-1.0048 z^{-1} + 0.3561 z^{-2}}$$

Example 3.7.6 [Low pass filter] Design a low pass digital filter derived from a second order Butterworth analog filter with a 3 dB cut-off frequency of 50 Hz. The sampling rate of the system is 500 Hz.

Solution (This is similar to Example 10). The second order Butterworth analog filter is 1

$$H_{a}(s) = \frac{1}{s^{2} + \sqrt{2} s + 1}$$

Step 0 Convert the analog specs to digital

$$\Omega_c = 2\pi F_c = 2\pi 50 = 100 \pi \text{ rad/sec.}$$

The corresponding digital frequency is

$$\omega_c = \Omega_c T = 100 \ \pi \ \overline{500} = 0.2 \ \pi \ \text{rad.}$$

Step 1 Pre-warp (with $T = 1$)
$$\Omega'_c = \frac{1}{T} \ \tan\left(\frac{\omega_c}{2}\right) = 2 \ \tan\left(\frac{0.2 \ \pi}{2}\right) = 0.6498$$

Step 2 $H_a(s)$ is given to be of order 2. We only need to do the analog low pass to low pass transformation.

$$= \frac{1}{\left(\begin{array}{c} s \\ 0.6498 \end{array}\right)^2} \left(\begin{array}{c} 0.42229 \\ \overline{s^2 + 0.91901s + 0.42229} \\ 0.6498 \end{array}\right)} = \frac{0.42229}{\overline{s^2 + 0.91901s + 0.42229}}$$

Step 3 Applying the BZT (with T = 1)

$$H(z) = H_{a}(s) \Big|_{s \to \frac{2(1-z^{-1})}{T(1+z)}} = \frac{0.42229}{\begin{bmatrix} 2(1-z^{-1}) \end{bmatrix}^{2}} \begin{bmatrix} 2(1-z^{-1}) \end{bmatrix}$$
$$\frac{1}{\begin{bmatrix} 0 \\ (1+z^{-1}) \end{bmatrix}} + 0.91901 \Big|_{(1+z^{-1})} + 0.42229 \Big|_{(1+z^{-1})} \Big| + 0.42229 \Big|_{(1+z^{-1})} \Big|_{(1+z^{-1})}$$

 $= \frac{0.42229(1+z^{-1})^2}{4(1-z^{-1})^2 + (0.91901)(2)(1-z^{-1})(1+z^{-1}) + 0.42229(1+z^{-1})^2}$ is

The denominator is

 $Dr = 4(1 - 2z^{-1} + z^{-2}) + 1.83802(1 - z^{-2}) + 0.42229(1 + 2z^{-1} + z^{-2})$ $Dr = (4 + 1.83802 + 0.42229) + z^{-1}(-8 + 0.84458) + z^{-2}(4 - 1.83802 + 0.42229)$ $Dr = 6.26031 - 7.15542 z^{-1} + 2.58427 z^{-2}$

$$H(z) = \frac{0.42229(1+2z^{-1}+z^{-2})}{6.26031-7.15542z^{-1}+2.58427z^{-2}}$$
$$= _{0.06746(1+2z^{-1}+z^{-2})}$$
$$= _{-6.26031-1.14298z^{-1}+0.4128z}$$

Summary of IIR filter design

u	
Bilinear	Bilinear
Digital specs given:	Analog specs given:
ω_1, ω_2 (rad.)	Ω_1, Ω_2 (rad./sec.),
$K_1, K_2(dB)$	K_1, K_2 (dB), and
$2 (\omega)$	Т
(1) Pre-warp $\Omega' = \frac{1}{T} \tan \left(\frac{1}{2} \right), T = 1$ sec.	Step 0. (Preparation) Use $\omega = \Omega T$ to convert Ω to ω
Calculate Ω_1 , Ω_2 and K_1 , K_2	Get ω_1, ω_2 , and K_1, K_2
(2) Find filter order N, cut-off frequency Ω_c	-, -, 1, 2
Find normalized filter $H_a(s)$	Then
Find $H_a(s)$	(1)
(3) $H(z) = H_a(s) \xrightarrow{ s \to (s/\Omega_c) }_{ s \to (s/\Omega_c) } \text{ with } T = 1 \text{ sec.}$	(2) As in the left side column(3)
Impulse invariance	
Digital specs given:	
$\omega_1, \omega_2,$	
K_1, K_2 , and	
T (or $T = 1$)	
(1) Determine Ω_1 , Ω_2 using $\omega = \Omega T$	
(2) Find filter order N, cut-off frequency Ω_c	
Find normalized filter $H_a(s)$	
Find $H_a(s) _{s \to (s/\Omega_c)}$	
(3) Find individual pole locations.	
Then $(s = s_I) \rightarrow (z = e^{s_I T})$ for given T	

Frequency Transformation

Frequency transformation is useful for converting a frequency-selective filter from one type to another.

Analog-to-analog transformations Suppose we are given a continuous-time normalized low pass filter G(s) with a cut-off frequency of 1 rad/sec. Then what is the effect of the

transformation $s' = s\Omega_c$ where s' represents the transformed frequency variable? In other words, we make the substitution $s \to s'/\Omega_c$ in the transfer function. Since $s' = s\Omega_c$ implies $\Omega' = \Omega \Omega_c$, it follows that the frequency range $0 \le \Omega \le 1$ is mapped into the range $0 \le \Omega' \le \Omega_c$. Thus G(s') represents a low pass filter with a cut-off frequency of Ω_c . In the rest of what follows, rather than use an explicitly different symbol s' we shall instead indicate such transformation by $s \to s/\Omega_c$

and the resulting transfer function by $H(s) = G(s) \Big]_{s \to s/\Omega_c}$.

1. Low-pass to low-pass transformation Given the prototype low pass filter G(s) with unit band width (i.e., cut-off frequency = 1 rad/sec.) and unity gain at $\Omega = 0$ (i.e., |G(j0)| = 1), the transformation $s \rightarrow s/\Omega_c$ gives us a new filter, H(s), with cut-off frequency of Ω_c . The filter H(s)is given by

$$H(s) = G(s) \int_{s \to s/\Omega_c} ds$$

The critical frequency Ω_r of the filter G(s) is transformed to Ω'_r of the filter H(s), given by $\Omega'_r = \Omega_r \Omega_c$. Both G(s) and H(s) are low pass filters.



More generally, the transformation $s \rightarrow s \frac{\Omega_c}{\Omega_c}$

transforms a low pass filter with a cut-off (or critical) frequency Ω_c to a low pass filter with a cut-off (or critical) frequency Ω'_c .

2. Low-pass to high-pass transformation Given G(s) as above the transformation $s \to \Omega_c/s$ will transform the low-pass G(s) into a high-pass filter H(s) with cut-off frequency of Ω_c : $H(s) = G(s) \Big]_{s \to \Omega_c/s}$ The critical frequency Ω_r of the filter G(s) is changed to Ω'_r of H(s), given by $\Omega'_r = \Omega_c/\Omega_r$.



There are similar transformations from low-pass to band-pass and low-pass to band-stop. Refer to table 3.2, page 128, Ludeman.

The Low-pass prototype analog unit bandwidth filter could be any analog filter such as Butterworth, Chebyshev etc. of any order, any ripple etc.

Digital-to-digital transformations Similarly, a set of transformations can be found that take a low-pass digital filter and turn it into another low-pass or high-pass or band-pass or band stop digital filter. Refer to Sec 4.3, p. 181, Ludeman.

The Chebyshev filter

The Butterworth filter provides a good approximation to the ideal low pass characteristics for values of Ω near zero, but has a low fall-off rate in the transition band. The Chebyshev filter has ripples in either the pass band or the stop band but has a sharper cut-off in the transition band. Thus, for filters of the same order, the Chebyshev filter has a smaller transition band than the Butterworth filter. We now look at the analog low pass Chebyshev filter. There are two types of the Chebyshev filter.

Type I (Chebyshev I) is an *all-pole filter*. It has equiripple behavior in the pass band and monotonically decreases in the stop band. For N = order of the filter, the magnitude response looks as below (N = 5 and N = 6 illustrated). The magnitude, $|H(j\Omega)|$, is an even symmetric function of Ω .



Type II (**Chebyshev II** or **Inverse Chebyshev**) filter has *both poles and zeros*. It has a monotonically decreasing shape in the pass band and an equiripple behavior in the stop band.



Design of the Chebyshev I filter A typical *magnitude* response specification is sketched below (shown for N = 5 and N = 6). The magnitudes at the critical frequencies Ω_1 and Ω_2 are A and B, respectively. Typically Ω_1 is in the pass band or is the edge of the pass band and Ω_2 is in the stop band or is the edge of the stop band. In terms of the *log-magnitude* the analog filter specifications are as below. Note that $(20 \log A) = K_1 dB$ and $(20 \log B) = K_2 dB$. If A and B are less than 1, K_1 and K_2 are negative.

$$0 \ge 20\log_{10}|H(j\Omega)| \ge K_1 \text{ for all } \Omega \le \Omega_1$$

$$20\log_{10}|H(j\Omega)| \le K_2 \text{ for all } \Omega \ge \Omega_2$$



The magnitude characteristic of the N^{th} order Chebyshev I filter is given by

. . .

$$|H(j\Omega)| = \frac{1}{\sqrt{1 + \varepsilon^2 C_N(\Omega / \Omega_1)}}, \qquad N = 1, 2,$$

where \langle has to do with pass band ripple and $C_N(x)$ is the N^{th} order Chebyshev cosine polynomial defined as

 $C_N(x) = \cos(N\cos^{-1}x), \qquad x \le 1 \text{ (pass band)}$

 $\cosh(N\cosh^{-1}x), \quad x > 1$ (outside the pass band) Chebyshev polynomials are also defined by the recursion formula

$$C_N(x) = 2xC_{N-1}(x) - C_{N-2}(x)$$

with $C_0(x) = 1$ and $C_1(x) = x$. Using this recursion formula we get, for N = 2, $C_2(x) = 1$ $2xC_1(x) - C_2(x) = 2x^2 - 1.$

Chebyshev Polynomials, $C_N(x)$				
N	$C_N(x)$			
0	1			
1	Х			
2	$2x^2 - 1$			
3	$4x^3 - 3x$			

[Aside If the frequencies are normalized, that is, for a normalized filter, $\Omega_1 = 1$ rad/sec and the magnitude characteristic of the N^{th} order Chebyshev I filter is given by

$$|H(j\Omega)| = \frac{1}{\sqrt{1 + \varepsilon^2 C_N(\Omega)}}, \qquad N = 1, 2, \dots$$

End of Aside]

At $\Omega = 0$ we have $C_N(0) = \cos\left(N\cos^{-1}0\right) = \cos\left(N\pi/2\right)$ = 0, Nodd $\pm 1, N$ even

As a consequence, on the vertical axis ($\Omega = 0$) the magnitude curve starts at |H(j0)| = 1 for odd N and at $|H(j0)| = A = \frac{1}{\sqrt{1+c^2}}$ for even N.

At $\Omega = \Omega_I$ we have $C_N(1) = \cos(N \cos^{-1} 1) = \cos(0) = 1$ for all N. The corresponding magnitude is

$$|H(j\Omega_{I})| = A = \frac{1}{\sqrt{1 + \varepsilon^{2}}}, \quad \text{for all } N$$

This equation is used to compute \langle from the given $|H(j\Omega)|$.

The order, N, of the filter is given by

$$N = \left| \frac{\cosh^{-1} \sqrt{\frac{10^{-K_{2}/10} - 1}{10^{-K_{1}/10} - 1}}}{\cosh \frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \right|$$

The symbol $\lceil \rceil$ means that the computed result is rounded to the next larger integer. For example, if N = 3.2 by the above calculation then it is rounded up to 4, and the order of the required filter is N = 4. In such a case the resulting filter would exceed the specification at both Ω_1 and Ω_2 .

Example 3.9.1 [Filter order] Determine the order of a Chebyshev I filter to have an attenuation of no more than 1 dB for $|\Omega| \le 1000$ rad/sec and at least 10 dB for $|\Omega| \ge 5000$ rad/sec.

Solution The specifications as given are

$$0 \ge 20\log_{10}|H(j\Omega)| \ge -1 \text{ dB for all } \Omega \le 1000 \text{ rad/sec}$$
$$20\log_{10}|H(j\Omega)| \le -10 \text{ dB all } \Omega \ge 5000 \text{ rad/sec}$$

The design parameters are identified as

 $K_1 = -1 \text{ dB}, \quad \Omega_1 = 1000 \text{ rad/sec}$ $K_2 = -10 \text{ dB}, \ \Omega_2 = 5000 \text{ rad/sec}$

The filter order is given by

$$N = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_2/10} - 1}{10^{-K_1/10} - 1}} \\ \frac{\cos h^{-1} \sqrt{\frac{10^{-K_2/10} - 1}{10^{-K_1/10} - 1}}}{\left| \frac{\cos h^{-1} \sqrt{\frac{10^{10/10} - 1}{10^{1/10} - 1}}}{\left| \frac{\cos h^{-1} \sqrt{\frac{10^{-1} - 1}{10^{1/10} - 1}}}}{\left| \frac{\cos h^{-1} \sqrt{\frac{10^{-1} - 1}{10^{1/10} - 1}}}{\left| \frac{\cos h^{-$$

[On the HP 15C, $\cosh^{-1} 5.8956$ is obtained by: (1) Enter Radian mode, (2) Enter 5.8956, (3) g, (4) HYP^{-1} , and (5) COS]

Pole locations and transfer function The poles of the Chebyshev I filter are related to those of the Butterworth filter of the same order and are located on an ellipse in the *s*-plane. If *N* is odd there will be a pole on the negative real axis. In order to find the pole locations and hence the transfer function we introduce the parameter β



Note that if the sinh β and cosh β terms were not present we would have the pole locations of the normalized Butterw₍₀₂rt_kh_f₁ilt_{π}er₎(on the unit circle), that i_s,

$$\zeta_{k} = \underbrace{ \begin{array}{c} \sin\left(\begin{array}{c} 2 & \kappa & 1 \end{array}\right)}_{k} \operatorname{and} & \Omega = \cos\left(\begin{array}{c} 2 & \kappa & 1 \end{array}\right)}_{k} \left(\begin{array}{c} 1 & \chi \\ 1 & N & 2 \end{array}\right)$$

with $\sigma_k^2 + \Omega_k^2 = 1$ which is the unit circle. Thus, the hyperbolic sine and cosine terms are scale factors which, when applied to the Butterworth pole coordinates, give the pole coordinates of a Chebyshev I filter of the same order. The Chebyshev poles are located on an ellipse in the *s*-plane described by

$$\frac{O_{k}^{2}}{\sinh \beta} + \frac{\Omega_{k}^{2}}{\cosh^{2} \beta} = 1$$

The major axis of the ellipse is on the imaginary $(j\Omega)$ axis and the minor axis is on the real axis and the foci are at $\Omega = \pm 1$. The 3 dB cut-off frequency occurs at the point where the ellipse intersects the $j\Omega$ axis, that is, at $\Omega = \cosh \beta$.

In putting together the transfer function, H(s), we rely on the symmetry of pole positions and make use of the left half plane poles only. Finally, the pole positions are scaled by the actual "cut-off frequency" Ω_1 . This last step amounts to $s \rightarrow s/\Omega_1$ (in the case of the Butterworth design this was $s \rightarrow s/\Omega_c$).

Example 3.9.2 [Pole locations and transfer function] Find the pole locations and the transfer function of the Chebyshev I filter designed in above example, that is, with an attenuation of no more than 1 dB for $|\Omega| \le 1000$ rad/sec and at least 10 dB for $|\Omega| \ge 5000$ rad/sec.

Solution The filter order has been determined above as N = 2. Further, we know that $|H(j\Omega_1)| = 1/\sqrt{1 + \varepsilon^2}$ and 20 log $H(j\Omega) \models -1$ dB. Thus $20 \log(1/\sqrt{1 + \varepsilon^2}) = -1$

Solving for $\langle we \text{ get } \langle = 0.5088$. Since *N* is even $|H(j0)| = 1/\sqrt{1 + \varepsilon^2} = 0.8913$ is the starting point on the vertical axis.

With regard to the pole locations, if it were a Butterworth filter of order 2 the poles are located at $(2k-1\pi)$ and $\Omega = \cos \left| \frac{2k-1\pi}{2k-1\pi} \right|$, k = 0, 1

$$\zeta_{k} = \frac{\sin \left| \begin{array}{c} N_{\pi} \right\rangle 2}{\left(-N_{\pi} \right) 2} \right) \quad (-\pi) \quad k \quad \left(\begin{array}{c} \square \\ N & 2 \end{array} \right)$$
$$s_{0,1} = \sin \left| \begin{array}{c} \square \\ 4 \end{array} \right| \pm j \cos \left| \begin{array}{c} \square \\ 4 \end{array} \right| = \left(-1 \quad \sqrt{2} \right) \pm j \left(\frac{1}{\sqrt{2}} \right)$$

The Chebyshev I poles are then obtained from the Butterworth poles by scaling the real and imaginary parts, respectively, by sinh β and cosh β and then scaling both parts by Ω_1 :

$$s_{0,1} = \Omega_1 \left\{ \left(-\frac{1}{\sqrt{2}} \right) \sinh \beta \pm j \left(\frac{1}{\sqrt{2}} \right) \cosh \beta \right\}$$

where $\Omega_1 = 1000 \text{ rad/sec.}$, and

$$\beta = \frac{1}{N} \sinh^{-1}(1/\varepsilon) = \frac{1}{2} \sinh^{-1}(1/0.5088) = 0.7140$$

Thus the Chebyshev I pole locations are

$$s_{0,1} = -\frac{1000}{\sqrt{2}} \sin 0.714 \ 2 \quad \pm j \quad \frac{1000}{\sqrt{2}} \cosh 0.714$$
$$= -\frac{1000}{\sqrt{2}} \sinh 0.714 \ \pm j \quad \frac{1000}{\sqrt{2}} \cosh 0.714$$
$$= -(707.11) \ (0.7762) \ \pm j \ (707.11) \ (1.2659) = -548.86 \ \pm j \ 895.15$$

Hence

$$H(s) = \frac{K}{(s - s_0)(s - s_1)} = \frac{K}{(s - (-548.86 + j895.15))(s - (-548.86 - j895.15))}$$
$$= \frac{K}{(s + 548.86)^2 + (895.15)^2}$$

Since *N* is even the constant *K* will be adjusted to achieve $|H(j0)| = 1/\sqrt{1 + \varepsilon^2} = 0.8913$. (If *N* were odd, *K* would be adjusted to achieve |H(j0)| = 1.)

$$|H(j0)| = \frac{K}{(548.86)^2 + (895.15)^2} = 0.8913$$

which yields $K = 982694.6$. The filter then is
 $H(s) = \frac{982694.6}{2}$

$$s = \frac{1}{(s + 548.86)^2 + (895.15)^2}$$

Example 3.9.3 [Ramesh Babu] Determine the order of a Chebyshev I filter to have a gain of -3 dB or better for $|F| \le 1000$ Hz and a gain of -16 dB or less for $|F| \ge 2000$ Hz. Find the pole locations and the transfer function.

Solution The specifications as given are

$$0 \ge 20\log_{10}|H(j\Omega)| \ge -3 \text{ dB for all } \Omega \le 2\pi \text{ (1000) rad/sec}$$

$$20\log_{10}|H(j\Omega)| \le -16 \text{ dB all } \Omega \ge 2\pi \text{ (2000) rad/sec}$$

The design parameters are identified as

 $K_1 = -3 \text{ dB}, \quad \Omega_1 = 2000 \ \pi \text{ rad/sec}$ $K_2 = -16 \text{ dB}, \quad \Omega_2 = 4000 \ \pi \text{ rad/sec}$

The filter order is given by

$$N = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{2}/10} - 1}{10^{-K_{1}/10} - 1}} \\ \frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1} \\ \frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{16/10} - 1}{10^{3/10} - 1}} \\ \frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{39.81 - 1}{-1}} \\ \frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{16/10} - 1}{10^{3/10} - 1}} \\ \frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{16/10} - 1}{10^{3/10} - 1}} \\ \frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{16/10} - 1}{10^{3/10} - 1}} \\ \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{16/10} - 1}{10^{3/10} - 1}} \\ \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{16/10} - 1}{10^{3/10} - 1}} \\ \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{3/10} - 1}} \\ \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \\ \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}}} \\ \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{10^{-K_{1}/10} - 1}} \end{vmatrix} = \begin{vmatrix} 10^{-K_{1}/10} - 1 \end{vmatrix} = \begin{vmatrix} 10^{-K_{1}/10} - 1 \end{vmatrix} \end{vmatrix} = \begin{vmatrix} 10^{-K_{1}/10} - 1 \end{vmatrix} = \begin{vmatrix} 10^{-K_{1}/$$

Further $|H(j\Omega_1)| = 1/\sqrt{1+\varepsilon^2}$ and $20 \log|H(j\Omega_1)| = -3$ dB implies that $|H(j\Omega_1)| = 1/\sqrt{2}$.

From these two conditions it follows that $\epsilon = 1$. Since N is even $|H(j0)| = 1/\sqrt{1+\varepsilon^2} = 1/\sqrt{2}$. The poles of $H(s^{2k})$ is by $s_{P} = (i + iQ_{k}) k = 0$. 1 + 2k - 1. $\pi + 1 = 0$.

The poles of
$$\mathcal{H}(s)$$
 for griven by $s\beta = \zeta_k \pm \beta \Omega_k$, $k2 = 0, 1 \text{ cosl} | 2k - 1 \pi | \cosh \beta$
 $\zeta_k = | | | | | \cos \beta$
 $\beta = \frac{1}{s} \sinh^{-1} (1/\varepsilon) = -\frac{1}{s} \sinh^{-1} (1/1) = -\frac{0.88137}{2} = 0.44069$
we of the left half plane poles is

Thus, one of the left half plane poles is

$$s_{1} = \sigma_{1} + j\Omega_{1} = \sin\left(\frac{-\pi}{4}\right) \sin 0.44069 + j\cos\left(\frac{-\pi}{4}\right) \cos 0.44069$$
$$= -\frac{1}{\sqrt{2}} 0.45509 + j \quad \frac{1}{\sqrt{2}} 1.09868 = -0.3218 + j0.77689$$

Scaling this by $\Omega_1 = 2000 \pi$ rad/sec results in

 $s_1 = 2000 \pi (-0.3218 + j0.77689) = -2021.9 + j4881.3 = -a + jb$ where a = 2021.9 and b = 4881.3. The other left half plane pole is the conjugate of the above, at

$$s_2 = -2021.9 - j4881.3 = -a - jk$$

The transfer function is H(s) = 1/D(s), where the denominator, D(s), is put together from the poles as follows

$$D(s) = (s - s_1)(s - s_2) = (s - (-2021.9 + j4881.3))(s - (-2021.9 - j4881.3))$$

= (s + 2021.9)² + (4881.3)²

For convenience we shall write this as $D(s) = (s + a)^2 + b^2 = s^2 + 2sa + a^2 + b^2$, so that

$$H(s) = \frac{K}{D(s)} = \frac{K}{s^2 + 2sa + a^2 + b^2}$$

where the constant *K* is adjusted so as to make $|H(j0)| = 1/\sqrt{2}$:

$$|H(j0)| = \frac{K}{a^2 + b^2} = 1/\sqrt{2} \longrightarrow K = \frac{a^2 + b^2}{\sqrt{2}}$$

Thus

$$H(s) = \frac{a^2 + b^2}{\sqrt{2}} \frac{1}{s^2 + 2sa + a^2 + b^2}$$

= 19739005.5 1
 $s^2 + 4043.8s + 27915169.3$

Design of the Chebyshev II filter A typical *magnitude* response specification is sketched below. The magnitudes at the critical frequencies Ω_1 and Ω_2 are *A* and *B*, respectively. Typically Ω_1 is in the pass band or is the edge of the pass band and Ω_2 is in the stop band or is the edge of the stop band. In terms of the *log-magnitude* the analog filter specifications are as below. Note that (20 $\log A$) = K_1 dB and (20 $\log B$) = K_2 dB. If *A* and *B* are less than 1, K_1 and K_2 are negative.

$$0 \ge 20\log_{10}|H(j\Omega)| \ge K_1 \text{ for all } \Omega \le \Omega_1$$

$$20\log_{10}|H(j\Omega)| \le K_2 \text{ for all } \Omega \ge \Omega_2$$



The magnitude characteristic of the N^{th} order Chebyshev II filter is given by

$$|H(j\Omega)| = \frac{1}{\sqrt{1 + \varepsilon^2 \frac{C_N^2(\Omega_2 / \Omega_1)}{C_N^2 \Omega_2 / \Omega_1}}}, \qquad N = 1, 2, ...$$

where ϵ has to do with pass band attenuation and $C_N(x)$ is the N^{th} order Chebyshev polynomial. At $\Omega = \Omega_I$ we have

$$|H(j\Omega_1)| = \frac{1}{\sqrt{1 + \varepsilon^2 \frac{C_N^2(\Omega_2/\Omega_1)}{C_N^2(\Omega_2/\Omega_1)}}} = \frac{1}{\sqrt{1 + \varepsilon^2}} = A, \quad \text{for all } N$$

At $\Omega = 0$, $C_N(\Omega_2/\Omega) = C_N(\infty) \rightarrow \infty$, and the magnitude becomes

$$|H(j0)| = \frac{1}{\sqrt{1 + \varepsilon^2 \frac{C_N^2(\Omega_2 / \Omega_1)}{C_N^2}}} = 1 \quad \text{for all } N$$

The magnitude curve starts at $|H(j0)| = 1$ on the vertical axis for all N .
At $\Omega = \Omega_2$ we have $C_N(\Omega_2 / \Omega_2) = C_N(1) = \cos(N \cos 1) = \cos(0) = 1$
$$|H(j\Omega_2)| = \frac{1}{\sqrt{1 + \varepsilon^2 \frac{C_N^2(\Omega_2 / \Omega_1)}{C_N^2}}} = \frac{1}{\sqrt{1 + \varepsilon^2 \frac{C_N^2(\Omega_2 / \Omega_1)}{1}}}$$
$$= \frac{1}{\sqrt{1 + \varepsilon^2 \frac{C_N^2(\Omega_2 / \Omega_1)}{C_N^2}}} = B, \quad \text{for all } N$$

The order, N_{τ} , of the filter is given by

$$N = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_2/10} - 1}{10^{-K_1/10} - 1}} \\ \frac{10^{-K_1/10} - 1}{\cosh^{-1} \left| \frac{2}{\Omega_1} \right|} \end{vmatrix}$$

The symbol $\lceil \rceil$ means that the computed result is rounded to the next larger integer. For example, if N = 3.2 by the above calculation then it is rounded up to 4, and the order of the required filter is N = 4. In such a case the resulting filter would exceed the specification at both Ω_1 and Ω_2 .

Example 3.9.4 If an analog low pass filter is to have an attenuation of 1 dB at cut-off frequency of 1 kHz, and a maximum stop band ripple of 0.01 for |F| > 5 kHz, determine the required the filter order for (a) a Butterworth filter, (b) a Chebyshev I filter, and (c) a Chebyshev II filter. **Solution** The specifications are the same for all three cases but the magnitude characteristic differs from one case to the next.

(a) The Butterworth magnitude (absolute value) characteristic is sketched below.

$\Omega_I = 2\pi F_I = 2\pi 1000 \text{ rad/sec.},$	$K_I = -1 \text{ dB}$	$\rightarrow A = 0.8912$
$\Omega_2 = 2\pi F_2 = 2\pi 5000 \text{ rad/sec.},$	B = 0.01	$\rightarrow K_2 = -40 \text{ dB}$



The relation between the absolute values and the dB figures ($K_1 = 20 \log A$ and $K_2 = 20 \log B$) is used to compute $A = 10^{K_1/20} = 10^{-1/20} = 0.8912$ and $K_2 = 20 \log B = 20 \log 0.01 = -40$ dB. The Butterworth filter order is given by

$$N = \begin{vmatrix} 10^{-K_{1}/10} - 1 \\ 10^{-K_{2}/10} - 1 \\ 10^{-K_{2}/10} - 1 \end{vmatrix} \begin{vmatrix} \log_{10} \left(\frac{10^{-(-1)/10} - 1}{10^{-(-40)/10} - 1} \right) \\ -\frac{10}{10^{-(-40)/10} - 1} \end{vmatrix} = \begin{vmatrix} \log_{10} \left(\frac{10^{0.1} - 1}{10^{4} - 1} \right) \\ -\frac{10}{10^{4} - 1} \end{vmatrix} \end{vmatrix} = \begin{vmatrix} \log_{10} \left(\frac{10^{0.1} - 1}{10^{4} - 1} \right) \\ -\frac{10}{10^{4} - 1} \end{vmatrix} = \begin{vmatrix} \log_{10} \left(\frac{10^{0.1} - 1}{10^{4} - 1} \right) \\ -\frac{1000}{10^{4} - 1} \end{vmatrix} = \begin{vmatrix} 1000 \\ -\frac{1000}{10^{4} - 1} \end{vmatrix}$$
$$= \begin{vmatrix} 1000 \\ -\frac{10000 - 1}{10^{4} - 1} \\ -\frac{10000 - 1}{10^{4} - 1} \end{vmatrix} = \begin{vmatrix} 1000 \\ -\frac{10000 - 1}{10^{4} - 1} \end{vmatrix}$$

(b) The Chebyshev I specs and magnitude (absolute value) characteristic are diagrammed below. As earlier we compute $A = 10^{K_1/20} = 10^{-1/20} = 0.8912$ and $K_2 = 20 \log B = 20 \log 0.01 = -40 \text{ dB}$.

$\Omega_I = 2\pi F_I = 2\pi 1000 \text{ rad/sec.},$	$K_1 = -1 \text{ dB}$	$\rightarrow A = 0.8912$
$\Omega_2 = 2\pi F_2 = 2\pi 5000 \text{ rad/sec.},$	B = 0.01	$\rightarrow K_2 = -40 \text{ dB}$



(c) The Chebyshev II specs and magnitude (absolute value) characteristic are shown below. As earlier we compute $A = 10^{K_1/20} = 10^{-1/20} = 0.8912$ and $K_2 = 20 \log B = 20 \log 0.01 = -40$ dB.

$$\Omega_1 = 2\pi F_1 = 2\pi 1000 \text{ rad/sec.}, \qquad K_1 = -1 \text{ dB} \quad \to A = 0.8912$$

 $\Omega_2 = 2\pi F_2 = 2\pi 5000 \text{ rad/sec.},$





The Chebyshev II filter order is given by

$$\begin{vmatrix} \cos h^{-1} & 10^{-K_{2}/10} & -1 \\ & \cos h^{-1} & \sqrt{10^{-K_{1}/10} - 1} \\ & \sqrt{10^{-K_{1}/10} - 1} \\ & \sqrt{10^{-(-1)/10} - 1} \\ &$$

Example 3.9.5 Determine the system function H(z) of the lowest order Chebyshev filter that meets the following specs. Use the impulse invariance method.

(a) 0.5 dB ripple in the pass band, $0 \le |\omega| \le 0.24\pi$

(b) At least 50 db attenuation in the stop band, $0.35\pi \le |\omega| \le \pi$

Solution We assume the Chebyshev I filter. The procedure is similar for the Chebyshev II. The specs are:

$$\omega_1 = 0.24 \pi \, \text{rad.}, \qquad K_1 = -0.5 \, \text{dB} \rightarrow \qquad A = 10^{K_1/20} = 0.94406$$

 $\omega_2 = 0.35 \pi \, \text{rad.}, \qquad K_2 = -50 \, \text{dB} \rightarrow \qquad B = 10^{K_2/20} = 0.0031622$

The sampling time, T, is not specified. Since this is impulse invariant design, T should be very small – the smaller the better. Strictly for convenience we shall use T = 1 sec., and convert the digital specs to analog using the relation $\omega = \Omega T$.

$$\Omega_{I} = \omega_{I}/T = 0.24 \pi \,\text{rad/sec}, \quad K_{I} = -0.5 \,\text{dB} \rightarrow A = 10^{K_{1}/20} = 0.94406$$

$$\Omega_{2} = \omega_{2}/T = 0.35 \pi \,\text{rad/sec}, \quad K_{2} = -50 \,\text{dB} \rightarrow B = 10^{K_{2}/20} = 0.0031622$$

The magnitude (absolute value) characteristic is diagrammed below.



The Chebyshev I filter order is given by

$$N = \begin{vmatrix} \cosh^{-1} \sqrt{\frac{10^{-K_{2}/10} - 1}{10^{-K_{1}/10} - 1}} \\ \frac{10^{-K_{1}/10} - 1}{\cosh^{-1} \left| \frac{2}{2} \right|^{1}} \\ \frac{10^{-K_{1}/10} - 1}{\cos^{-1} \left| \frac{2}{2} \right|^{1}} \\ = \left[\cosh^{-1} \sqrt{\frac{10^{-(-50)/10} - 1}{\cos^{-1} \left| \frac{0.35}{2} \right|^{1}}} \right] = \left[\cosh^{-1} \sqrt{\frac{100000 - 1}{10^{-(-0.5)/10} - 1}} \\ \frac{100000 - 1}{\cos^{-1} \left| \frac{100000 - 1}{10^{-(-0.5)/10} - 1} \right|^{1}} \\ = \left[\cosh^{-1} \sqrt{\frac{10^{-K_{1}/10} - 1}{\cos^{-1} \left| \frac{0.35}{\cos^{-1} \left(\frac{0.35}{2} \right) \right|^{1}}} \right] = \left[\cosh^{-1} \sqrt{\frac{10^{-(-50)/10} - 1}{\cos^{-1} \left| \frac{0.35}{2} \right|^{1}}} \\ = \left[\cosh^{-1} \sqrt{\frac{819539.96}{\cos^{-1} \left(\frac{1.4583}{2} \right)}} \right] = \left[\cosh^{-1} \left(905.28 \right) \right] = \left[7.5014 \right] \\ = \left[\cosh^{-1} \left(\frac{1.4583}{\cos^{-1} \left(\frac{1.4583}{2} \right)} \right] \\ = \left[\cosh^{-1} \left(\frac{1.4583}{2} \right) \right] = \left[8.117 \right] = 9$$

Determine the poles and transfer function H(s), h(t), h(n), and H(z).

The Elliptic (or Cauer) filter

An approximation to the ideal low-pass characteristic, which, for a given order of filter, has an even smaller transition band than the Chebyshev filter, can be obtained in terms of Jacobi elliptic sine functions. The resulting filter is called an elliptic filter. The magnitude characteristic of the elliptic filter has ripples in both the pass band and the stop band.


Digital Signal Processing – 4

IV. FIR digital filters

Characteristics of FIR digital filters, Frequency response, Design of FIR digital filters using Window techniques, Frequency sampling technique, Comparison of IIR and FIR filters.

Contents:

FIR – Recapitulation Characteristics if FIR digital filters Frequency response Design of FIR digital filters – The Fourier series and windowing method Choosing between FIR and IIR filters Relationship of the DFT to the *z*-transform

FIR - Recapitulation

Nomenclature With $a_0 = 1$ in the linear constant coefficient difference equation,

$$a_0 y(n) + a_1 y(n-1) + \dots + a_N y(n-N) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M), \qquad a_0 \neq 0$$

we have,

$$H(z) = \frac{\sum_{i=1}^{N} b^{i} z^{-i}}{1 + \sum_{i=1}^{N} a_{i} z^{-i}}$$

М

This represents an IIR filter if at least one of a_1 through a_N is nonzero, and all the roots of the denominator are not canceled exactly by the roots of the numerator. In general, there are M finite zeros and N finite poles. There is no restriction that M should be less than or greater than or equal to N. In most cases, especially digital filters derived from analog designs, $M \le N$. Systems of this type are called N^{th} order systems. This is the case with IIR filter design.

When M > N, the order of the system is no longer unambiguous. In this case, H(z) may be taken to be an N^{th} order system in cascade with an FIR filter of order (M - N).

When N = 0, as in the case of an FIR filter, according to our convention the *order* is 0. However, it is more meaningful in such a case to focus on M and call the filter an FIR filter of M stages or (M+1) coefficients.

Example The system $H(z) = (1 - z^{-8}) (1 - z^{-1})$ is an FIR filter. Why (verify)?

An FIR filter then has only the "b" coefficients and all the "a" coefficients (except a_0 which equals 1) are zero. An example is the three-term moving average filter y(n) = (1/3) x(n) + (1/3) x(n-1) + (1/3) x(n-2). In general the difference equation of an FIR filter can be written

$$y(n) = \sum_{r=0}^{\infty} b_r x(n-r) = b_0 x(n) + b_1 x(n-1) + \dots + b_M x(n-M) \longrightarrow (1)$$

There are (M + 1) coefficients; some use only *M* coefficients. This equation describes a **nonrecursive** implementation. Its impulse response h(n) is made up of the coefficients $\{b_r\} = \{b_0, b_1, ..., b_M\}$

$$h(n) = b_n, \text{ for } 0 \le n \le M \qquad = \{b_0, b_1, \dots, b_M\}$$

0, elsewhere

Equivalently, the finite length impulse response can also be written in the form of a weighted sum of δ functions as was done in Unit I (for example, $x(n) = \sum_{k=1}^{\infty} \frac{x(k)\delta(n-k)}{k}$)

$$h(n) = \sum_{r=0}^{M} b_r \delta(n-r) = b_0 \,\delta(n) + b_1 \delta(n-1) + \dots + b_M \,\delta(n-M)$$

The difference equation (1) is also equivalent to a direct convolution of the input and the impulse response:

$$y(n) = \sum_{r=0}^{M} b_r x(n-r) = \sum_{r=0}^{M} b(r) x(n-r)$$

where we have written b_r as b(r), i.e., the subscript in b_r is written as an index in b(r).

The transfer function H(z) of the FIR filter can be obtained either from the difference equation or from the impulse response h(n):

$$H(z) = \sum_{n=0}^{m=0} h(n)z^{-n} = \bigcup_{0}^{m} + \bigcup_{1}^{m-1} z^{-1} + \bigcup_{2}^{m-2} + \dots + \bigcup_{M}^{m-1} z^{-M}$$

$$= \bigcup_{0}^{m-1} \sum_{1}^{m-1} + \dots + \bigcup_{M}^{m-1} z^{-1} + \bigcup_{M}^{m-1} z^{-$$

The transfer function has *M* nontrivial zeros and an M^{th} order (trivial) pole at z = 0. This is considered as an *all-zero system*.

We may obtain the frequency response $H(e^{j\omega})$ or $H(\omega)$ of the FIR filter either from H(z) as

$$H\left(e^{j\omega}\right) = H\left(z\right)$$

or, from the impulse response, h(n), as the discrete-time Fourier transform (DTFT) of h(n):

$$H(e^{j\omega}) = \sum_{n=-\infty} h(n) e^{-j\omega n} = \sum_{n=0}^{M} b e^{-j\omega n} = b + b e^{-j\omega} + b e^{-j2\omega} + \dots + b e^{-jM\omega}$$

The inverse DTFT of $H(\omega)$ is of course the impulse response, given by

$$h(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} H(\omega) e^{j\omega n} d\omega$$

The basic design problem is to determine the impulse response h(n), or, the coefficients b_r , for r = 0 to M, required to achieve a desired $H(\omega)$. These coefficients are of course the constants that appear in the numerator of the transfer function H(z). The various transformations used in IIR filter design cannot be used here since they usually yield IIR functions, i.e., with both numerator and denominator coefficients.

Characteristics of FIR digital filters

Illustration The equations of the three-term moving average filter are repeated below

$$y(n) = \frac{x(n) + x(n-1) + x(n-2)}{3}$$

$$\frac{Y(z)}{X(z)} = H(z) = \frac{1 + z^{-1} + z^{-2}}{3}$$

$$H(e^{j\omega}) = \frac{1 + e^{-\frac{1+z}{3}} + e^{-\frac{1+z}{3}}}{3} = \frac{e^{-\frac{1+z}{3}} + e^{-\frac{1+z}{3}}}{3} = \left| \left(\frac{1 + 2\cos\omega}{3} \right) \right|_{e^{-j\omega}} e^{-j\omega}$$

This is a crude low pass filter with linear phase, $\angle H(\omega) = -\omega$.

%Magnitude and phase response of 3-coefficient moving average filter %Filter coefficients: h(n) = {1/3, 1/3, 1/3} b3=[1/3, 1/3, 1/3], a=[1] w=-pi: pi/256: pi; Hw3=freqz(b3, a, w); subplot(2, 1, 1), plot(w, abs(Hw3)); legend ('Magnitude'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid





Normalized frequency We define the $r = \omega/\pi$. As ω goes from $-\pi$ to π the variable r goes from -1 to 1. This corresponds to a frequency range of $-F_s/2$ to $F_s/2$ Hz. In terms of the normalized frequency the frequency response of the three-term moving average filter becomes

$$H(r) = \frac{1 + e^{-j\pi r} + e^{-j2\pi r}}{3}$$

% Magnitude and phase response of 3-coefficient moving average filter % Filter coefficients: $h(n) = \{1/3, 1/3, 1/3\}$ subplot(2,1,1); fplot('abs((1/3)*(1+exp(-j*pi*r)+exp(-j*2*pi*r)))', [-1, 1], 'k'); legend ('Magnitude'); xlabel('Normalized frequency, r'); ylabel('Magnitude of H(r)'); grid subplot(2,1,2); fplot('angle((1/3)*(1+exp(-j*pi*r)+exp(-j*2*pi*r)))', [-1, 1], 'k'); legend ('Phase');

xlabel('Normalized frequency');ylabel('Phase of H(r)'); grid



We illustrate below the characteristics of several types of FIR filter. The filter length N may be an odd (preferred) or an even number. Further, we are typically interested in linear phase. This requires the impulse response to have either even or odd symmetry about its center.

Example 4.2.1 Find the frequency response of the following FIR filters

A.	$h(n) = \{0.25, 0.5, 0.25\}$	Even symmetry
B.	$h(n) = \{0.5, 0.3, 0.2\}$	No symmetry
C.	$h(n) = \{0.25, 0.5, -0.25\}$	No symmetry
D.	$h(n) = \{0.25, 0, -0.25\}$	Odd symmetry

Solution

(A) The sequence $h(n) = \{0.25, 0.5, 0.25\}$ has even symmetry.

%Magnitude and phase response of h(n) = {0.25, 0.5, 0.25} %Filter coefficients – Even symmetry b3=[0.25, 0.5, 0.25], a=[1] w=-pi: pi/256: pi; Hw3=freqz(b3, a, w); subplot(2, 1, 1), plot(w, abs(Hw3)); legend ('Magnitude'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw3)); legend ('Phase = -\omega'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



(B) The sequence $h(n) = \{0.5, 0.3, 0.2\}$ is not symmetric.

%Magnitude and phase response of $h(n) = \{0.5, 0.3, 0.2\}$ %Filter coefficients – No symmetry b3=[0.5, 0.3, 0.2], a=[1] w=-pi: pi/256: pi; Hw3=freqz(b3, a, w); subplot(2, 1, 1), plot(w, abs(Hw3)); legend ('Magnitude'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw3)); legend ('Nonlinear Phase'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



(C) The sequence $h(n) = \{0.25, 0.5, -0.25\}$ is not symmetric.

%Magnitude and phase response of $h(n) = \{0.25, 0.5, -0.25\}$ %Filter coefficients – This is *not* odd symmetry b3=[0.25, 0.5, -0.25], a=[1] w=-pi: pi/256: pi; Hw3=freqz(b3, a, w); subplot(2, 1, 1), plot(w, abs(Hw3)); legend ('Magnitude'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw3)); legend ('Nonlinear Phase'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid





%Magnitude and phase response of $h(n) = \{0.25, 0, -0.25\}$ %Filter coefficients – *This* is odd symmetry b3=[0.25, 0, -0.25], a=[1] w=-pi: pi/256: pi; Hw3=freqz(b3, a, w); subplot(2, 1, 1), plot(w, abs(Hw3)); legend ('Magnitude'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw3)); legend ('Phase = -\omega + \pi/2 '); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



Frequency response

Realization of linear phase FIR filters An important special subset of FIR filters has a *linear phase* characteristic. Linear phase results if the *impulse response is symmetric about its center*. For a causal filter whose impulse response begins at 0 and ends at *N*–1, this symmetry is expressed thus

Even: h(n) = h(N-1-n), for n = 0, 1, ..., (N-1) – a total of N points Odd: h(n) = -h(N-1-n), for n = 0, 1, ..., (N-1) – a total of N points

This symmetry allows the transfer function to be rewritten so that only half the number of multiplications is required for the resulting realization.

Linear phase – phase and delay distortion Assume a low pass filter with frequency response $H(e^{j\omega})$ given by

$$H(e^{j\omega}) = 1e^{-j\omega k}, \quad |\omega| < \omega_c$$

$$0, \qquad \omega_c < |\omega| < \pi$$

where *k* is an integer. This is a linear phase filter with the slope of the phase "curve" in the pass band being -k. Let $X(e^{j\omega})$ represent the Fourier transform of an input sequence x(n). Then the transform of the output sequence y(n) is given by $Y(e^{j\omega}) = X(e^{j\omega}) \cdot H(e^{j\omega})$. If $X(e^{j\omega})$ is entirely within the pass band of $H(e^{j\omega})$ then

$$Y(e^{j\omega}) = X(e^{j\omega}) \cdot e^{-j\omega k}$$

So the output signal y(n) can be obtained as the inverse F-transform of $Y(e^{j\omega})$ as





Thus the linear phase filter did not alter the shape of the original signal, simply translated (delayed) it by k samples. If the phase response had not been linear, the output signal would have been a distorted version of x(n).

It can be shown that a causal IIR filter cannot produce a linear phase characteristic and that only special forms of causal FIR filters can give linear phase.

Theorem If h(n) represents the impulse response of a discrete time system, a necessary and sufficient condition for linear phase is that h(n) have a finite duration N, and that it be symmetric about its midpoint.

Example 4.3.1 (a) For the FIR filter of length N = 7 with impulse response h(n) let h(n) = h(N - 1 - n). Show that the filter has a linear phase characteristic. (b) Repeat for N = 8.



Solution (a) For N = 7, the positive symmetry relation h(n) = h(N-1-n) leads to h(n) = h(6-n) which means that h(0) = h(6), h(1) = h(5), and h(2) = h(4), as shown in figure above.

$$H(z) = \sum_{n=0}^{6} h(n)z^{-n} \text{ and } H(e^{j\omega}) = H(z)\Big|_{z=e^{j\omega}} = \sum_{n=0}^{6} h(n)e^{-j\omega n}$$
$$H(e^{j\omega}) = h(0) + h(1) e^{-j\omega} + h(2) e^{-j2\omega} + h(3)e^{-j3\omega} + h(4) e^{-j4\omega} + h(5)e^{-j5\omega} + h(6) e^{-j6\omega}$$
$$= e^{-j3\omega} \{ h(0)e^{j3\omega} + h(1)e^{j2\omega} + h(2)e^{j\omega} + h(3) + h(4)e^{-j\omega} + h(5) e^{-j2\omega} + h(6)e^{-j3\omega} \}$$

Since h(0) = h(6), etc., we can write

$$H(e^{j\omega}) = e^{-j3\omega} \{ h(0)(e^{j3\omega} + e^{-j3\omega}) + h(1)(e^{j2\omega} + e^{-j2\omega}) + h(2)(e^{j\omega} + e^{-j\omega}) + h(3) \}$$

= $e^{-j3\omega} \{ 2h(0)\cos 3\omega + 2h(1)\cos 2\omega + 2h(2)\cos \omega + h(3) \}$
= $a(3)$ = $a(2)$ = $a(1)$ = $a(0)$
= $e^{-j3\omega} \sum_{k=0}^{3} a(k) \cos k\omega$, with $a(0) = h(3)$ and $a(k) = 2h(3-k)$, $k = 1, 2, 3$

The coefficients, in general, are given by

$$a(0) = h \left(\begin{array}{c} N-1 \\ \Box \\ 2 \end{array} \right)^{a(m)} a(k) = 2h \left| \begin{array}{c} N-1 \\ \Box \\ -k \\ 2 \end{array} \right)^{a(m)} for \ k = 1, 2, ..., (N-1)/2$$
$$H(e^{j\omega}) = \pm \left| H(e^{j\omega}) \right|_{e^{j \angle H(e^{j\omega})}} = e^{-j3\omega} \sum_{k=0}^{a(k)} \cos k\omega$$

where $\pm |H(e^{j\omega})| = \sum_{k=0}^{3} a(k) \cos k\omega$ and $\angle H(e^{j\omega}) = \Theta(\omega) = -3\omega$. The phase response is obviously

linear, with slope = -3 = -(N-1)/2 which means that the delay is an integer number of samples.



(b) For N = 8, the positive symmetry relation h(n) = h(N-1-n) leads to h(n) = h(7-n), which means h(0) = h(7), h(1) = h(6), h(2) = h(5), and h(3) = h(4) as shown in figure below.



$$H(z) = \sum_{n=0}^{7} h(n) z^{-n} \quad \text{and} \quad H(e^{j\omega}) = H(z) \Big|_{z=e^{j\omega}} = \sum_{n=0}^{7} h(n) e^{-j\omega n}$$

$$H(e^{j\omega}) = h(0) + h(1) e^{-j\omega} + h(2) e^{-j2\omega} + h(3) e^{-j3\omega} + h(4) e^{-j4\omega} + h(5) e^{-j5\omega} + h(6) e^{-j6\omega} + h(7) e^{-j7\omega} = e^{-j7\omega/2} \{ h(0) e^{j7\omega/2} + h(1) e^{j5\omega/2} + h(2) e^{j3\omega/2} + h(3) e^{j\omega/2} + h(4) e^{-j\omega/2} + h(5) e^{-j3\omega/2} + h(6) e^{-j5\omega/2} + h(7) e^{-j7\omega/2} \}$$

Since h(0) = h(7), etc., we can write

$$H(e^{j\omega}) = e^{-j7\omega/2} \{ h(0)(e^{j7\omega/2} + e^{-j7\omega/2}) + h(1) (e^{j5\omega/2} + e^{-j5\omega/2}) \\ \begin{bmatrix} & +h(2)(e^{j3\omega/2} + e^{-j3\omega/2}) + h(3)(e^{j\omega/2} + e^{-j\omega/2}) \\ & +h(3)(e^{j\omega/2} + e^{-j\omega/2}) \\ & +$$

With b(k) = 2h((N/2) - k), for k = 1, 2, ..., N/2, we can write

$$H(e^{j\omega}) = \pm \left| H(e^{j\omega}) \right| e^{j \angle H(e^{j\omega})} = e^{-j7\omega/2} \sum_{k=1}^{\infty} b(k) \cos[(k-1/2)\omega]$$

where $\pm |H(e^{j\omega})| = \sum_{k=1}^{4} b(k) \cos[(k-1/2)\omega]$ and $\angle H(e^{j\omega}) = \Theta(\omega) = -7\omega/2$. The phase, $\angle H(e^{j\omega})$,

is clearly linear. However, the slope of the phase curve is (-7/2), which is not an integer. The non-integer delay will cause the values of the sequence to be changed, which, in some cases, may be undesirable.



Implementation For a causal filter whose impulse response has even symmetry:

$$h(n) = h(N-1-n),$$
 for $n = 0, 1, ..., (N-1) - a$ total of N points

the transfer function $H(z) = \Im\{h(n)\} = \sum_{n=0}^{\infty} h(n) z^{-n}$ can be written, depending on whether N is even or odd, as follows.

For even *N* The difference equation is derived starting from H(z),

$$H(z) = \sum_{n} h(n) \left(z^{-n} + z^{-(N-1-n)} \right)$$

Since
$$Y(z) = H(z) X(z)$$
, we can write

$$Y(z) = \left(\sum_{n=0}^{\binom{(N/2)-1}{n}} h(n) \left(z^{-n} + z^{-(N-1-n)} \right) \right) X(z)$$

$$= h(0) \left(1 + z^{-(N-1)} \right) X(z) + h(1) \left(z^{-1} + z^{-(N-2)} \right) X(z)$$

$$+ \dots + h \left| \frac{N}{n} - 1 \right| z^{\binom{(2-)}{2}} + z^{-2} \left| X(z) \right|$$

$$(2 + M) \left| \frac{N}{2} \right|$$

Taking the inverse *z*-transform of the above we get y(n) as

$$y(n) = h(0) \begin{bmatrix} x(n) + x(n - N - 1) \\ (N) \end{bmatrix} + h(1) \begin{bmatrix} x(n - 1) + x(n - N - 2) \end{bmatrix} \\ + \dots + h \begin{bmatrix} 2 - 1 \|x\| & n - \begin{bmatrix} -2 - 1 \|x\| & n - 2 \end{bmatrix}$$

The delayed versions of x(n) are added in pairs and then multiplied by coefficients h(.). This is shown in figure below for N = 8. Note that there are an odd number (= 7) of delay elements. There are N/2 = 4 multiplications and (N/2) + 1 = 4 + 1 = 5 adders (actually the number of two-operand additions is 4 + 3 = 7).

Figure for N = 8



For odd *N* We need not derive the equations (they would be necessary if we were writing a computer program to automate it). For N = 7, there are N - 1 = 6 delay elements – an even number of delay elements. There are (N + 1)/2 = (7 + 1)/2 = 4 multiplications and 4 adders (the number of two-operand additions is 6).



Properties of FIR digital filters The sinusoidal steady state transfer function of a digital filter is periodic in the sampling frequency. We have

$$H(e^{j\omega}) = H(z)\Big|_{z=e^{j\omega}} = \sum_{n=-\infty} h(n)e^{-j\omega n}$$

in which h(n) represents the terms of the unit pulse response. The above expression can be decomposed into real and imaginary components by writing

$$H(e^{j\omega}) = \sum_{n = -\infty}^{\infty} h(n) \cos \omega n - j \sum_{n = -\infty}^{\infty} h(n) \sin \omega n = H_R(\omega) + j H_I(\omega)$$

where the real and imaginary parts of the transfer function are given by

$$H_R(\omega) = \sum_{n = -\infty} h(n) \cos \omega n$$
 and $H_I(\omega) = -\sum_{n = -\infty} h(n) \sin \omega n$

These expressions for $H_R(\omega)$ and $H_I(\omega)$ show that

- 1. $H_R(\omega)$ is an even function of frequency and $H_I(\omega)$ is an odd function of frequency.
- 2. If h(n) is an even sequence, the imaginary part of the transfer function, $H_I(\omega)$, will be zero. (The even sequence, h(n), multiplied by the odd sequence sin ωn will yield an odd sequence. An odd sequence summed over symmetric limits yields zero.) In this case

$$H(e^{j\omega}) = \sum_{n = -\infty}^{\infty} h(n) \cos \omega n = H_R(\omega)$$

3. Similarly, if h(n) is an odd sequence, the real part of the transfer function, $H_R(\omega)$, will be zero

$$H(e^{j\omega}) = -j\sum_{n=-\infty}^{\infty}h(n)\sin\omega n = jH_{I}(\omega)$$

Thus an even unit pulse response yields a real-valued transfer function and an odd unit pulse response yields on imaginary-valued transfer function. Recall that a real transfer function has a phase shift of 0 or $\pm \pi$ radians, while an imaginary transfer function has a phase shift of $\pm \pi/2$ radians as shown in figures below. So, by making the unit pulse response either even or odd, we can generate a transfer function that is either real or imaginary.



Two types of applications In designing digital filters we are usually interested in one of the following two situations:

- 1. **Filtering** We are interested in the amplitude response of the filter (e.g., low pass, band pass, etc.) without phase distortion. This is realized by using a real valued transfer function, i.e., $H(e^{j\omega}) = H_R(\omega)$, with $H_I(\omega) = 0$.
- 2. Filtering plus quadrature phase shift These applications include integrators, differentiators, and Hilbert transform devices. For all of these the desired transfer function is imaginary, i.e., $H(e^{j\omega}) = j H_I(\omega)$, with $H_R(\omega) = 0$

FIR Filter Design Procedure

- 1. Decide whether $H_R(\omega)$ or $H_I(\omega)$ is to be set equal to zero. Typically,
 - $H_d(\omega) = H_R(\omega) + j 0$ for filtering, and
 - $H_d(\omega) = 0 + i H_l(\omega)$ for integrators, differentiators and Hilbert transformers
- 2. Expand $H_d(\omega)$ into Fourier series $h_d(n)$. This is the desired impulse response.
- 3. Decide on the length N of the impulse response duration. Truncate the sequence $h_d(n)$ to N samples { $h_t(n)$, n = -(N-1)/2 to (N-1)/2}. Even values of N result in delays of half-sample periods; odd values of N avoid this problem.
- 4. Apply window function {w(n), n = (N-1)/2 to (N-1)/2}
 5. Find the transfer function H(z) = z^{-(N-1)/2} H_t(z) and the frequency response $H(\omega)$. If not satisfactory the value of N may have to be increased or a different window function may be tried.

Phase delay and group delay If we consider a signal that consists of several frequency components (such as a speech waveform or a modulated signal) the **phase delay** of the filter is the amount of time delay each frequency component of the signal suffers in going through the filter. Mathematically, the phase delay τ_p is given by secant

$$\tau_p = - \frac{\Theta(\omega)}{\omega}$$

The **group delay** on the other hand is the average time delay the composite signal suffers at each frequency. The group delay τ_g is given by the slope (tangent) at ω

$$\tau_g = -\frac{d\Theta(\omega)}{d\omega}$$

where $\Theta(\omega) = \angle H(e^{j\omega})$ of the filter.

A nonlinear phase characteristic will cause phase distortion, which is undesirable in many applications, for example, music, data transmission, video and biomedicine.

A filter is said to have a linear phase response if its phase response satisfies one of the following relationships:

 $\Theta(\omega) = -k\omega \rightarrow (A)$ or $\Theta(\omega) = \beta - k\omega \rightarrow (B)$

where k and β are constants. If a filter satisfies equation (A) its group delay and phase delay are the same constant k. It can be shown that for condition (A) to be satisfied the impulse response of the filter must have positive symmetry (aka even symmetry or just symmetry). The phase response in this case is simply a function of the filter length N:

$$h(n) = h(N-1-n), \quad n = 0, 1, 2, ..., (N-1)/2 \text{ for } N \text{ odd}$$

 $n = 0, 1, 2, ..., (N/2) - 1 \text{ for } N \text{ even}$
 $k = (N-1)/2$

If equation (B) is satisfied the filter will have a constant group delay only. In this case, the impulse response h(n) has negative symmetry (aka odd symmetry or antisymmetry):

$$h(n) = -h(N-1-n)$$

 $k = (N-1)/2$
 $\beta = \pi/2$

Analog filter background of phase and group delay Phase delay "at a given frequency" is the slope of the secant line from dc to the particular frequency and is a sort of overall average delay parameter. Phase delay is computed over the *frequency range* representing the major portion of the input signal spectrum (0 to F_1 in the figure below).



The group delay at a given frequency represents the slope of the tangent line at the particular frequency and represents a *local* or *narrow range* (neighborhood of F_1 in the figure) delay parameter.

A case of significance involving both phase delay and group delay is that of a *narrow band modulated signal*. When a narrow band modulated signal is passed through a filter, the carrier is delayed by a time equal to the phase delay, while the envelope (or intelligence) is delayed by a time approximately equal to the group delay. Since the intelligence (modulating signal) represents the desired information contained in such signals, strong emphasis on good group delay characteristics is often made in filters designed for processing modulated waveforms.

Summary of symmetry

Positive symmetry (or just "symmetry" or even symmetry about the middle) is characterized by h(n) = h(N-1-n). Show that for positive symmetry





Summary of symmetry, cont'd

Negative symmetry (or "antisymmetry" or odd symmetry about the middle) is characterized by h(n) = -h(N-1-n). Show that for negative symmetry

a) For N odd (Type III):

$$H(e^{j\omega}) = e^{-\int_{-2}^{-1} \left[\frac{\omega}{2}\right]^{N-1}/2} \sum_{k=0}^{2} a(k) \sin k\omega$$

$$a(0) = h \prod_{k=0}^{-1} k a(k) = 2h \left(\frac{N-1}{2}, k \neq 0\right)$$
b) For N even (Type IV):

$$H(e^{j\omega}) = e^{-\int_{-2}^{-1} \left[\frac{\omega}{2}\right]^{N/2}} \sum_{k=1}^{2} b(k) \sin[(k-1/2)\omega]$$

$$b(k) = 2h \left(\frac{N}{2} - k\right)$$



Qualitative nature of symmetry

Type I Positive symmetry, N is odd. To illustrate take N = 5:

$$H(e^{j\omega}) = e^{-j\omega(N-1)/2} \sum_{k=0}^{\infty} a(k) \cos k\omega = e^{-j2\omega} \sum_{k=0}^{\infty} a(k) \cos k\omega$$
$$= e^{-j2\omega} [a(0) + a(1) \cos \omega + a(2) \cos 2\omega]$$

We have to add up a(0), and the two cosine terms. It is clear that at $\omega = 0$ all the cosine terms are at their positive peak, so that when added the response of $H(e^{i\omega})$ vs. ω would indicate a low pass filter. Consider

$$y(n) = \frac{x(n) + x(n-1) + x(n-2) + x(n-3) + x(n-4)}{9(n) = 4}$$

$$H(e^{j\omega}) = H(z)\Big|_{z=\frac{j\omega}{e}} = \frac{1 + e^{-j\omega} + e^{-j2\omega} + e^{-j4\omega}}{1 + e^{-j\omega} + e^{-j2\omega} + e^{-j4\omega}}$$

$$= e^{-j2\omega}\Big|\left(\frac{e^{j2\omega} + e^{j\omega} + 1 + e^{-j\omega} + e^{-j2\omega}}{5}\right) = \left(\frac{1 + 2\cos\omega + 2\cos 2\omega}{\Box}\right)$$

%Frequency response of moving average filter $h(n) = \{0.2, 0.2, 0.2, 0.2, 0.2\}$ b5 = [0.2, 0.2, 0.2, 0.2], a = [1]

w=-pi: pi/256: pi; Hw5=freqz(b5, a, w);

subplot(2, 1, 1), plot(w, abs(Hw5)); legend ('Magnitude'); title ('Type I, N is odd'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw5)); legend ('Phase');

xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



Type II Positive symmetry, *N* is even. Take N = 6:

$$H(e^{j\omega}) = e^{-j\omega(N-1)/2} \sum_{k=1}^{N/2} b(k) \cos[(k-1/2)\omega] = e^{-j5\omega/2} \sum_{k=1}^{3} b(k) \cos[(k-1/2)\omega]$$
$$= e^{-j5\omega/2} [b(1) \cos \omega/2 + b(2) \cos 3\omega/2 + b(3) \cos 5\omega/2]$$

At $\omega = \pi$, corresponding to half the sampling frequency (maximum possible frequency), all the cosine terms will be zero. Thus this type of filter is unsuitable as a high-pass filter. It should be ok as a low pass filter. Consider

$$y(n) = \frac{x(n) + x(n-1) + x(n-2) + x(n-3) + x(n-4) + x(n-5)}{6}$$

$$H(e^{j\omega}) = H(z) \Big|_{z=e^{j\omega/2}6} = \frac{1 + e^{-j\omega} + e^{-j2\omega} + e^{-j3\omega} + e^{-j4\omega} + e^{-j5\omega}}{e^{-j(1/2)\omega} + e^{-j(1/2)\omega} + e^{-j(3/2)\omega} + e^{-j(5/2)\omega}} \Big|_{z=e^{j\omega/2}6} = \frac{e^{-j(5/2)\omega} \left(\left| \frac{e^{j(5/2)\omega} + e^{j(3/2)\omega} + e^{j(1/2)\omega} + e^{-j(1/2)\omega} + e^{-j(5/2)\omega}}{6} \right|_{z=e^{j(5/2)\omega}} \right)}{6} = \left(\frac{\cos(\omega/2) + \cos(3\omega/2) + \cos(5\omega/2)}{3} \right)^{2} e^{-j(5/2)\omega}}{3}$$

title ('Type II, N is even');

xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw6)); legend ('Phase');

xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



Type III Negative symmetry, *N* is odd. This introduces a 90⁰ (= $\pi/2$) phase shift. Because of the sine terms |H| is always zero at $\omega = 0$ and at $\omega = \pi/2$ (half the sampling frequency). Therefore the filter is unsuitable as a low pass or a high pass filter. To illustrate take *N* = 5 and

$$h(n) = \{0.2, 0.2, 0, -0.2, -0.2\}$$

$$= e^{-j\left(\omega^{N-1} - \frac{\pi}{2}\right)(N-1)/2} \sum_{k=0}^{-j\left(\omega^{N-1} - \frac{\pi}{2}\right)(5-1)/2} \sum_{k=0}^{-j\left(\omega^{N-1} - \frac{\pi}{2}\right)(5-1)/2} \sum_{k=0}^{2} a(k) \sin k\omega$$

$$= e^{-j2\omega + j(\pi/2)} \sum_{k=0}^{2} a(k) \sin k\omega$$

$$a(0) = h\left(\frac{N-1}{2}\right) = h(2) = 0$$

$$a(k) = 2 \left(\frac{N-1}{2} - k\right) = 2h(2-k), k \neq 0$$
Etc.

Etc.

% Frequency response of Type III filter, h(n) = {0.2, 0.2, 0, -0.2, -0.2} b5 = [0.2, 0.2, 0, -0.2, -0.2], a = [1] w=-pi: pi/256: pi; Hw5=freqz(b5, a, w); subplot(2, 1, 1), plot(w, abs(Hw5)); legend ('Magnitude'); title ('Type III, N is odd'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw5)); legend ('Phase'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



Type IV Negative symmetry, *N* is even. This introduces a 90⁰ (= $\pi/2$) phase shift. Because of the sine terms |H| is always zero at $\omega = 0$. Therefore the filter is unsuitable as a low pass filter. To illustrate take N = 6 and

$$h(n) = \{\frac{1/6, \frac{1}{6}, \frac{1}{6}, \frac{-1}{6}, \frac{-1}{6}, \frac{-1}{6}, \frac{-1}{6}, \frac{-1}{6}\}$$

$$= e^{-j(\omega - 1) | \frac{\pi}{2} | \frac{$$

Etc.

%Frequency response of Type IV filter $h(n) = \{1/6, 1/6, 1/6, -1/6, -1/6, -1/6\}$ b6 = [1/6, 1/6, 1/6, -1/6, -1/6], a = [1]w=-pi: pi/256: pi; Hw6=freqz(b6, a, w); subplot(2, 1, 1), plot(w, abs(Hw6)); legend ('Magnitude'); title ('Type IV, N is even'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw6)); legend ('Phase'); xlabel('Enguency \omega, rad/sample'), ylabel('Dhase of H(\omega)'); grid

xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



Types III and IV are often used to design **differentiators** and **Hilbert transformers** because of the 90^{0} phase shift that each one can provide.

The phase delay for Type I and II filters or group delay for all four types of filters is expressible in terms of the number of coefficients of the filter and so can be corrected to give a zero phase or group delay response.

Types I and II: $\tau_p = \tau_g = -\Theta(\omega)/\omega = \left(\frac{N-1}{2}\right)T$ Types III and IV: $\tau_g = -\frac{d\Theta(\omega)}{2} T \left(\frac{N-1}{2}\right)$

	Magnitude (H(
	$\omega = 0$ rad.	$\omega = \pi$ rad.	
Type I	Max		Low pass
Type II		Zero	OK as LP
			Not OK as HP filter
Type III	Zero	Zero	90 [°] Phase shift
Type IV	Zero		90 [°] Phase shift

Design of FIR digital filters – The Fourier series and windowing method

This method of filter design originates with the observation that the sinusoidal steady state transfer function of a digital filter is periodic in the sampling frequency. Since $H(\omega)$ is a continuous and periodic function of ω we can expand it into a Fourier series. The resulting Fourier coefficients are the impulse response, h(n). The major disadvantage is that one cannot easily specify in advance the exact values for pass band and stop band attenuation/ripple levels, so it may be necessary to check several alternate designs to get the required one.

Consider the ideal low pass filter with frequency response $H_d(e^{j\omega})$ or $H_d(\omega)$ as shown below. The subscript *d* means that it is the *desired* or *ideal* filter.



The impulse response is given by $(e^{j\omega}) e^{j\omega n} d\omega = \begin{pmatrix} 1 & \omega_c \\ & & \end{pmatrix}^{j\omega n} \omega$

$$h_{d}(n) = \underbrace{\int H_{d}}_{2\pi} \underbrace{\int 1e}_{-\pi} d$$

$$= \begin{cases} \frac{\sin n\omega_{c}}{\pi n}, & -\infty \leq n \leq \infty, n \neq 0 \\ \frac{\omega_{c}}{\pi}, & n = 0 \end{cases}$$

This is a non-causal infinite impulse response sequence. It is made a finite impulse response sequence by truncating it symmetrically about n = 0; it is made causal by shifting the truncated sequence to the right so that it starts at n = 0. The shifting results in a time delay and we shall ignore it for now. The truncation results in the sequence $h_t(n)$ where the subscript *t* means truncation but we shall ignore the subscript

$$h(n) = h_d(n), \qquad -(N-1)/2 \le n \le (N-1)/2$$

0, otherwise

In general h(n) can be thought of as obtained by multiplying $h_d(n)$ with a window function w(n) as follows

$$h(n) = h_d(n) \cdot w(n)$$

For the h(n) obtained by simple truncation as above the window function is a rectangular window given by

$$w(n) = 1, -(N-1)/2 \le n \le (N-1)/2$$

0, otherwise

Let $H(e^{j\omega})$, $H(e^{j\omega})$ and $W(e^{j\omega})$ represent the Fourier transforms of h(n), $h_d(n)$ and w(n) respectively. Then the frequency response $H(e^{j\omega})$ of the resulting filter is the convolution of $H_d(e^{j\omega})$ and $W(e^{j\omega})$ given by

$$H(e^{j\omega}) = \underbrace{1}_{2\pi} \int_{d}^{\pi} H(e^{j\theta}) W(e^{j(\omega-\theta)} d\theta = H(e^{j\omega}) * W(e^{j\omega})$$

The convolution produces a smeared version of the ideal low pass filter. In other words, $H(e^{j\omega})$ is a smeared version of $H(e^{j\omega})$. In general, the wider the main lobe of $W(e^{j\omega})$, the more spreading or smearing, whereas the narrower the main lobe (larger *N*), the closer

 $H(e^{j\omega})$ comes to $H_{d}(e^{j\omega})$.

For any arbitrary window the *transition band* of the filter is determined by the width of the *main lobe* of the window. The *side lobes* of the window produce *ripples* in both pass band and stop band.

In general, we are left with a trade-off of making N large enough so that smearing is minimized, yet the number of filter coefficients (= N) is not too large for a reasonable implementation. Some commonly used windows are the rectangular, Bartlett (triangular), Hanning, Hamming, Blackman, and Kaiser windows.



Example 4.4.1 [Design of 9-coefficient LP FIR filter] Design a nine-coefficient (or 9-point or 9-tap) FIR digital filter to approximate an ideal low-pass filter with a cut-off frequency $\omega_c = 0.2\pi$. The magnitude response, $|H_d(\omega)|$, is given below. Take $\angle H_d(\omega) = 0$.



Solution The impulse response of the desired filter is

$$\frac{1}{4} \frac{\pi}{(e^{j\omega}) e^{j\omega n} d\omega} = \frac{1}{2\pi} \frac{0.2\pi}{\int_{j\omega n}^{j\omega n} \omega} \frac{1}{\int_{j\omega n}^{j\omega n} \omega} \frac{1}{\int_{j\omega n}^{j\omega n} \int_{j\omega n}^{j\omega n} \omega} \frac{1}{\int_{j\omega n}^{j\omega n} \int_{j\omega n}^{j\omega n} \omega} \frac{1}{\int_{j\omega n}^{j\omega n} \int_{j\omega n}^{j\omega n} \omega} \frac{1}{\int_{j\omega n}^{j\omega n} \int_{-0.2\pi}^{j\omega n} \omega} \frac{1}{\int_{j\omega n}^{j\omega n} \omega} \frac{1}{\int_{j$$

Since $h_d(n) \neq 0$ for n < 0 this is a noncausal filter (also, it is not BIBO stable – see Unit IV). The rest of the design is aimed at coming up with a noncausal approximation of the above impulse response.

For a rectangular window of length 9, the corresponding impulse response is obtained by evaluating $h_d(n)$ for $-4 \le n \le 4$ on a calculator. In the MATLAB segment below, which generates the $h_d(n)$ coefficients for $-50 \le n \le 50$, division by zero for n = 0 causes "NaN" (*Not a Number*), while all the other coefficients are correct. In generating the frequency

response $|H(\omega)|$ we copy and paste all the coefficients except for $h_d(0)$ which is entered by hand.

$$h_d(0) = \frac{\frac{d (\sin 0.2\pi n)}{dn}}{\frac{d \underline{x} n}{dn}} \bigg|_{n=0} = \frac{0.2\pi \cos 0.2\pi n}{\pi} \bigg|_{n=0} = 0.2$$

Aside (MATLAB) The segment below generates and stem-plots the $h_d(n)$ coefficients.

%Calculate hdn = $(\sin (0.2*pi*n)) / (pi*n)$ and stem plot n = -50: 50, hdn = $(\sin(0.2*pi*n)) ./(pi*n)$, stem(n, hdn) xlabel('n'), ylabel('hd(n)'); grid; title ('hd(n) = $(\sin (0.2*pi*n)) / (pi*n)')$

n = -50 to 50

Warning: Divide by zero.

hdn = -0, -0.0038, -0.0063, -0.0064, -0.0041, 0, 0.0043, 0.0070, 0.0072,0.0046 -0. -0.0048 -0.0080 -0.0082 -0.0052 0,0.0055 0.0092, 0.0095 0.0060 -0, -0.0065 -0.0108 -0.0112 -0.0072 0, 0.0078, 0.0132 0.0138 0.0089 -0, -0.0098 -0.0168 -0.0178 -0.0117 0, 0.0134 0.0233 0.0252 0.0170 -0, -0.0208 -0.0378 -0.0432 -0.0312, 0,0.0468 0.1009 0.1514 0.1871 NaN 0.1871 0.1514 0.1009, 0.0468 0.0000 -0.0312 -0.0432 -0.0378 -0.0208 -0.0000 0.0170 0.0252, 0.0233 0.0134 0.0000 -0.0117 -0.0178 -0.0168 -0.0098 -0.0000 0.0089.0.0138 0.0132 0.0078 0.0000 -0.0072 -0.0112 -0.0108 -0.0065 -0.0000, 0.0060 0.0095 0.0092 0.0055 0.0000 -0.0052 -0.0082 -0.0080 -0.0048, -0.0000 0.0046 0.0072 0.0070 0.0043 0.0000 -0.0041 -0.0064 -0.0063, -0.0038 -0.0000



Note 1 The segment below is used to get a quick look at the $h_d(n)$ coefficients and their symmetry. The MATLAB problem with n = 0 may be avoided by replacing n with (n - 0.001). This will affect the other coefficients very slightly (which is not a serious problem as far as demonstrating the even symmetry of $h_d(n)$) but the accuracy of the coefficients is somewhat compromised in the third or fourth significant digit.

%Calculate hdn = $(\sin (0.2*pi*n)) / (pi*n)$ and stem plot n = -4: 4, hdn = $(\sin(0.2*pi*(n-0.001))) ./(pi*(n-0.001))$, stem(n, hdn), xlabel('n'), ylabel('hd(n)'); grid; title ('hd(n) = $(\sin (0.2*pi*n)) / (pi*n)')$

hdn = 0.0467, 0.1009, 0.1513, 0.1871, 0.2, 0.1871, 0.1514, 0.1010, 0.0468



Note 2 As an alternative to the above, one could write a custom program to calculate *all* coefficients exactly including $h_d(0)$.

End of Aside

The values are shown in table below:

$$n = -4 -3 -2 -1 0 1 2 3 4$$

 $h_t(n) = \{0.047, 0.101, 0.151, 0.187, 0.2, 0.187, 0.151, 0.101, 0.047\}$

By a rectangular window of length 9 we mean that we retain the above 9 values of $h_d(.)$ and *truncate* the rest outside the window. Thus

$$h_t(-4) = 0.047, h_t(-3) = 0.101, \dots, h_t(0) = 0.2, \dots, h_t(4) = 0.047$$



The transfer function of this filter is

$$H_{t}(z) = \sum_{n=-4}^{\infty} h_{d}(n) z_{-n} = 0.047 z^{4} + 0.101 z^{3} + 0.151 z^{2} + 0.187 z^{1} + 0.2 z^{0} + 0.187 z^{-1} + 0.151 z^{-2} + 0.101 z^{-3} + 0.047 z^{-4}$$

The causal filter is then given by *delaying* the sequence $h_t(n)$ by 4 samples. That is, $h(n) = h_t(n-4)$, and the resulting transfer function is $H(x) = x^{-4} H(x)$

$$H(z) = z^{-4} H_t(z)$$

= 0.047 (1+ z^{-8}) + 0.101 (z^{-1} + z^{-7}) + 0.151 (z^{-2} + z^{-6})
+ 0.187 (z^{-3} + z^{-5}) + 0.2 z^{-4}

We may obtain the frequency response of this realizable (causal) filter by setting $z = e^{j\omega}$, that is, $H(e^{j\omega}) = H(z)|_{z^{\frac{j\omega}{2}}}$. Because of the *truncation* the magnitude |H| will only be approximately equal to $|H_d|$ – Gibbs phenomenon, see comparison of 9 coefficients versus 101 coefficients below. Further, because of the *delay* the phase $\angle H = -4\omega$ whereas, as originally specified, $\angle H_d = 0$. The slope $\frac{d \angle H(\omega)}{d\omega} = -4$, showing that the filter introduces a delay of 4 samples.



%Magnitude and phase response of 9-coefficient LP filter %Filter coefficients b9=[0.0468, 0.1009, 0.1514, 0.1871, 0.2, 0.1871, 0.1514, 0.1009, 0.0468], a=[1] w=-pi: pi/256: pi; Hw9=freqz(b9, a, w); subplot(2, 1, 1), plot(w, abs(Hw9)); legend ('9 coefficients'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw9)); legend ('9 coefficients'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid







Rectangular window In this example the sequence $h_d(n)$, which extends to infinity on both sides, has been truncated to 9 terms. This truncation process can be thought of as multiplying the

. . . .

infinitely long sequence by a *window function* called the rectangular window, $w_R(n)$. The figure below shows both $h_d(n)$ and $w_R(n)$ in the undelayed form, that is, symmetrically disposed about n = 0.

Specifying the window function The interval over which the window function is defined



depends on whether we first delay $h_d(n)$ and then truncate it or the other way around. If, instead of truncating first and then delaying, we adopt the procedure of first delaying $h_d(n)$ and then truncating it, the window function may be defined over the interval $0 \le n \le N-1$, where N is the number of terms retained. With this understanding the rectangular window, $w_R(n)$, is given below.

$$w_R(n) = 1, 0 \le n \le N-1$$

0, elsewhere

If, however, we define $w_R(n)$ symmetrically about n = 0, as we do later, we have

$$w_R(n) = 1, -(N-1)/2 \le n \le (N-1)/2$$

0, elsewhere

In this case we have implied that *N* is odd.

Hamming window As an alternative to the rectangular window we shall apply the Hamming window, defined over the interval $0 \le n \le N-1$, by

 $w_{Ham}(n) = 0.54 - 0.46 \cos[2\pi n / (N-1)], \quad 0 \le n \le N-1$ 0, elsewhere

For N = 9, the Hamming window is given by $w_{Ham}(n) = 0.54 - 0.46 \cos(\pi n / 4)$, $0 \le n \le 8$

$$n = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8$$

$$w_{Ham}(n) = \{0.08, \quad 0.215, \quad 0.54, \quad 0.865, \quad 1, \quad 0.865, \quad 0.54, \quad 0.215, \quad 0.08\}$$



Imagine that we line up this sequence alongside the $h_d(n)$ given earlier. (This means we should imagine $w_{Ham}(n)$ is moved to the left by 4 samples). We then multiply the two sequences at each point to get the windowed sequence, $h_t(n)$:

<i>n</i> =	-4	-3	-2	-1	0	1	2	3	4
$h_d(n) =$	{0.047,	0.101,	0.151,	0.187,	0.2,	0.187,	0.151,	0.101,	0.047}
w(n) =	{0.08,	0.215,	0.54,	0.865,	1,	0.865,	0.54,	0.215,	0.08}
$h_t(n) =$	{0.00382,	0.0216,	0.0815,	0.1617,	0.2,	0.1617,	0.0815,	0.0216,	0.00382}

The transfer function is given by

$$H_t(z) = 0.00382 z^{4} + 0.0216 z^{3} + 0.0815 z^{2} + 0.1617 z^{1} + 0.2 + 0.1617 z^{-1} + 0.0815 z^{-2} + 0.0216 z^{-3} + 0.00382 z^{-4}$$

Delaying by 4 sample periods we get

$$H(z) = z^{-4} H_t(z) = 0.00382 (1 + z^{-8}) + 0.0216 (z^{-1} + z^{-7}) + 0.0815 (z^{-2} + z^{-6}) + 0.1617 (z^{-3} + z^{-5}) + 0.2 z^{-4}$$

HW Write the corresponding difference equation and show the filter structure.

We compare below the 9-tap Hamming windowed filter to the filter without the window.

%Comparison of no window vs. Hamming window
%Nine-tap filter coefficients, no window
b9=[0.0468, 0.1009, 0.1514, 0.1871, 0.2, 0.1871, 0.1514, 0.1009, 0.0468],
%
%Nine-tap, Hamming-windowed filter coefficients

b9Ham=[0.00382, 0.0216, 0.0815, 0.1617, 0.2, 0.1617, 0.0815, 0.0216, 0.00382] a=[1] w=-pi: pi/256: pi; Hw9=freqz(b9, a, w); Hw9Ham=freqz(b9Ham, a, w); plot(w, abs(Hw9), w, abs(Hw9Ham), 'k') legend ('9 coefficients, No window', '9 coefficients, Hamming window'); title('Magnitude Response'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid



Example 4.4.2 [Low pass filter] [2003] Design a low pass FIR filter that approximates the following frequency response

 $H(F) = 1, \quad 0 \le F \le 1000 \text{ Hz}$ 0, elsewhere in $0 \le F \le F_s/2$

where the sampling frequency F_s is 8000 sps. The impulse response duration is to be limited to msec. Draw the filter structure.

Solution Note the specs are given in Hertz. On the digital frequency scale ω goes from 0 to 2π , with 2π corresponding to the sampling frequency of $F_s = 8000$ Hz or to $\Omega_s = 2\pi 8000$ rad/sec. Based on this 1000 Hz corresponds to $(1/8)2\pi = \pi/4$. Or we may use the relation $\omega = \Omega T$ to convert the analog frequency 1000 Hz to the digital frequency. Thus

$$\omega = \Omega T = 2\pi F T = 2\pi 1000(1/8000) = \pi/4$$

Thus the specifications are restated on the ω scale as



$$2\pi jn$$
 πn

We need to decide the number of coefficients, that is, the filter length, needed. The problem specifies the impulse response duration is to be limited to 2.5 msec. At 8000



samples/sec., the sampling period T = 1/8000 = 0.125 msec. So the duration of 2.5msec translate to 2.5/0.125 = 20 sample periods. Arranging these on a linear scale we see that the filter length N = 21 or we need 21 coefficients. Therefore determine the values $h_d(-10)$ through $h_d(10)$,

$$h_t(n) = \frac{\sin 0.25 \pi n}{\pi n}, \quad -10 \le n \le 10$$
$$h_t(0) = 0.25 \text{ by L'Hopital's rule}$$

%Calculate hdn = $(\sin (0.25*pi*n)) / (pi*n)$ and stem plot n = -50: 50, hdn = $(\sin(0.25*pi*n)) ./(pi*n)$, stem(n, hdn)

n = -50 to 50 Warning: Divide by zero.
	hdn = [0]	0.0064 ().0046 -0	.0000 -().0048 -().0069 -	0.0050	0.0000	0.0052
	0.0076	0.0055	-0.0000	-0.0058	-0.0084	-0.0061	0.0000	0.006	4
	0.0094	0.0068	-0.0000	-0.0073	-0.0106	-0.0078	0.0000	0.008	3
	0.0122	0.0090	-0.0000	-0.0098	-0.0145	-0.0107	0.0000	0.011	8
	0.0177	0.0132	-0.0000	-0.0150	-0.0227	-0.0173	0.0000	0.020	5
	0.0318	0.0250	-0.0000	-0.0322	-0.0531	-0.0450	0.0000	0.075	0
	0.1592	0.2251	NaN	0.2251	0.1592	0.0750	0.0000	-0.0450) -
	0.0531	-0.0322	-0.0000	0.0250	0.0318	0.0205	0.0000	-0.017	3 -
	0.0227	-0.0150	-0.0000	0.0132	0.0177	0.0118	0.0000	-0.0107	7
	-0.014	5 -0.009	98 -0.000	0.009	0.012	0.008	0.000	0 -0.007	78 -
	0.0106	-0.0073	-0.0000	0.0068	0.0094	0.0064	0.0000	-0.006	1 -
	0.0084	-0.0058	-0.0000	0.0055	0.0076	0.0052	0.0000	-0.005	- 0
	0.0069	-0.0048	-0.0000	0.0046	0.0064]				
					-				
0	4	•	0	4	-		- 0	0	10

Then take the *z* transform $H_t(z) = \sum_{n=-10}^{10} (n) z^{-n}$. Then determine the transfer function, H(z), of the realizable FIR filter as $H(z) = z^{-10} H_t(z)$ from which the filter structure can be drawn. The frequency response is compared below for 21 and 101 coefficients. The 21-tap filter is not that bad.

%21-tap filter coefficients b21=[0.0318 0.0250 -0.0000 -0.0322 -0.0531 -0.0450 0.0000 0.0750 0.1592 0.2251 0.25 0.2251 0.1592 0.0750 0.0000 -0.0450 -0.0531 -0.0322 -0.0000 0.0250 0.0318], a=[1] %101-tap filter coefficients b101=[0.0064 0.0046 -0.0000 -0.0048 -0.0069 -0.0050 0.0000 0.0052 0.0055 -0.0000 -0.0058 -0.0084 -0.0061 0.0076 0.0000 0.0064 0.0094 0.0068 -0.0000 -0.0073 -0.0106 -0.0078 0.0000 0.0083 0.0122 0.0090 -0.0000 -0.0098 -0.0145 -0.0107 0.0000 0.0118 0.0177 0.0132 -0.0000 -0.0150 -0.0227 -0.0173 0.0000 0.0205 0.0318 0.0250 -0.0000 -0.0322 -0.0531 -0.0450 0.0000 0.0750 0.1592 0.2251 0.25 0.2251 0.1592 0.0750 0.0000 -0.0450 -0.0531 -0.0322 -0.0000 0.0250 0.0318 0.0205 0.0000-0.0173 -0.0227 -0.0150 -0.0000 0.0132 0.0177 0.0118 0.0000 -0.0107 -0.0145 -0.0098 -0.0000 0.0090 0.0122 0.0083 0.0000 -0.0078 -0.0106 -0.0068 0.0094 0.0064 0.0000 -0.0061 -0.0084 -0.0073 -0.0000 0.0058 -0.0000 0.0076 0.0052 0.0000 -0.0050 -0.0069 -0.0055 0.0048 -0.0000 0.0046 0.0064],w=-pi: pi/256: pi; Hw21=freqz(b21, a, w); Hw101=freqz(b101, a, w); plot(w, abs(Hw21), w, abs(Hw101), 'k') legend ('21 coefficients', '101 coefficients'); title('Magnitude Response'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)');



Example 4.4.3 [Very narrow band pass filter] [2003] Design a band pass FIR filter that approximates the following frequency response:

 $H(F) = 1, \quad 160 \le F \le 200 \text{ Hz}$ 0, elsewhere in $0 \le F \le F_s/2$

when the sampling frequency is 8000 sps. Limit the duration of impulse response to 2 msec. Draw the filter structure.

Solution The sampling frequency $F_s = 8000$ Hz corresponds to $\omega = 2\pi$ rad. Thus

160 Hz corresponds to $\omega = (160/8000)2\pi = 0.04 \ \pi$ rad., and 200 Hz corresponds to $\omega = (200/8000)2\pi = 0.05 \ \pi$ rad.



$$H_{d}(e^{j\theta}) = 1, \qquad -0.05\pi \le \omega \le -0.04\pi \quad \text{and} \qquad 0.04\pi \le \omega \le 0.05\pi$$

$$h_{d}(n) = \frac{1}{-1} \int H_{d}(e^{j\theta}) e^{j\theta m} d\omega = \frac{1}{2\pi} \left\{ \int e^{j\theta n} \int e^{j\theta n} d\omega + \int e^{j\theta n} \int e^{j\theta n} d\omega \right\}$$

$$= \frac{1}{2\pi} \left\{ \int e^{j\theta n} \int e^{j\theta n} \int e^{j\theta n} d\omega + \int e^{j\theta n} \int e^{j\theta n} d\omega + \int e^{j\theta n} \int e^{j\theta n} d\omega \right\}$$

$$= \frac{1}{2\pi} \left\{ \int e^{j\theta n} d\omega + \int e^{j\theta n} \int e^{j\theta n} d\omega + \int e^{j\theta n} \int e^{j\theta n} d\omega \right\}$$

$$= \frac{1}{2\pi} \left\{ \int e^{j\theta n} d\omega + \int e^{j\theta n} \int e^{j\theta n} d\omega + \int e^{j\theta n} d\omega \right\}$$

$$= \frac{1}{2\pi} \left\{ \int e^{j\theta n} d\omega + \int e^{j\theta n}$$

Filter length *N* is determined by the duration of the impulse response. Two milliseconds corresponds to 0.002 / (1/8000) = 16 sample periods. This means that the filter length N = 17, and there will be 17 coefficients. Determine the values of $h_d(n)$ for $-8 \le n \le 8$, so that

$$h_t(n) = \{ \underset{k=-8}{\text{h}_d} (-8), \underset{k=-8}{\text{h}_d} (-7), \ldots, \underset{k=0}{\text{h}_d} (0), \ldots, \underset{k=0}{\text{h}_d} (8) \}, \text{ and}$$

 $H_t(z) = \sum_{n=-8}^{8} h_t(n) z^{-n}$

Delay the impulse response by 8 sample periods so that

$$h(n) = h_t(n-8)$$
 and $H(z) = z^{-8}H_t(z)$

 $h_d(0) = 0.01$ by L'Hopital's rule

%Calculate hdn = $(\sin(0.05*pi*n) - \sin(0.04*pi*n)) / (pi*n)$ and stem plot n = -50: 50, hdn = $(\sin(0.05*pi*n) - \sin(0.04*pi*n)) . / (pi*n)$, stem(n, hdn)

 $\begin{array}{l} n=-50\ to\ 50\\ Warning:\ Divide\ by\ zero.\\ hdn=& [0.0064\quad 0.0072\quad 0.0080\quad 0.0085\quad 0.0089\quad 0.0092\quad 0.0092\quad 0.0091\\ 0.0087\quad 0.0082\quad 0.0076\quad 0.0067\quad 0.0058\quad 0.0047\quad 0.0035\quad 0.0022\quad 0.0009\\ -0.0005\quad -0.0018\quad -0.0031\quad -0.0044\quad -0.0056\quad -0.0066\quad -0.0076\quad -0.0084\quad -\\ 0.0090\quad -0.0095\quad -0.0097\quad -0.0098\quad -0.0097\quad -0.0094\quad -0.0088\quad -0.0082\quad -\\ \end{array}$

	0.0073	-0.0063	-0.0052	-0.0039	-0.002	.6 -0.00	0.0	0002 0.	0016	
	0.0029	0.0042	0.0055	0.0066	0.0076	6 0.00	84 0.00	091 0.0	096 0.0	099
	NaN	0.0099	0.0096	0.0091	0.0084	0.0076	0.006	6 0.005	55 0.004	<i>42</i>
	0.0029	0.0016	0.0002	-0.0012	-0.002	6 -0.00	0.0-0.0	052 -0.	0063 -	
	0.0073	-0.0082	-0.0088	-0.0094	4 -0.009	-0.00	0.0- 0.0	0097 -0.	.0095 -	
	0.0090	-0.0084	-0.0076	-0.0066	5 -0.005	6 -0.00	0.0-0.0	0031 -0.	.0018 -	
	0.0005	0.0009	0.0022	0.0035	0.0047	7 0.00	58 0.00	0.0	076 0.0	082
	0.0087	0.0091	0.0092	0.0092	0.0089	9 0.00	85 0.00	0.0 080	072	
	0.0064]								
<i>n</i> =	0	±1	±2	± 3	±4	± 5	±6	±7	± 8	
$h_d(n) =$	0.01 (0.0099 ().0096 0	.0091 0	.0084 ().0076	0.0066	0.0055	0.0042	

The frequency responses 17-tap and 101-tap filters are shown below. The 17-tap filter looks more like a low pass filter! Owing to the very narrow pass band a very large number of coefficients is needed before the pass band becomes discernible. In general FIR filters are characterized by a large number of coefficients compared to IIR filters.

%Filter coefficients

b17=[0.0042 0.0055 0.0066 0.0076 0.0084 0.0091 0.0096 0.0099 0.01 0.0099 0.0096 0.0091 0.0084 0.0076 0.0066 0.0055 0.0042], a=[1] $b101 = [0.0064 \quad 0.0072 \quad 0.0080 \quad 0.0085 \quad 0.0089 \quad 0.0092 \quad 0.0092 \quad 0.0091$ 0.0087 0.0082 0.0076 0.0067 0.0058 0.0047 0.0035 0.0022 0.0009 -0.0018 -0.0031 -0.0044 -0.0056 -0.0066 -0.0076 -0.0084 --0.0005 0.0090 -0.0095 -0.0097 -0.0098 -0.0097 -0.0094 -0.0088 -0.0082 -0.0073 -0.0063 -0.0052 -0.0039 -0.0026 -0.0012 0.0002 0.0016 0.0029 0.0042 0.0055 0.0066 0.0076 0.0084 0.0091 0.0096 0.0099 0.01 0.0099 0.0096 0.0091 0.0084 0.0076 0.0066 0.0055 0.0042 0.0029 0.0016 0.0002 -0.0012 -0.0026 -0.0039 -0.0052 -0.0063 -0.0073 -0.0082 -0.0088 -0.0094 -0.0097 -0.0098 -0.0097 -0.0095 -0.0090 -0.0084 -0.0076 -0.0066 -0.0056 -0.0044 -0.0031 -0.0018 -0.0005 0.0009 0.0022 0.0035 0.0047 0.0058 0.0067 0.0076 0.0082 0.0087 0.0091 0.0092 0.0092 0.0089 0.0085 0.0080 0.0072 0.00641 w=-pi: pi/256: pi; Hw17=freqz(b17, a, w); Hw101=freqz(b101, a, w); subplot(2, 1, 1), plot(w, abs(Hw17)); legend ('17 coefficients');

xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, abs(Hw101)); legend ('101 coefficients');

 $xlabel('Frequency \omega, rad/sample'), \ ylabel('Magnitude \ of \ H(\omega)'); \ grid$



Example 4.4.4 [Band pass filter] [2003] Design a band pass filter to pass frequencies in the range 1 to 2 rad/sec using **Hanning window** with N = 5. Draw the filter structure and plot its spectrum.

Solution The Hanning window is also known as the **Hann window**. It is a raised cosine, is very similar to the Hamming window, and is given by

$$w_{Han}(n) = 0.5 \{ 1 - \cos(2\pi n/(N-1)) \}, \quad 0 \le n \le N-1$$

0, elsewhere

Note In order to convert the analog frequencies to digital we need the sampling time *T*. The sampling frequency F_s (or the sampling time *T*) is not specified. We assume T = 1 sec or, what amounts to the same, we assume that the frequencies given are actually digital, that is, 1 to 2 rad/sample instead of 1 to 2 rad/sec. However, the solution below assumes that the frequencies are given correctly, that is, that they are analog, and uses a sampling frequency of $\Omega_s = 4$ rad / sec.

Although there is a specialized version of the sampling theorem for band pass signals we shall simply take the sampling frequency to be twice the highest frequency which is 2 rad / sec. Thus we shall take $\Omega_s = 4$ rad / sec. This then gives us a *high pass* filter rather a band pass. However, if we take, say, $\Omega_s = 8$ rad / sec., we shall have a *band pass* filter. We shall next convert the analog frequency specs to digital (ω):

 $\Omega_s = 4 \text{ rad} / \text{sec corresponds to } \omega = 2\pi \text{ rad}$ 1 rad / sec corresponds to $\omega = 2\pi/4 = \pi/2 \text{ rad}$ 2 rad / sec corresponds to $\omega = (2\pi/4) 2 = \pi \text{ rad}$ Thus

$$H_{d}(\omega) = 1, \quad -\pi \le \omega \le -\pi/2 \text{ and } \pi/2 \le \omega \le \pi$$

0, elsewhere in $[-\pi, \pi]$

$$I = \frac{|H_{d}(\omega)|}{-\pi - \pi/2 - 0} \quad Take \angle H_{d}(\omega) = 0$$

$$I = \frac{1}{-\pi} \int_{-\pi}^{\pi/2} H(\omega) e^{-j\omega n} d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi/2} \frac{1}{1 - 1} e^{-j\omega n} d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi/2} \frac{1}{1 - 1} e^{-j\omega n} d\omega = \frac{1}{2\pi} \int_{-\pi/2}^{\pi/2} \frac{1}{1 - 1} e^{-j\omega n} d\omega = \frac{1}{2\pi} \int_{-\pi/2}^{\pi/2} \frac{1}{1 - 1} e^{-j\pi/2} \int_{-\pi/2}^{\pi/2} \frac{1}{1 - 1} e^{-j\pi/2} e^{-j\pi/2} \int_{-\pi/2}^{\pi/2} \frac{1}{1 - 1} e^{-j\pi/2} e^{-j\pi/2} \int_{-\pi/2}^{\pi/2} \frac{1}{1 - 1} e^{-j\pi/2} e^{-$$

 $h_d(0) = 0.5$ by L'Hopital's rule

Determine $h_d(n)$ for $-2 \le n \le 2$, so that

$$h_t(n) = \{h_d(-2), h_d(-1), h_d(0), h_d(1), h_d(2)\}$$

For N = 5 the Hanning window is given by

$$w_{Han}(n) = 0.5 \left\{ 1 - \cos\left(2\pi n 75 - 1\right) \right\}, \ 0 \le n \le 5 - 1$$

= 0.5 \{ 1 - \cos(\pi n / 2) \}, \quad 0 \le n \le 4
Thus w(n) = \{0, 0.5, 1, 0.5, 0\}. Multiplying h_t(n) and w_{Han}(n) point by point we get
h_t(n) = h_t(n) w_{Han}(n) = \{0, \frac{d}{2}, h_d(0), \frac{d}{2}, 0\} and
H_t(z) = \sum_{n=-2}^2 h_t(n) z^{-n}

Delay by 2 samples to get $h(n) = h_t(n-2)$ and $H(z) = z^{-2} H_t(z)$. Now draw the direct form structure for H(z). The spectrum is given by $H(e^{j\omega}) = H(z)|_{z^{j\omega}e}$. We compare below filters lengths of 5 and 101 without the Hanning window.

Generate filter coefficients:

%Calculate hdn = $(\sin(pi*n) - \sin(pi*n/2))/(pi*n)$ and stem plot n = -50: 50, hdn = $(\sin(pi*n) - \sin(pi*n/2))./(pi*n)$, stem(n, hdn)

n = -50 to 50										
Warning: Divide by zero.										
hdn = [0]).0000 -(0.0065 -().0000 ().0068 -0	.0000 -0).0071 -(0.0000	0.0074		
0.0000	-0.0078	-0.0000	0.0082	-0.0000	-0.0086	-0.0000	0.0091	l		
0.0000	-0.0096	-0.0000	0.0103	-0.0000	-0.0110	-0.0000	0.0118	3		
0.0000	-0.0127	-0.0000	0.0138	-0.0000	-0.0152	-0.0000	0.0168	3 -		
0.0000	-0.0187	-0.0000	0.0212	-0.0000	-0.0245	-0.0000	0.0289) -		
0.0000	-0.0354	-0.0000	0.0455	-0.0000	-0.0637	-0.0000	0.1061	_		
0.0000	-0.3183	NaN	-0.3183	-0.0000	0.1061	-0.0000	-0.0637	-		
0.0000	0.0455	-0.0000	-0.0354	-0.0000	0.0289	-0.0000	-0.0245	5 -		
0.0000	0.0212	-0.0000	-0.0187	-0.0000	0.0168	-0.0000	-0.0152	-		
0.0000	0.0138	-0.0000	-0.0127	0.0000	0.0118	-0.0000	-0.0110	-		
0.0000	0.0103	-0.0000	-0.0096	0.0000	0.0091	-0.0000	-0.0086			
0.0000	0.0082	-0.0000	-0.0078	0.0000	0.0074	-0.0000	-0.0071	-		
0.0000	0.0068	-0.0000	-0.0065	0.0000]						

Generate frequency responses:

%5-tap filter coefficients b5=[0.0000 -0.3183 0.5 -0.3183 -0.0000], a=[1] %101-tap filter coefficients b101=[0.0000 -0.0065 -0.0000 0.0068 -0.0000 -0.0071 -0.0000 0.0074 0.0000 -0.0078 -0.0000 0.0082 -0.0000 -0.0086 -0.0000 0.0091 0.0000 -0.0096 -0.0000 0.0103 -0.0000 -0.0110 -0.0000 0.0118 0.0000 -0.0127 -0.0000 0.0138 -0.0000 -0.0152 -0.0000 0.0168 -0.0000 -0.0187 -0.0000 0.0212 -0.0000 -0.0245 -0.0000 0.0289 -0.0000 -0.0354 -0.0000 0.0455 -0.0000 -0.0637 -0.0000 0.1061 -0.0000 -0.3183 0.5 -0.3183 -0.0000 0.1061 -0.0000 -0.0637 -0.0000 0.0455 -0.0000 -0.0354 -0.0000 0.0289 -0.0000 -0.0245 -0.0000 0.0212 -0.0000 -0.0187 -0.0000 0.0168 -0.0000 -0.0152 -0.0000 0.0138 -0.0000 -0.0127 0.0000 0.0118 -0.0000 -0.0110 -0.0000 0.0103 -0.0000 -0.0096 0.0000 0.0091 -0.0000 -0.0086 -0.0000 -0.0000 -0.0078 0.0000 0.0074 -0.0000 -0.0071 -0.0000 0.0082 0.0068 -0.0000 -0.0065 0.0000], w=-pi: pi/256: pi; Hw5=freqz(b5, a, w); Hw101=freqz(b101, a, w); plot(w, abs(Hw5), w, abs(Hw101), 'k') legend ('5 coefficients', '101 coefficients'); title('Magnitude Response'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid



We compare below the frequency responses of the 5-tap filter with and without the Hanning window. It can be seen that the Hanning window aggravates what is already a poor (short) filter length. There may be more to gain by increasing the filter length than by windowing.

Generate Hanning window:

%Generate Hanning window wn = $0.5*(1-\cos(pi*n/2))$ and stem plot n = -2: 2, wn = $0.5*(1-\cos(pi*n/2))$, stem(n, wn)

 $n = -2 \text{ to } 2 \\ wn = 1.0000 \quad 0.5000 \qquad 0 \quad 0.5000 \quad 1.0000$

Generate frequency responses:

%5-tap filter coefficients b5=[0.0000 -0.3183 0.5 -0.3183 -0.0000], %Hanning window coefficients wn = [1.0000 0.5000 0 0.5000 1.0000], %Windowed coefficients b5Han = b5 .*wn, a=[1] w=-pi: pi/256: pi; Hw5=freqz(b5, a, w); Hw5Han=freqz(b5Han, a, w); plot(w, abs(Hw5), w, abs(Hw5Han), 'k') legend ('5 coefficients, No window', '5 coefficients, Hanning window ');



Example 4.4.5 [High pass filter] [2008] Design a high pass linear phase filter with frequency response

$$H_{d}(e^{j\omega}) = 1 e^{-j3\omega}, \ \omega_{c} \le \omega |\le |\pi|$$

$$0, \qquad \text{elsewhere in the range } [-\pi \text{to } \pi]$$

The number of filter coefficients is N = 7 and $\omega_c = \pi/4$. Use (a) rectangular window and (b) Hamming window.



This sequence is centered at n = 3. Since the filter length is N = 7 we can calculate the 7 coefficients as $\{h_d(n), 0 \le n \le 6\}$. In other words we truncate it outside the interval $0 \le n \le 6$; moreover, there is no need to right-shift the truncated sequence.

In general, one may not know the filter length with certainty and there is no special advantage in specifying the phase, $\angle H_d(\omega)$, as anything but zero. We calculate 101 coefficients of the sequence centered about n = 0, that is, $h_d(n) = \frac{\sin(\pi n) - \sin(\pi n / 4)}{\sin(\pi n) - \sin(\pi n / 4)}$:

 πn

 $h_d(0) = 0.75$ by L'Hopital's rule

% Generate hdn = $(\sin(pi*n) - \sin(pi*n/4)) / (pi*n)$ and stem plot n = -50: 50, hdn = $(\sin(pi*n) - \sin(pi*n/4)) . / (pi*n)$, stem(n, hdn)

n = -50 to 50

Warning: Divide by zero.

-		-						
hdn = [-	0.0064 -0	.0046 -0.	0000	0.0048	0.0069	0.0050 -0	.0000 -0.0	052
-0.0076	-0.0055	-0.0000	0.0058	0.0084	0.0061	-0.0000	-0.0064 -	
0.0094	-0.0068	-0.0000	0.0073	0.0106	0.0078	-0.0000	-0.0083 -	
0.0122	-0.0090	-0.0000	0.0098	0.0145	0.0107	-0.0000	-0.0118	-
0.0177	-0.0132	-0.0000	0.0150	0.0227	0.0173	-0.0000	-0.0205	-
0.0318	-0.0250	-0.0000	0.0322	0.0531	0.0450	-0.0000	-0.0750	-
0.1592	-0.2251	NaN	-0.2251	-0.1592	-0.0750	-0.0000	0.0450	
0.0531	0.0322	-0.0000	-0.0250	-0.0318	-0.0205	-0.0000	0.0173	
0.0227	0.0150	-0.0000	-0.0132	-0.0177	-0.0118	-0.0000	0.0107	
0.0145	0.0098	-0.0000	-0.0090	-0.0122	-0.0083	-0.0000	0.0078	
0.0106	0.0073	-0.0000	-0.0068	-0.0094	-0.0064	-0.0000	0.0061	
0.0084	0.0058	-0.0000	-0.0055	-0.0076	-0.0052	-0.0000	0.0050	
0.0069	0.0048	-0.0000	-0.0046	-0.0064	1			

The 7-tap and 101-tap filter responses are shown below

%7-tap filter coefficients b7=[-0.0750 -0.1592 -0.2251 0.75 -0.2251 -0.1592 -0.0750], a=[1] %101-tap filter coefficients b101 = [-0.0064 - 0.0046 - 0.0000 0.0048 0.0069]0.0050 -0.0000 -0.0052 -0.0076 -0.0055 -0.0000 0.0058 0.0084 0.0061 -0.0000 -0.0064 -0.0094 -0.0068 -0.0000 0.0073 0.0106 0.0078 -0.0000 -0.0083 -0.0122 -0.0090 -0.0000 0.0098 0.0145 0.0107 -0.0000 -0.0118 0.0177 -0.0132 -0.0000 0.0150 0.0227 0.0173 -0.0000 -0.0205 -0.0318 -0.0250 -0.0000 0.0322 0.0531 0.0450 -0.0000 -0.0750 --0.0000 0.0450 0.1592 -0.2251 0.75 -0.2251 -0.1592 -0.0750 0.0531 0.0322 -0.0000 -0.0250 -0.0318 -0.0205 -0.0000 0.0173 0.0227 0.0150 -0.0000 -0.0132 -0.0177 -0.0118 -0.0000 0.0107 0.0145 0.0098 -0.0000 -0.0090 -0.0122 -0.0083 -0.0000 0.0078 0.0073 -0.0000 -0.0068 -0.0094 -0.0064 -0.0000 0.0106 0.0061 0.0084 0.0058 -0.0000 -0.0055 -0.0076 -0.0052 -0.0000 0.0050 0.0069 0.0048 -0.0000 -0.0046 -0.0064], w=-pi: pi/256: pi; Hw7=freqz(b7, a, w); Hw101=freqz(b101, a, w);plot(w, abs(Hw7), w, abs(Hw101), 'k') legend ('7 coefficients', '101 coefficients'); title('Magnitude Response');

xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid



Example 4.4.6 [Band-stop filter] Design a band-stop linear phase filter with the following frequency response. The number of filter coefficients is N = 31.



Example 4.4.7 [Design of 9-coefficient narrow-band LP FIR filter] This is a repeat of Example 2.1 with the bandwidth reduced to 0.02π . The band width is *narrower* by a factor of 10 from the band width in that example. The objective is to show that the narrower the band width the larger the number of coefficients needed to achieve the specified frequency response.

Design a nine-coefficient (or 9-point or 9-tap) FIR digital filter to approximate an ideal low-pass filter with a cut-off frequency $\omega_c = 0.02\pi$. The magnitude response, $|H_d(\omega)|$, is given below. Take $\angle H_d(\omega) = 0$.



$$h_{d}(0) = \frac{\frac{d (\sin 0.02 \pi n)}{dn}}{\frac{d (\pi n)}{dn}} \bigg|_{n=0} = \frac{0.02 \pi \cos 0.02 \pi n}{\pi} \bigg|_{n=0} = 0.02$$

The MATLAB calculation of coefficients follows (good for all *n* except 0 where we fill in the value 0.02):

% Generate hdn = $(\sin (0.02*pi*n)) / (pi*n)$ and stem plot n = -50: 50, hdn = $(\sin(0.02*pi*n)) ./(pi*n)$, stem(n, hdn) xlabel('n'), ylabel('hd(n)'); grid; title ('hd(n) = $(\sin (0.02*pi*n)) / (pi*n)')$

n = -50 to 50 Warning: Divide by zero. hdn = -0.0000 0.0004 0.0008 0.0013 0.0017 0.0022 0.0027 0.0032 0.0037 0.0042 0.0047 0.0052 0.0057 0.0063 0.0068 0.0074 0.0079 0.0085 0.0090 0.0095 0.0101 0.0106 0.0122 0.0127 0.0112 0.0117 0.0132 0.0137 0.0142 0.0147 0.0151 0.0156 0.0160 0.0164 0.0168 0.0182 0.0184 0.0172 0.0175 0.0178 0.0187 0.0190 0.0192 0.0194 0.0195 0.0200 0.0197 0.0198 0.0199 0.0199 0.0200 NaN 0.0199 0.0199 0.0198 0.0197 0.0195 0.0194 0.0192 0.0190 0.0187 0.0184 0.0182 0.0178 0.0175 0.0172 0.0168 0.0164 0.0160 0.0156 0.0151 0.0137 0.0132 0.0147 0.0142 0.0127 0.0122 0.0117 0.0112 0.0106 0.0101 0.0095 0.0090 0.0085 0.0079 0.0074 0.0068 0.0063 0.0057 0.0032 0.0052 0.0047 0.0042 0.0037 0.0027 0.0022 0.0017 0.0013 0.0008 0.0004 -0.0000

The following MATLAB plot of coefficients is good for all n except at n = 0.



MATLAB plots of magnitude and phase response of the 9-coefficient narrower band LP filter follow:

%Magnitude and phase response of 9-coefficient LP filter %Filter coefficients b9ex7= [0.0198 0.0199 0.0199 0.0200 0.02 0.0200 0.0199 0.0199 0.0198], a=[1] w=-pi: pi/256: pi; Hw9ex7=freqz(b9ex7, a, w); subplot(2, 1, 1), plot(w, abs(Hw9ex7)); legend ('9 coefficients'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw9ex7)); legend ('9 coefficients'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid



The MATLAB plots below enable us to compare (only visually) the magnitude responses of the two 9-coefficient filters, the only difference being that Ex 1 has a band width of 0.2π while Ex 2 is very narrow at 0.02π .

% Magnitude responses compared: Ex 1 and Ex 7 (9-coefficient LP filters) %Filter coefficients b9ex1=[0.0468, 0.1009, 0.1514, 0.1871, 0.2, 0.1871, 0.1514, 0.1009, 0.0468], b9ex7= [0.0198 0.0199 0.0199 0.0200 0.02 0.0200 0.0199 0.0199 0.0198], a=[1] w=-pi: pi/256: pi; Hw9ex1=freqz(b9ex1, a, w); Hw9ex7=freqz(b9ex1, a, w); plot(w, abs(Hw9ex1), w, abs(Hw9ex7), 'k') legend ('Ex1, 9 coefficients', 'Ex7, 9 coefficients'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid

From the plots it clear that the narrow-band filter of Example 6 is no where near either the band width or the gain specified. In contrast, the wider band width filter of Example 1 is relatively much closer to specification.



In the MATLAB plots below we have increased the number of coefficients to 51 for both filters, the only difference being that Ex 1 has a band width of 0.2π while Ex 2 is very narrow at 0.02π .

% Magnitude responses compared: Ex 1 and Ex 7 (51-coefficient LP filters) %Filter coefficients b51ex1=[0.0, 0.0078 0.0132 0.0138 0.0089 -0.0, -0.0098 -0.0168 -0.0178 -0.0117 0.0, 0.0134 0.0233 0.0252 0.0170 -0.0, -0.0208 -0.0378 -0.0432 -0.0312 0.0, 0.0468 0.1009 0.1514 0.1871 0.2 0.1871 0.1514 0.1009 0.0468 0.0 -0.0312 -0.0432 -0.0378 -0.0208 -0.0 0.0170 0.0252, 0.0233 0.0134 0.0 -0.0117 -0.0178 -0.0168 -0.0098 -0.0 0.0089, 0.0138 0.0132 0.0078 0.0], % b51ex7= [0.0127 0.0132 0.0137 0.0142 0.0147 0.0151 0.0156 0.0160 0.0164 0.0168 0.0172 0.0175 0.0182 0.0184 0.0178 0.0187 0.0190 0.0192 0.0194 0.0195 0.0197 0.0198 0.0199 0.0199 0.0200 $0.02 \quad 0.0200 \quad 0.0199 \quad 0.0199 \quad 0.0198 \quad 0.0197 \quad 0.0195 \quad 0.0194 \quad 0.0192$ 0.0190 0.0187 0.0184 0.0182 0.0178 0.0175 0.0172 0.0168 0.0164 0.0160 0.0156 0.0151 0.0147 0.0142 0.0137 0.0132 0.0127],a=[1] w=-pi: pi/256: pi; Hw51ex1=freqz(b51ex1, a, w); Hw51ex7=freqz(b51ex7, a, w); plot(w, abs(Hw51ex1), w, abs(Hw51ex7), 'k') legend ('Ex1, 51 coefficients', 'Ex7, 51 coefficients'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid

From the plots it clear that the narrow-band filter of Example 6 still has a long way to go. In contrast, the wider band width filter of Example 1 is almost there.



We turn now to **filtering plus quadrature phase shift**. These applications include integrators, differentiators, and Hilbert transform devices. For all of these the desired transfer function is imaginary, i.e., $H(e^{j\omega}) = j H_l(\omega)$, with $H_R(\omega) = 0$

Example 4.4.8 [The ideal differentiator] In the analog situation the ideal differentiator is given by the transfer function H(s) = s with the frequency response $H(\Omega) = j\Omega$, for $0 \le \Omega \le \infty$. The digital version of the ideal differentiator may be defined as

$$H_d(e^{j\omega})$$
 or $H_d(\omega) = j\omega$, $-\pi \le \omega \le \pi$

Note $H(\omega) = j\omega = \omega e^{j\pi/2}$ has a magnitude of $H|(\omega) = \omega$, and a phase $\angle H(\omega)_d = \pm \pi/2$ (see figure).



Since $H_d(\omega)$ is periodic in ω with period 2π , we can expand in into a Fourier series as

$$H_{d}(\omega) = \sum_{n = -\infty}^{\infty} h_{d}(n) e^{-j\omega n}$$

where the Fourier coefficients $h_d(n)$ (the impulse response) are given by

$$h_d(n) = \underbrace{\int}_{2\pi} H_d(\omega) e^{-d\omega} d\omega$$

Substituting $H_d(\omega) = j\omega$, we get $j\omega e^{j\omega n}d\omega = j \int \omega$

$$h_d(n) = \underbrace{\int}_{2\pi} \underbrace{\frac{e^{j\omega n}}{jn}}_{-\pi} - \underbrace{\int}_{-\pi} \underbrace{\frac{e^{j\omega n}}{jn}}_{-\pi} d\omega \Big|_{-\pi}$$

٦

$$= \frac{i}{\left| \frac{\pi e^{-(-\pi)e}}{\pi^{n}} - \frac{e^{-j\pi n}}{\pi^{n}} - \frac{i\pi}{\pi^{n}} \right|^{\pi}} \left| \frac{2\pi \left| \frac{jn}{\pi^{n}} - \frac{jn}{\pi^{n}} - \frac{j\pi^{n}}{\pi^{n}} \right|^{\pi}}{(jn)^{2}} \right|_{-\pi} \left| \frac{2\pi \left| \frac{jn}{\pi^{n}} - e^{-j\pi n} \right|}{2\pi \left| \frac{j^{2}n^{2}}{\pi^{n}} \right|^{2}} \right|^{2\pi}}$$
$$= \frac{\cos n\pi}{n} - \frac{\sin n\pi}{\pi^{n}} = \frac{\cos n\pi}{n} - 0, \quad n = \text{non-zero integers}$$

where we have used the fact that $\frac{1}{\pi n^2} = 0$ for non-zero integer values of *n*. For n = 0 the value

of $h_d(0)$ is evaluated from the defining equation, viz.,

$$h_{d}(0) = \frac{1}{2\pi} \int_{-\pi}^{\pi} j\omega \ e^{j\omega 0} d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi} j\omega \ d\omega = \frac{J}{2\pi} \frac{J}{2\pi} \frac{\omega^{2}}{2} \Big|_{-\pi} = 0$$

Thus

$$h_d(n) = \begin{cases} \frac{\cos n\pi}{n}, & n = \text{non-zero integers} \\ 0, & n = 0 \end{cases}$$

As a specific case choose a filter length N = 9 and evaluate $h_d(n)$ for $-4 \le n \le 4$.

%Ideal Differentiator % Generate hdn = $(\cos(n^*pi)) / (n)$ and stem plot n = -50: 50, hdn = (cos(n*pi)) ./(n), stem(n, hdn) xlabel('n'), ylabel('hd(n)'); grid; title ('hd(n) = $(\cos(n*pi)) / (n)')$ n = -50 to 50 Warning: Divide by zero. $hdn = [-0.0200 \quad 0.0204 \quad -0.0208 \quad 0.0213 \quad -0.0217 \quad 0.0222 \quad -0.0227 \quad 0.0233$ $-0.0238 \quad 0.0244 \quad -0.0250 \quad 0.0256 \quad -0.0263 \quad 0.0270 \quad -0.0278 \quad 0.0286 \quad -0.0278 \quad 0.0286 \quad -0.0278 \quad 0.0286 \quad -0.0286 \quad$ $0.0294 \quad 0.0303 \quad -0.0313 \quad 0.0323 \quad -0.0333 \quad 0.0345 \quad -0.0357 \quad 0.0370 \quad -0.0370 \quad$ $0.0385 \quad 0.0400 \quad -0.0417 \quad 0.0435 \quad -0.0455 \quad 0.0476 \quad -0.0500 \quad 0.0526 \quad -0.0500 \quad 0.0500 \quad -0.0500 \quad 0.0500 \quad -0.0500 \quad -0.0500$ 0.0556 0.0588 -0.0625 0.0667 -0.0714 0.0769 -0.0833 0.0909 -0.1000 0.1111 -0.1250 0.1429 -0.1667 0.2000 -0.2500 0.3333 -0.5000 1.0000 *0 -1.0000 0.5000 -0.3333 0.2500 -*0.2000 0.1667 -0.1429 0.1250 -0.1111 0.1000 -0.0909 0.0833 -0.0769 0.0714 -0.0667 0.0625 -0.0588 0.0556 -0.0526 0.0500 -0.0476 0.0455 -0.0435 0.0417 -0.0400 0.0385 -0.0370 0.0357 -0.0345 0.0333 -0.0323 0.0313 -0.0303 0.0294 -0.0286 0.0278 -0.0270 0.0263 -0.0256 0.0250 -0.0244 0.0238 -0.0233 0.0227 -0.0222 0.0217 -0.0213 0.0208 -0.0204 0.0200]

We can see the odd symmetry of $h_d(n)$ from the following stem plot (note that the correct value of $h_d(0) = 0$; the MATLAB segment used here has $h_d(0) = \infty$ which is not correct).



The frequency response

%Ideal Differentiator %9-tap filter coefficients 0 -1.0000 0.5000 -0.3333 b9=[-0.2500 0.3333 -0.5000 1.0000 0.2500], a=[1] %101-tap filter coefficients b101=[-0.0200 0.0204 -0.0208 0.0213 -0.0217 0.0222 -0.0227 0.0233 -0.0238 0.0244 -0.0250 0.0256 -0.0263 0.0270 -0.0278 0.0286 -0.0294 0.0303 -0.0313 0.0323 -0.0333 0.0345 -0.0357 0.0370 -0.0385 0.0400 -0.0417 0.0435 -0.0455 0.0476 -0.0500 0.0526 -0.0556 0.0588 -0.0625 0.0667 -0.0714 0.0769 -0.0833 0.0909 -0.1000 0.1111 -0.1250 0.1429 -0.1667 0.2000 -0.2500 0.3333 -0 -1.0000 0.5000 -0.3333 0.5000 1.0000 0.2500 -0.2000 0.1667 -0.1429 0.1250 -0.1111 0.1000 -0.0909 0.0833 -0.0769 0.0714 -0.0625 -0.0588 0.0556 -0.0526 0.0500 -0.0476 0.0667 0.0455 -0.0435 0.0385 -0.0370 0.0357 -0.0345 0.0333 -0.0417 -0.0400 0.0323 0.0313 -0.0303 0.0294 -0.0286 0.0278 -0.0270 0.0263 -0.0256 0.0250 -0.0244 0.0238 -0.0233 0.0227 -0.0222 0.0217 -0.0213 0.0208 -0.0204 0.0200], w=-pi: pi/256: pi; Hw9=freqz(b9, a, w); Hw101=freqz(b101, a, w);plot(w, abs(Hw9), w, abs(Hw101), 'k') legend ('9 coefficients', '101 coefficients');





Example 4.4.9 [The Hilbert transformer] This is used to generate signals that are in *phase* quadrature to an input sinusoidal signal (or, more generally, an input narrowband waveform). That is, if the input to a Hilbert transformer is the signal $x_a(t) = \cos \Omega_0 t$, the output is $y_a(t) = \sin \theta$ $\Omega_0 t$. The Hilbert transformer is used in communication systems in various modulation schemes. The frequency response of the Hilbert transformer is (Figure)

 $H_{\mathcal{A}}(e^{j\omega}) = -j \operatorname{sgn}(\omega), -\pi \le \omega \le \pi$ where the $sgn(\omega)$ is the signum function defined as $|H_d(\omega)|$ 1 if ω is positive $sgn(\omega) =$ -1 if ω is negative ω -π π $\angle H_d(\omega)$ $\pi/2$ (1)π π Since $j = e^{j\pi/2}$ we may also express $H(e^{j\omega})$ also as $-\pi/2$

$$H_{d}(e^{j\omega}) = -e^{j\pi/2} = -j = 1 \text{ at } \angle -\pi/2, \qquad 0 \le \omega \le \pi$$
$$e^{j\pi/2} = j = 1 \text{ at } \angle \pi/2, \qquad -\pi \le \omega \le 0$$

Note that the magnitude $H(e^{j\omega}) = 1$ for all ω and $\angle H_d(\omega)$ changes from $\pi/2$ to $-\pi/2$ at $\omega = 0$. The filter coefficients are the Fourier coefficients given by

$$h_{d}(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} H_{d}(\omega) e^{j\omega n} d\omega = \frac{1}{2\pi} \left[\int_{-\pi}^{0} j e^{j\omega n} d\omega + \int_{0}^{0} (-j) e^{j\omega n} d\omega \right]$$

$$= \frac{j}{2\pi} \left[e^{j\omega n}_{jn} \right]_{-\pi}^{0} + \frac{(-j)}{2\pi} \left[e^{j\omega n}_{jn} \right]_{0}^{0}$$

$$= \frac{1}{2\pi n} \left((1 - e^{-j\pi n}) - (e^{j\pi n} - 1) \right) = \frac{1}{2\pi n} \left(2 - e_{-j\pi n} - e^{j\pi n} \right)$$

$$= \frac{2 - 2\cos n\pi}{2\pi n} = \frac{1 - \cos n\pi}{n\pi}, \qquad n \neq 0$$

For n = 0, $h_d(0)$ may be evaluated from the defining equation:

$$h_{d}(0) = \frac{1}{2\pi} \int_{-\pi}^{\pi} H_{d}(\omega) e^{j\omega 0} d\omega = \frac{1}{2\pi} \left(\int_{-\pi}^{0} j1 d\omega + \int_{0}^{0} (-j) 1 d\omega \right)$$
$$= \frac{j}{2\pi} \left(\omega \Big|_{-\pi}^{0} \right) - \frac{j}{2\pi} \left(\omega \Big|_{0}^{\pi} \right) = 0$$
$$2\pi \int_{-\pi}^{\pi} 2\pi \int_{0}^{0} e^{j\omega 0} d\omega = \frac{1}{2\pi} \left(\int_{-\pi}^{0} j1 d\omega + \int_{0}^{0} (-j) 1 d\omega \right)$$

Thus

 $h_d(n) = 2/n\pi$ *n* odd 0. *n* even (including 0)

As a specific case choose a filter length N = 9 and evaluate $h_d(n)$ for $-4 \le n \le 4$.

Generate the filter coefficients:

%Hilbert Transform %Generate hdn = 2/(n*pi) for odd n & hdn = 0 for even n, and stem plot n = -50: 50, hdn = 2 ./(n*pi), stem(n, hdn)

n = -50 to 50 (Entered hdn = 0 by hand for even n) Warning: Divide by zero. hdn = [0 -0.0130 0 -0.0135 0 -0.0141 0 -0.0148 0 -0.0155 0 - $0.0163 \quad 0 \quad -0.0172 \quad 0 \quad -0.0182 \quad 0 \quad -0.0193 \quad 0 \quad -0.0205 \quad 0 \quad -0.0220 \quad -0.0220 \quad 0 \quad -0.0220 \quad -0.0220 \quad 0 \quad -0.0220 \quad -0.020 \quad -0.020$ 0.0236 0 -0.0255 0 -0.0277 0 -0.0303 0 -0.0335 0 -0.0374 0 -0.0424 0 -0.0490 0 -0.0579 0 -0.0707 0 -0.0909 0 -0.1273 θ -0.2122 0 -0.6366 Inf 0.6366 0 0.2122 0 0.1273 0 0.0909 0 0.0707 0 0.0579 0 0.0490 0 0.0424 0 0.0374 0 0.0335 0 0.0303 0 0.0277 0 0.0255 0 0.0236 0 0.0220 0 0.0205 0 0.0193 0 0.0182 0 0.0172 0 0.0163 0 0.0155 0 0.0148 0 0.0141 0 0.0135 0 0.0130 0]

We copy and paste the above coefficients and do a stem plot:

%Hilbert Transform

%Stem plot hdn = 2/(n*pi) for odd n & hdn = 0 for even n n = -50: 50,hdn =[0 -0.0130 0 -0.0135 0 -0.0141 0 -0.0148 0 -0.0155 0 -0.0163 0 -0.0172 0 -0.0182 0 -0.0193 0 -0.0205 0 -0.0220 0 -0.0236 0 -0.0255 0 -0.0277 0 -0.0303 0 -0.0335 0 -0.0374 0 -0.0424 0 -0.0490 0 -0.0579 0 -0.0707 0 -0.0909 0 -0.1273 0 -0.2122 0 -0.6366 0 0.6366 0 0.2122 0 0.1273 0 0.0909 0 0.0707 0 0.0579 0 0.0490 0 0.0424 0 0.0374 0 0.0335 0 0.0303 0 0.0205 0.0277 0 0.0255 0 0.0236 0 0.0220 0 0 0.0193 0 0 0.0172 0 0.0163 0 0.0155 0 0.0148 0.0141 0.0182 0 0 0.0135 0 0.0130 0], stem(n, hdn) xlabel('n'), ylabel('hd(n)'); grid; title ('hd(n) = 2/(n*pi) for odd n & hdn = 0 for even n') $hd(n) = 2/(n^*pi)$ for odd n & hdn = 0 for even n 0.8 G 0.6 0.4 0.2 (u)pq 7068-07 0 ტ -0.2 ſŤ -0.4 -0.6 φ -0.8 -50 -40 -30 -20 -10 0 10 20 30 40 50 n

Frequency response

%Hilbert Transform
%Comparison of 9 coefficients vs. 101 coefficients
%9-tap filter coefficients
b9=[0 -0.2122 0 -0.6366 0 0.6366 0 0.2122 0],
a=[1]
b101 =[0 -0.0130 0 -0.0135 0 -0.0141 0 -0.0148 0 -0.0155 0 0.0163 0 -0.0172 0 -0.0182 0 -0.0193 0 -0.0205 0 -0.0220 0 0.0236 0 -0.0255 0 -0.0277 0 -0.0303 0 -0.0335 0 -0.0374 0 0.0424 0 -0.0490 0 -0.0579 0 -0.0707 0 -0.0909 0 -0.1273 0 0.2122 0 -0.6366 0 0.6366 0 0.2122 0 0.1273 0 0.0909 0



Window functions [Ref. S.K. Mitra] We want to see the frequency behavior of window functions by themselves. Note that depending on convenience we shall define the window functions either over $-(N-1)/2 \le n \le (N-1)/2$ or over $0 \le n \le (N-1)$. In the former case the phase function will be zero; in the latter case, used especially in MATLAB, the phase has a negative slope.

There are fixed windows and adjustable windows. Among the fixed windows we have (in addition to the rectangular window) the following tapered windows:

- 1. Bartlett (triangular) window
- 2. Hann (aka Hanning or von Hann) window
- 3. Hamming window
- 4. Blackman window

These windows result in a fixed amount of ripple in the frequency response of the designed filter. The Kaiser window is an adjustable window which allows some control over the ripple.

Bartlett:
$$w(n) = 1 - \frac{2|n|}{N-1}, \qquad -(N-1)/2 \le n \le (N-1)/2$$

Hann: $w(n) = 0.5 + 0.5 \cos |\square|, -(N-1)/2 \le n \le (N-1)/2$
Hamming: $w(n) = 0.54 + 0.46 \cos |\square|, -(N-1)/2 \le n \le (N-1)/2$
Blackman: $w(n) = 0.42 + 0.5 \cos |\square| + 0.08 \cos |\square|, -(N-1)/2 \le n \le (N-1)/2$

The following MATLAB multiplot gives a graphical comparison of the above window functions. Note that all are discrete sequences; to make it easy on the eyes some are plotted as continuous lines and some as discrete.

% Window functions defined over n = -(N-1)/2 to (N-1)/2 N = 31; n = -(N-1)/2: (N-1)/2; wR = n-n+1; % Rectangular wBa = 1 - 2* abs(n)/(N-1); % Bartlett wHn = 0.5 + 0.5 * cos(2*pi*n/(N-1)); % Hamming wHm = 0.54 + 0.46 * cos(2*pi*n/(N-1)); % Hamming wBl = 0.42 + 0.5 * cos(2*pi*n/(N-1)) + 0.08 * cos(4*pi*n/(N-1)); % Blackman % Multiplot. All are discrete sequences. % To make it easy on the eyes some are plotted with continuous lines. plot (n, wR, 'o', n, wBa, 'b', n, wHn, 'k', n, wHm, 'b*', n, wBl, 'k--'); legend ('Rectangular', 'Bartlett', 'Hanning', 'Hamming', 'Blackman'); xlabel('n'), ylabel('w(n)'); grid; title ('Window functions')







The **rectangular window** defined over $-(N-1)/2 \le n \le (N-1)/2$ is given by

$$w_R(n) = 1, \quad -(N-1)/2 \le n \le (N-1)/2$$

0, elsewhere



%Rectangular window defined over n = -(N-1)/2 to (N-1)/2% w(n) = {1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1} N = 11; n = -(N-1)/2; (N-1)/2; wn = n-n+1;stem (n, wn); xlabel('n'), ylabel('w(n)'); grid



The Fourier transform (spectrum) of the window is $W(e^{j\omega}) = \sum_{n=-\infty}^{\infty} w(n) e^{-j\omega n} = \sum_{n=-\infty} 1 e^{-j\omega n}$

We can simplify the above in one of two ways. One possibility is

$$= e^{i\binom{N-1}{j\omega} \left[1 + e^{-j\omega} + 1 + e^{-j\omega} + 1 + e^{-j\omega} + \dots + e^{-j\omega} \right] + \dots + e^{i\binom{N-1}{2}} + \dots + e^{-j\binom{N-1}{2}} + \dots + e^{-j\binom{N-1}{2}} = e^{i\binom{N-1}{2}} \left[\frac{1-e_{-j\omega N}}{2} \right]$$

$$= e^{i\omega} \left[\frac{1+e}{N} + e^{-i\omega N/2} - e^{-j\omega N/2} \right] = e^{i\binom{N-1}{2}} \left[\frac{1-e_{-j\omega N}}{2} \right] = e^{i\binom{N-1}{2}} \left[\frac{1-e_{-j\omega N}}{2} \right]$$

$$= e^{i\omega} \left[\frac{e^{-j\omega N/2}}{N} + \frac{e^{j\omega N/2}}{2} - e^{-j\omega N/2} \right] = (e^{j\omega N/2} - e^{-j\omega N/2})$$

$$= \frac{\left[e^{-j\omega N/2} + e^{j\omega N/2} - e^{-j\omega N/2} \right]}{\sin(\omega/2)}$$

Using L'Hopital's rule,

$$W(e^{j\omega})\Big]_{\omega=0} = \frac{(N/2)\cos(\omega N/2)}{(1/2)\cos(\omega/2)} = N$$
$$20\log_{10} |W(e^{j\omega})]_{\omega=0} = 20\log_{10}N$$

In frequency response plots this is normally regarded as a reference point, that is, as the 0 dB level. That is, the expression for $W(e^{j\omega})$ is *normalized* by dividing it by its value at dc, N in this case:

$$W(e^{j\omega}) = \frac{1}{N} \frac{\sin(\omega N/2)}{\sin(\omega/2)}$$

The second method to simplify $W(e^{j\omega})$ is to combine pairs of terms, one from each end of the expression, into a cosine (or sine, in the case of odd symmetry):

$$W(e^{j\omega}) = \sum_{\substack{n=-\infty\\ [N-1]\\ |\omega| = 1\\ |\omega| = 1}}^{\infty} w(n)e^{-j\omega n} = \sum_{\substack{n=-\\ |\omega| = 1\\ |\omega| = 1\\ |\omega| = 1}}^{-j\omega n} e^{j\omega n} + 1 + e^{j\omega n}$$

This consists of $\begin{pmatrix} 1 & N & -1 \\ 1 & + & -2 \end{pmatrix}$ terms and should be normalized by dividing by N.

Using the equation $W(e^{j\omega}) = \frac{1}{N} \frac{\sin(\omega N/2)}{\sin(\omega/2)}$, its zero-crossings occur when $(\partial V/2)$ equals

integer multiples of π , that is,

 $(\omega N/2) = \pm k\pi$ or $\omega = k(2\pi/N), k \neq 0$

The spectrum between the zero-crossings at $-(2\pi/N)$ and $(2\pi/N)$ is called the *main lobe*, the remaining lobes are called *side lobes*. The width of the main lobe is

Width of the main lobe = $2(2\pi/N) = 4\pi/N$ Width of each side lobe = $2\pi/N$

As the length of the window, N, is increased the lobes become narrower; also the height of the main lobe increases (= N). However, with reference to the *normalized* frequency response

$$W\left(e^{j\omega}\right) = \frac{1}{N} \frac{\sin(\omega N/2)}{\sin(\omega/2)}$$



the height of the main lobe, $W(e^{j\omega})\Big]_{\omega \to 0}$ stays at 1 (or 0 dB) while the side lobes keep getting smaller with increasing *N*.

In the MATLAB segment below the window is defined over $0 \le n \le (N-1)$ rather than over $-(N-1)/2 \le n \le (N-1)/2$.

%Frequency response of rectangular window defined over n = 0 to N-1 % h(n) = {1, 1, 1, 1, 1, 1, 1, 1, 1, 1} b11= [1, 1, 1, 1, 1, 1, 1, 1, 1, 1], a= [1] w=-pi: pi/256: pi; Hw11=freqz(b11, a, w); subplot(2, 1, 1), plot(w, abs(Hw11)); legend ('Magnitude (Length = 11)'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid subplot(2, 1, 2), plot(w, angle(Hw11)); legend ('Phase'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid

Note in the plot below that the height of the main lobe is 11 (= N). The width of the main lobe is taken as the separation between the zero crossings on either side of $\omega = 0$:

Width of main lobe = $4\pi/N = 4\pi/11$

From the plot, by eyeballing, we can gather:

- 1. For a given window length the side lobes have the same width.
- 2. For a given window length the magnitude of the side lobes decreases with increasing frequency



Normalized frequency *r* In the MATLAB segment below the window is defined over $0 \le n \le (N-1)$. Further we define the normalized frequency $r = \omega/\pi$. As ω varies from $-\pi$ to π the normalized frequency varies from -1 to 1.

% Magnitude response of rectangular window defined over n = 0 to N–1 % h(n) = {1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1} b11=[1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1], a=[1] w=-pi: pi/256: pi; r = w/pi; Hw11=freqz(b11, a, w); subplot(2, 1, 1), plot(w, abs(Hw11)); legend ('Length = 11'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid % %Normalized frequency $r = \omega/\pi$ Hr11=freqz(b11, a, pi*r); subplot(2, 1, 2), plot(r, abs(Hr11)); legend ('Length = 11'); xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid

Note in the plot below that as ω goes from $-\pi$ to π the normalized frequency *r* goes from -1 to 1.



Normalized magnitude In the MATLAB segment below we go another step: we normalize the magnitude by dividing it by *N*. The window is defined over $0 \le n \le (N-1)$ and the normalized frequency is $r = \omega/\pi$. As ω varies from $-\pi$ to π the normalized frequency varies from -1 to 1.

% Magnitude response rectangular window defined over n = 0 to N-1 % $h(n) = \{1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1\}$ N = 11; b11 = ones(1, 11); a = [1]; w=-pi: pi/256: pi; r = w/pi; % %Normalized magnitude $H(\omega)/N$ Hw11n= freqz(b11, a, w)/N; subplot(2, 1, 1), plot(w, abs(Hw11n)); legend ('Length = 11'); . xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)');grid; title ('Normalized magnitude') % %Normalized magnitude and normalized frequency Hr11n = freqz(b11, a, pi*r)/N;subplot(2, 1, 2), plot(r, abs(Hr11n)); legend ('Length = 11'); xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid; title ('Normalized magnitude')

Note in the plot below that the height of the main lobe is 1 since it is normalized.



Comparison of two rectangular windows of lengths 11 and 31 In the MATLAB segment below we go another step: we normalize the magnitude by dividing it by *N*. The window is defined over $0 \le n \le (N-1)$ and the normalized frequency is $r = \omega/\pi$. As ω varies from $-\pi$ to π the normalized frequency varies from -1 to 1.

%Magnitude response rectangular window defined over n = 0 to N-1 % Comparison of two rectangular windows of lengths 11 and 31 N1 = 11; b11= ones(1, N1); N2 = 31; b31= ones(1, N2); a=[1]; w=-pi: pi/512: pi; r = w/pi; % %Length 11, normalized magnitude and normalized frequency Hr11n= freqz(b11, a, pi*r)/N1; subplot(2, 1, 1), plot(r, abs(Hr11n)); legend ('Length = 11'); xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid % %Length 31, normalized magnitude and normalized frequency Hr31n= freqz(b31, a, pi*r)/N2; subplot(2, 1, 2), plot(r, abs(Hr31n)); legend ('Length = 31'); xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid

Only the *first* side lobe is of concern since all the other side lobes are smaller. From the plot below the maximum of the first side lobe is a little over 20% of the height of the main lobe in both windows.



Two rectangular windows of lengths 11 and 31 compared on a multi-plot In the MATLAB segment below we compare the two window lengths on the same multi-plot. As before, the window is defined over $0 \le n \le (N-1)$ and the normalized frequency is $r = \omega/\pi$. As ω varies from $-\pi$ to π the normalized frequency varies from -1 to 1.

%Magnitude response rectangular window defined over n = 0 to N–1 % Comparison of two rectangular windows of lengths 11 and 31 N1 = 11; b11= ones(1, N1); N2 = 31; b31= ones(1, N2); a=[1]; w=-pi: pi/768: pi; r = w/pi; % %Length 11, normalized magnitude and normalized frequency Hr11n= freqz(b11, a, pi*r)/N1; % %Length 31, normalized magnitude and normalized frequency Hr31n= freqz(b31, a, pi*r)/N2; % plot(r, abs(Hr11n), r, abs(Hr31n), 'k'); legend ('Length = 11', 'Length = 31'); title('Comparison window lengths 11 vs. 31'); xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid

From the plot, by eyeballing, we can gather that for a given window length the magnitude of the side lobes decreases with increasing frequency. Further, as the window length increases the height of a specific side lobe (such as the first side lobe) decreases: for instance the height of the first side lobe for N = 31 is smaller than that of the first side lobe for N = 11.


The **Hamming window** defined over $-(N-1)/2 \le n \le (N-1)/2$ is

 $w_{Ham}(n) = 0.54 + 0.46 \cos[2\pi n / (N-1)], -(N-1)/2 \le n \le (N-1)/2$ 0, elsewhere

% Hamming window defined over n = -(N-1)/2 to (N-1)/2% w(n) = 0.54 + 0.46 cos(2*pi*n/(N-1)) N = 31; n = -(N-1)/2: (N-1)/2; wn = 0.54 + 0.46 * cos(2*pi*n/(N-1)), stem (n, wn); xlabel('n'), ylabel('w(n)'); grid; title ('Hamming window')

n = -5 -4 -3 -2 -1 2 0 1 3 4 5 $w(n) = \{0.0800$ 0.0901 0.1198 0.1679 0.2322 0.3100 0.3979 0.4919 0.5881 0.6821 0.7700 0.8478 0.9121 0.9602 0.9899 1.0000 0.9899 0.9602 0.9121 0.8478 0.7700 0.6821 0.5881 0.4919 0.3979 0.3100 0.2322 0.1679 0.1198 0.0901 0.0800



The Fourier transform (spectrum) of the window is $\binom{N-1}{2}$

$$W(e^{j\omega}) = \sum_{\substack{n = -\infty \\ (N-1)/2}}^{\infty} w(n)e^{-j\omega n} = \sum_{\substack{n = -(N-1)/2 \\ (N-1)/2}}^{(N-1)/2} \{0.54 + 0.46\cos 2\pi n/(N-1)\}e^{-j\omega n}$$

= $\sum_{\substack{n = -(N-1)/2 \\ n=-(N-1)/2}}^{(N-1)/2} 0.46\{\cos 2\pi n/(N-1)\}e^{-j\omega n}$
= $0.54 \frac{\sin(\omega N/2)}{\sin(\omega/2)} + \sum_{\substack{n = -(N-1)/2 \\ (N-1)/2}}^{(N-1)/2} 0.46\left\{\frac{j2\pi n/(N-1)}{2} + \frac{j2\pi n/(N-1)}{2}\right\}e^{-j\omega n}$
= ...

$$= 0.54 \frac{\sin(\omega N/2)}{\sin(\omega/2)} + 0.23 \frac{\sin[(\omega N/2) - \pi N/(N-1)]}{\sin[(\omega/2) - \pi/(N-1)]} + 0.23 \frac{\sin[(\omega N/2) - \pi N/(N-1)]}{\sin[(\omega/2) + \pi N/(N-1)]}$$

The magnitude at dc is

$$W(e^{j\omega})\Big]_{\omega=0} = 0.54 N + 2 (0.23) \frac{\sin[\pi N / (N-1)]}{\sin[\pi / (N-1)]}$$

This is calculated below for N = 11, 21, 31 and 41:

%Magnitude at DC N = 11:10:41, WdcN = 0.54*N + 0.46* sin(pi*N./(N-1))./sin(pi./(N-1))

The width of the main lobe is taken as the separation between the zero crossings on either side of $\omega = 0$. This is obtained by setting $W \left(e^{j\omega} \right) = 0$ and solving for ω ; it is given as

Width of main lobe (Hamming) = $8\pi/N$

twice that of the rectangular window.

The Hamming window, defined over the interval $0 \le n \le N-1$, is given by

 $w_{Ham}(n) = 0.54 - 0.46 \cos[2\pi n / (N - 1)], \quad 0 \le n \le N - 1$ 0, elsewhere

For N = 11, the Hamming window is given by $w_{Ham}(n) = 0.54 - 0.46 \cos(\pi n / 5)$, $0 \le n \le 10$.

Frequency response of Hamming window:

% Magnitude response of 11-point Hamming window defined over n = 0 to N-1 % WHam = b11 = 0.54 - 0.46 *cos((2*pi/(N-1) .*n)) N = 11; n = 0: N-1; b11 = 0.54 - 0.46 *cos((2*pi/(N-1) .*n)); a=[1] w=-pi: pi/256: pi; Hw11=freqz(b11, a, w); subplot(2, 1, 1), plot(w, abs(Hw11)); legend ('Magnitude (Length = 11)'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of H(\omega)'); grid title('Hamming window, Frequency response'); subplot(2, 1, 2), plot(w, angle(Hw11)); legend ('Phase'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of H(\omega)'); grid

In the phase plot below the phase reaches $-\pi$ in the interval $0 \le \omega \le 1$ and is therefore adjusted by adding 2π ; this is not due to a zero-crossing. The same applies to the phase adjustment in the interval $-1 \le \omega \le 0$.



%Magnitude response of Hamming window defined over n = 0 to N-1 %Length 11 N = 11; n = 0: N-1; b11 = 0.54 - 0.46 *cos((2*pi/(N-1) .*n)); a=[1]; w=-pi: pi/768: pi; r = w/pi; % normalized frequency Hr11n= freqz(b11, a, pi*r); % Frequency response % Plot plot(r, abs(Hr11n)); legend ('Length = 11'); title('Hamming window, Frequency response');

xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid





% Magnitude response of Hamming window defined over n = 0 to N-1 % Length 31 N = 31; n = 0: N-1; b31 = 0.54 - 0.46 *cos((2*pi/(N-1) .*n)); a=[1]; w=-pi: pi/768: pi; r = w/pi; % normalized frequency Hr31n= freqz(b31, a, pi*r); % Frequency response % Plot plot(r, abs(Hr31n)); legend ('Length = 31'); title('Hamming window, Frequency response');

xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid



% Magnitude response of Hamming window defined over n = 0 to N-1 %Length 31 N = 31; n = 0: N-1; $b31 = 0.54 - 0.46 \cos((2 \frac{pi}{N-1}).*n)); a=[1];$ w=-pi: pi/2048: pi; r = w/pi; % normalized frequency %Magnitude at dc $WdcN = 0.54*N + 0.46* \sin(pi*N/(N-1))/\sin(pi/(N-1)),$ Hr31n= freqz(b31, a, pi*r)/WdcN; % normalized frequency response % Plot plot(r, abs(Hr31n)); legend ('Length = 31');title('Hamming window, Normalized frequency response'); xlabel('Normalized Frequency r'), ylabel('Magnitude of H(r)'); grid

Observations:

- In the plot below the magnitude of the side lobes appears to stay about the same 1. with increasing frequency for the first few lobes, but beyond r = 0.6 there is a discernible (to the eye) fall-off in the height of the side lobes. Relative to the rectangular window, however, we may say the side lobe height stays about the same with increasing frequency.
- The width of the side lobe appears to stay constant with increasing frequency. 2.



Hanning window

% Hanning window defined over n = -(N-1)/2 to (N-1)/2% w(n) = 0.5 + 0.5 * cos(2*pi*n/(N-1)) N = 31; n = -(N-1)/2: (N-1)/2; wn = 0.5 + 0.5 * cos(2*pi*n/(N-1)), stem (n, wn); xlabel('n'), ylabel('w(n)'); grid; title ('Hanning window')



Blackman window

%Blackman window defined over n = -(N-1)/2 to (N-1)/2% w(n) = 0.5 + 0.5 * cos(2*pi*n/(N-1)) N = 31; n = -(N-1)/2: (N-1)/2; wn = 0.42 + 0.5 * cos(2*pi*n/(N-1)) + 0.08 * cos(4*pi*n/(N-1)), stem (n, wn); xlabel('n'), ylabel('w(n)'); grid; title ('Blackman window')



The width of the main lobe is taken as the separation between the zero crossings on either side of $\omega = 0$. This is obtained by setting $W \left(e^{j\omega} \right) = 0$ and solving for ω ; it is given as

Width of main lobe (Blackman) = $12\pi/N$

thrice that of the rectangular window.

Bartlett window

% Bartlett window defined over n = -(N-1)/2 to (N-1)/2% w(n) = 1 - 2* abs(n)/(N-1) N = 31; n = -(N-1)/2: (N-1)/2; wn = 1 - 2* abs(n)/(N-1), stem (n, wn); xlabel('n'), ylabel('w(n)'); grid; title ('Bartlett window')



4.5 Choosing between FIR and IIR filters

IIR Filter	FIR Filter
Use IIR when the only important	Use FIR if
requirements are	1. the number of coefficients is not a
1. sharp cutoff frequency and	problem and
2. high throughput as IIR filters,	2. in particular, if very little or no
especially those using elliptic	phase distortion is desired
characteristics, will have fewer	Some DSP processors have architectures
coefficients than FIR	that are tailored/designed for FIR filters.

The choice between FIR and IIR filters depends largely on the relative advantages of the two types of filters.

(1) Linear phase FIR filters can have exactly linear phase response. This means no phase distortion. This is an important requirement in, for example, data transmission, biomedicine, digital audio, and image processing.

The phase response of IIR filters is nonlinear, especially at the band edges.

(2) **Stability** FIR filters realized non-recursively are always stable. The stability of IIR filters cannot always be guaranteed.

(3) The effects of finite word length such as round off noise and coefficient quantization errors are much less severe in FIR than in IIR.

(4) **FIR requires more coefficients** for sharp cut-off filters than IIR. Thus for a given amplitude response specification, more processing time and storage will be required for FIR implementation. However, one can readily take advantage of the computational speed of the FFT and multirate techniques to improve significantly the efficiency of FIR implementation.

(5) Analog filters can be readily transformed into equivalent IIR digital filters meeting similar specifications. This is not possible with FIR filters as they have no analog counterpart. However, with FIR it is easier to synthesize filters of arbitrary frequency response.

(6) In general, FIR is algebraically more difficult to synthesize, if CAD support is not available.

Digital Signal Processing – 5

5 Multirate DSP

Decimation, Interpolation, Sampling rate conversion, Filter design and-Implementation of sampling rate conversion.

Contents:

Time and frequency scaling in continuous-time systems
Transformation of the independent variable
Down-sampling
Up-sampling
Cascading sampling rate converters
Identities
FIR implementation of sampling rate conversion
Polyphase structures
Polyphase structure for a decimator

Discrete-time systems with different sampling rates at various parts of the system are called **multirate systems**. They are *linear and time-varying systems*. Integer sampling rate converters change the sampling frequency by an integer factor and rational sampling rate converters by a rational number. Here is a sampling of sampling rates in commercial applications (Mitra):

Sampling Rates	
Digital Audio	Video
Broadcasting – 32 kHz	Composite Video Signal
CD – 44.1 kHz	• NTSC – 14.3181818 MHz
DAT – 48 kHz	• PAL – 17.734475 MHz
	Digital Component Video Signal
	• Luminance – 13.5 MHz
	• Color difference – 6.75 MHz

Time and frequency scaling in continuous-time systems

Illustration An audio signal recorded on cassette tape at a certain speed could be played back at a higher speed than that at which it was recorded. This is called *time scaling*, in particular, *compression* in the time domain, and results in an inverse effect in the frequency domain, i.e., an expansion of the frequency spectrum. Similarly when the audio signal is played back at a slower speed than the recording speed we have *expansion* in the time domain resulting in a corresponding compression of the spectrum in the frequency domain.

Given the signal x(t) and its Fourier transform $X(\Omega)$, represented notationally by

$$x(t) - X(\Omega)$$

then time scaling results in

$$x(at) = -\frac{1}{a}X(\Omega/a)$$

If a > 1 the scaling corresponds to **compression in time**. If, for instance, a = 2, we may visualize a new signal $y_I(t) = x(2t)$; with t = 1, for instance, the value of x, that is, x(2) that occurred at 2 seconds occurs at 1 second in the case of y_I , that is, $y_I(1)$ – which is compression in time.



If x(t) is an audio signal *recorded* on tape then x(2t) could be the signal x(t) played back at twice the speed at which x(t) was recorded. The signal x(2t) varies more rapidly than x(t) and the playback frequencies are higher.

If a < 1 the scaling corresponds to **expansion in time**. If, for instance, a = 1/2, then x(t/2) is the signal x(t) played back at half the speed at which x(t) was recorded. The signal x(t/2) varies slower than x(t) and the playback frequencies are lower. Again, we may visualize this as a new signal $y_2(t) = x(t/2)$; the value of x(.) that occurred at t/2 occurs at t in the case of $y_2(.)$ – which is expansion in time.

Time expansion and frequency compression is found in data transmission from space probes to receiving stations on earth. To reduce the amount of noise superposed on the signal, it is necessary to keep the bandwidth of the receiver as small as possible. One means of doing this is to reduce the bandwidth of the signal: store the data collected by the probe, and then transmit it at a slower rate. Since the time-scaling factor is known, the signal can be reproduced at the receiver.

The corresponding operations in the case of discrete-time systems are not quite so straight forward owing to

- 1. The need to band limit the continuous-time signal prior to sampling, and
- 2. The need to avoid aliasing in the process of sampling

Example 5.1 Consider the 4 Hz signal $x(t) = \cos 2\pi 4t$ which is obviously band-limited to $F_{max} = 4$ Hz. It is sufficient to sample it at 8 Hz. Alternatively, the signal can be sampled at, say, 16 Hz or 20 Hz etc. Suppose that it has been *over-sampled* by a factor of, say, 6 at $F_s = 48$ Hz to give $x(n) = \cos 2\pi 4n(1/48) = \cos (\pi n/6)$.

- (a) If it is desired subsequently to generate from x(n) another signal $x_I(n)$ that is a discrete-time version of x(t) sampled at $F_{sI} = 16$ Hz (sampling rate reduced by a factor of 3), then can we do this by simply dropping two samples of x(n) for every sample that we keep? That is $x_I(n) = x(3n)$. This is called **down-sampling**.
- (b) How do we generate from x(n) another signal $x_2(n)$ that is a discrete-time version of x(t) sampled at, say, $F_{s2} = 96$ Hz (sampling rate doubled)? This is called **up-sampling**.
- (c) Can we generate from x(n) another signal $x_3(n)$ that is a discrete-time version of x(t) sampled at $F_{s3} = 6$ Hz?

We pick up on this problem again after covering transformation of the independent variable.

Transformation of the independent variable

Time scaling (Refer also to Section 7.5 of Signals and Systems, Oppenheim and Willsky.) Given the sequence x(n), the sequence y(n) = x(2n) is obtained by skipping odd-numbered samples in x(n) and retaining the even-numbered ones. The extension to y(n) = x(Mn) means we retain sample numbers 0, M, 2M, 3M, ..., and skip the intervening M-1 samples between those we keep. The original sequence x(n) is obtained by sampling a continuous signal x(t) at a certain rate (perhaps over-sampling). The signal y(n) = x(Mn) is then obtained by reducing the sampling rate by a factor of M on the continuous-time signal x(t). This is known as **down-sampling** or **decimation** or sampling rate compression.

Similarly the process of constructing the sequence y(n) = x(n/L) from the sequence x(n) means we derive y(n) by inserting (*L*-1) sequence points with zero value between points of x(n). This is called **up-sampling** or **interpolation** or sampling rate expansion. (Inserting (*L*-1) zeros is

just one way of interpolating. It is also possible for the up-sampler to be followed by a digital system that replaces the inserted zeros with more appropriate values based on a linear combination of the x(n) samples.)

In general, the result of time scaling a discrete-time signal is not just a stretched or compressed version of the original but possibly a totally different sequence/waveform.

Example 5.2.1 Given that $x(t) = e^{-5t}u(t)$ is sampled at 50 Hz, find an expression for x(n). Plot x(t), x(n) and x(2n). Sketch the spectrum of x(n).

Solution The sampling time is T = 0.02 sec. Replacing t with nT we get $x(nT) = e^{-5nT}u(nT)$, or

$$x(n) = (e^{-0.1})^n u(n) = (0.905)^n u(n)$$

We show below three plots: (1) The continuous-time signal x(t), (2) The sampled (at 50 Hz) version x(n), and (3) x(2n), the 2-fold down-sampled version of x(n); this is equivalent to sampling x(t) at 25 Hz.

t = 0 : 1/512: 1; xt = exp (-5*t); %x(t) evaluated at 512 points subplot(3, 1, 1), plot(t, xt); legend ('x(t) = exp(-5t)'); xlabel ('time, sec.'), ylabel('x(t)'); grid; title ('x(t) – Continuous-time') % t1 = 0 : 0.02: 1; xn = exp (-5*t1); %Sampled at 50 Hz. subplot(3, 1, 2), stem(t1, xn); legend ('x(n) at 50 Hz'); xlabel ('time, sec.'), ylabel('x(n)'); grid; title ('x(nT) at T = 0.02 sec') % t2 = 0 : 0.04: 1; xt2 = exp (-5*t2); %Sampled at 25 Hz subplot(3, 1, 3), stem(t2, xt2); legend ('2-fold down-sampled');

xlabel ('time, sec.'), ylabel('x(2n)'); grid; title ('x(nT) at T = 0.04 sec.')



Note that $X(s) = \mathcal{L}(e^{-5t} u(t)) \neq 1$ (s + 5). Shown below is the MATLAB plot of the magnitude spectrum $|X(j\Omega)|$ of the continuous-time signal x(t) using the function **plot**. Omega is a vector, consequently we use "./" instead of "/" etc. The main point to be made here is that $X(j\Omega)$ extends asymptotically to ∞ , so, strictly speaking, x(t) is not band-limited. Consequently, the spectrum $X(\omega)$ of the sampled signal x(n) (shown later below) has some built-in aliasing.

t = 0 : 1/512: 1; xt = exp (-5*t); %x(t) evaluated at 512 points subplot(3, 1, 1), plot(t, xt); legend ('x(t) = exp(-5t)'); xlabel ('time'), ylabel('x(t)'); grid; title ('x(t) – Continuous-time') %

Omega = -6*pi: pi/256: 6*pi; X = 1./(5.+j.*Omega);

subplot(3, 1, 2), plot(Omega, abs(X), 'k'); legend ('Spectrum of x(t)'); xlabel ('Omega, rad/sec'), ylabel('|X(Omega)|'); grid; title ('Magnitude')% subplot(3, 1, 3), plot(Omega, angle(X), 'k'); legend ('Spectrum of x(t)'); xlabel ('Omega, rad/sec'), ylabel('Phase of X(Omega)'); grid; title ('Phase')



Coming to the discrete-time signal, the spectrum of $x(n) = a^n u(n) = (0.905)^n u(n)$ is its DTFT

$$X(e^{j\omega}) = \sum_{n=0}^{\infty} a^n e^{-j\omega n} = \frac{1}{1 - ae^{-j\omega}} = \frac{1}{1 - 0.905e^{-j\omega}}$$

The MATLAB segment is

t1 = 0: 0.02: 1; xn = exp (-5*t1); % Sampled at 50 Hz.subplot(3, 1, 1), stem(t1, xn); legend ('x(n) at 50 Hz'); xlabel ('time, sec.'), ylabel('x(n)'); grid; title ('x(nT) at T = 0.02 sec') % b = [1]; %Numerator coefficient a = [1, -0.905]; %Denominator coefficients w = -6*pi: pi/256: 6*pi; [Xw] = freqz(b, a, w); subplot(3, 1, 2), plot(w, abs(Xw)); legend ('Spectrum of x(n)'); xlabel('Frequency \omega, rad/sample'), ylabel('Magnitude of X(\omega)'); grid subplot(3, 1, 3), plot(w, angle(Xw)); legend ('Spectrum of x(n)'); xlabel('Frequency \omega, rad/sample'), ylabel('Phase of X(\omega)'); grid



Example 5.2.2 Given $x(n) = e^{-n/2}u(n)$, find (a) x(5n/3), (b) x(2n), (c) x(n/2). **Answer** The sequence

$$x(n) = e^{-n/2}u(n) = \left(e^{-1/2}\right)^n u(n) = (0.606)^n u(n) = a^n u(n)$$

where $a = e^{-1/2} = 0.606$, is sketched below:

$$x(n) = e^{-n/2}u(n) = a^{n}u(n)$$

$$1 + a = e^{-1/2} = 0.606$$

$$a + a^{2} + a^{3} + a^{4} + a^{$$

(a) With y(n) = x(5n/3), we evaluate y(n) for several values of *n* (we have assumed here that x(n) is zero if *n* is not an integer):

$$y(0) = x(5 \cdot 0/3) = x(0) = e^{-0/2} = 1$$

$$y(1) = x(5 \cdot 1/3) = x(5/3) = 0$$

$$y(2) = x(5 \cdot 2/3) = x(10/3) = 0$$

$$y(3) = x(5 \cdot 3/3) = x(5) = e^{-5/2} = a^{5}$$

...

$$y(6) = x(5 \cdot 6/3) = x(10) = e^{-10/2} = a^{10}$$

The general expression for y(n) can be written as

 $y(n) = x(5n/3) = e^{-(5n/3)/2}$, *n* as specified below

- $= e^{-5n/6}, \qquad n = 0, 3, 6, ...$ $0, \qquad \text{otherwise}$
 - $n = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9 \quad 10$ $y(n) = 1 \quad 0 \quad 0 \quad a^{5} \quad 0 \quad 0 \quad a^{10} \quad 0 \quad 0 \quad a^{15} \quad 0$

The sequence is sketched below:



(b) With y(n) = x(2n), we evaluate y(.) for several values of n:

 $\begin{array}{l} y(0) = x(2 \ . \ 0) = x(0) = 1 \\ y(1) = x(2 \ . \ 1) = x(2) = a^2 \\ y(2) = x(2 \ . \ 2) = x(4) = a^4 \\ y(3) = x(2 \ . \ 3) = x(6) = a^6 \\ \end{array}$

The general expression for y(n) can be written as

 $y(n) = x(2n) = e^{-(2n)/2}$, *n* as specified below

$$= e^{-n}, \quad n \ge 0$$

0, otherwise

The sequence y(.) is made up of every other sample of x(.). This is **down-sampling** or **decimation** by a factor of 2 (or, compression in time). Note that some of the original sample values have disappeared. The sequence is sketched below.



(c) With $y(n) = x(n/2) = e^{-n/4}u(n)$, we evaluate y(.) for several values of *n* (again, we have assumed here that x(n) is zero if *n* is not an integer):

y(0) = x(0/2) = x(0) = 1 y(1) = x(1/2) = x(0.5) = 0 y(2) = x(2/2) = x(1) = ay(3) = x(3/2) = x(1.5) = 0

The general expression for y(n) can be written as $y(n) = x(n/2) = e^{-(n/2)/2}$, *n* as specified below

=
$$e^{-n/4}$$
, $n = 0, 2, 4, ...$
0, otherwise

$$n = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6$$

$$y(n) = 1 \quad 0 \quad a \quad 0 \quad a^2 \quad 0 \quad a^3$$

The sequence y(.) is constructed by inserting one zero between successive samples of x(.). This is **up-sampling** or **interpolation** by a factor of 2 (or expansion in time). The sequence is sketched below:



To get back to the problem raised earlier, given the sequence x(n) obtained from x(t) at a rate (1/T)

$$x(t) < x(nT) < x(n)$$
, rate (1/*T*)

we want to obtain the sequence x'(n) which corresponds to a sampling rate $(1/T^+)$ where $T \neq T^+$

$$x(t) < x(n T') < x'(n)$$
, rate $(1/T')$

There are two approaches to do this:

- Convert x(n) to x(t) and resample at (1/T') to generate x'(n). This is not ideal because of the imperfections in the A/D-H(z)-D/A originally involved in generating x(n). Or,
- 1 Change the sampling rate entirely with discrete-time operations.

Example 5.2.3 Consider the 4 Hz signal $x(t) = \cos 2\pi 4t$ which is obviously band-limited to $F_{max} = 4$ Hz. It is sufficient to sample it at 8 Hz. Suppose that it has been *over-sampled* by, say, a factor of 6 at $F_s = 48$ Hz to give $x(n) = \cos 2\pi 4n(1/48) = \cos (\pi n/6)$.

If it is desired subsequently to generate from x(n) another signal $x_I(n)$ that is a discretetime version of x(t) sampled at $F_{sI} = 16$ Hz (sampling rate reduced or down-sampled by a factor of 3), then can we do this by dropping two samples of x(n) for every sample that we keep? In this specific example this is possible since a sampling rate of 16 Hz is clearly greater than $2F_{max}$ of 8 Hz. Thus the down-sampled version is obtained by replacing n in x(n) by 3n

 $x_{l}(n) = x(3n) = \cos(\pi 3n/6) = \cos(\pi n/2)$ < (1)

Let us compare this with what we would get if we were to sample $x(t) = \cos 2\pi 4t$ directly at 16 Hz. We simply replace t by nT = n(1/16)

 $x_l(n) = \cos 2\pi 4n(1/16) = \cos(\pi n/2)$ < (2)

The results in (1) and (2) are the same. (QED)

We show below three plots: (1) The continuous-time signal x(t), (2) The sampled (at 48 Hz) version x(n), and (3) x(3n), the 3-fold down-sampled version of x(n); this is equivalent to sampling x(t) at 16 Hz.

t = 0 : 1/128: 0.5; xt = cos (2*pi*4*t); %x(t) evaluated at 128 points subplot(3, 1, 1), plot(t, xt); legend ('4-Hz Cosine'); xlabel ('time, sec.'), ylabel('x(t)'); grid; title ('x(t) – Continuous-time') % t1 = 0 : 1/48: 0.5; xn = cos (2*pi*4*t1); %Sampled at 48 Hz subplot(3, 1, 2), stem(t1, xn); legend ('x(n) at 48 Hz'); xlabel ('time, sec.'), ylabel('x(n)'); grid; title ('x(nT) at T = 1/48 = 0.020 sec.') % t3 = 0 : 1/16: 0.5; x1n = cos (2*pi*4*t3); %Sampled at 16 Hz subplot(3, 1, 3), stem(t3, x1n); xlabel ('time, sec.'), ylabel('x(3n)'); grid; title ('3-fold down-sampled, x(nT) at T = 1/16 = 0.0625 sec.')



Alternatively, assuming x(t) is not available, $x_1(n)$ could be obtained as follows:

- 1. Recover x(t) by passing x(n) through a DAC
- 2. Sample the resulting x(t) at $F_{sl} = 16$ Hz

We take it, however, that this option is not desirable.

The above analysis assumes that we know the frequency content of the base band signal, x(t). Generally this is not the case. Given the sequence x(n) that was obtained by sampling at a rate, say F_s , we do not know what is the maximum frequency, F_{max} , contained in the underlying analog signal, x(t). Assuming it was originally band-limited and properly sampled, it is safest to assume that the base band signal was band-limited to $F_s/2$ (= F_{max}) and not lower. In such a case simply dropping one or more samples of x(n) for every sample we keep will not work. If we want to reduce the sampling rate by a factor of, say, K, then we would have to band-limit the precursor of x(n) to $(F_s/2)/K = (F_s/2K)$ and then sample it at the K-fold reduced sampling rate to achieve the required decimation. This amounts to down sampling x(n) by a factor of K. (If the signal x(t) originally actually contained a maximum frequency of $F_s/2$ then subsequent down sampling will result in unavoidable loss of information. But if it was band limited to significantly less than $F_s/2$ then down-sampling without loss of information is possible.)

The band-limiting mentioned above may be done either in the continuous-time domain or in the discrete-time domain. The procedure in the continuous time domain is as follows: Imagine that x(t) is recovered from x(n); x(t) is then band-limited to $F_s/2K$ by passing it through an ideal low pass filter described by

$$H(F) = 1, \quad 0 \le F < F_s/2K$$
$$0, \quad F_s/2K \le F \le F_s/2$$



The band-limited signal, denoted $x_I(t)$ is then sampled at the reduced rate of F_s/K to generate $x_I(n)$. This method is generally undesirable because of all the imperfections inherent in originally generating x(n) from x(t) at a sampling rate of F_s , converting x(n) back into x(t), then band-limiting x(t) to $F_s/2K$ to generate $x_I(t)$ and then sampling $x_I(t)$ at a sampling rate of F_s/K to generate $x_I(n)$.

Sampling rate decimation Reducing the sampling rate by an integer factor in the discrete-time domain is shown in the following block diagram. The down arrow in K indicates **down sampling** by a factor of *K*. The filter H(z) is a *digital anti-aliasing filter* whose output v(n) is a low pass filtered version of x(n).



If the filter H(z) is implemented as a linear phase FIR filter with (M+1) coefficients specified as $\{b_r, r = 0 \text{ to } M\}$, (some call it " M^{th} order"), then

$$v(n) = \sum_{r=0}^{m} b_r x(n-r)$$

We desire the output y(n) to be a down-sampled version of x(n), that is

$$y(n) = v(Kn) = \sum_{r=0}^{m} b_r x(Kn-r)$$

Example 5.2.4 Consider the 4 Hz signal $x(t) = \cos 2\pi 4t$ which is obviously band-limited to $F_{max} = 4$ Hz. It is sufficient to sample it at 8 Hz. Suppose instead that it has been over sampled, say, by a factor of 6 at $F_s = 48$ Hz to give $x(n) = \cos 2\pi 4n(1/48) = \cos (\pi n/6)$.

Can we generate from x(n) another signal $x_3(n)$ that is a discrete-time version of x(t) sampled at $F_{s3} = F_s/8 = 6$ Hz? This is down sampling by a factor of 8. We simply replace t by nT = n(1/6) to get

$$x_3(n) = \cos 2\pi 4n(1/6) = \cos (8\pi n/6) = x(8n) = \{x(0), x(8), x(16), \dots\}$$

In other words, $x_3(n)$ is made up of every 8th sample of x(n). For every sample value of x(n) we keep we discard the next 7 samples. We know, however, that a sampling frequency of 6 Hz does not satisfy the sampling theorem; in this case down sampling has been taken too far.

We show below three plots: (1) The sampled (at 48 Hz) version x(n) – this is repeated from above, (2) x(2n), the 2-fold down-sampled version of x(n); this is equivalent to sampling x(t) at 24 Hz, and (3) x(8n), the 8-fold down-sampled version of x(n); this is equivalent to sampling x(t) at the unacceptably low rate of 6 Hz.

t1 = 0 : 1/48: 0.5; xn = cos (2*pi*4*t1); % Sampled at 48 Hzsubplot(3, 1, 1), stem(t1, xn); legend ('x(n) at 48 Hz');xlabel ('time, sec.'), ylabel('x(n)'); grid; title ('x(nT) at T = 1/48')%<math>t2 = 0 : 1/24: 0.5; xt2 = cos (2*pi*4*t2); % Sampled at 24 Hz subplot(3, 1, 2), stem(t2, xt2); xlabel ('time, sec.'), ylabel('x(2n)'); grid; title ('2-fold down-sampled, x(nT) at T = 1/24 sec.') %

t4 = 0: 1/6: 0.5; x3n = cos (2*pi*4*t4); %Sampled at 6 Hzsubplot(3, 1, 3), stem(t4, x3n); xlabel ('time, sec.'), ylabel('x(8n)');grid; title ('8-fold down-sampled, x(nT) at T = 1/6 sec.')



Example 7.2.5 To show visually a case of down sampling that is not satisfactory, consider $x_4(n)$ generated from x(n) by down sampling by a factor of 12, i.e., $x_4(n) = x(12n)$. This is also obtained by sampling at 48/12 = 4 Hz:

$$x_4(n) = x(nT) = \cos 2\pi 4n(1/4) = \cos (2\pi n) = \cos (12\pi n/6) = x(12n)$$

In this case $\cos(2\pi n) = 1$ for all *n*, so that

$$x_4(n) = 1$$
 for all n

which has no resemblance to x(n), making it visually obvious that down sampling has been taken too far. Depending on at what point in the cycle the samples are taken, $x_4(n)$ equals a constant (including 0), for all n.

Down-sampling

Assume that x(n) is obtained from an underlying continuous-time signal x(t) by sampling at F_x Hz. Assume that x(t) was originally band limited to $F_x/2$ Hz. On the digital frequency (ω) scale this amounts to x(n) being band limited to π .

We now wish to generate a signal y(n) by down-sampling x(n) by a factor of M, that is, we are reducing the sampling rate by a factor of M. This amounts to:

- 1. Converting x(n) to x(t) using a D/A converter.
- 2. Band limiting x(t) to $F_x/2M$ Hz. Assume that no information is lost due to this band limiting.
- 3. Resampling x(t) at F_x/M Hz. to produce y(n).

Equivalently the above task is accomplished entirely in the digital domain by

- 1. Band limiting x(n) to π/M . Assume that no information is lost due to this step.
- 2. Down-sampling the above x(n) by a factor of *M* to produce y(n).

We may view y(n) as though it were generated by sampling an underlying analog signal y(t) at a rate $F_y = F_x/M$ Hz.

Given the signal x(n) that was obtained at a certain sampling rate the new signal y(n), the down-sampled version of x(n), with a sampling rate that is (1/M) of that of x(n), obtained from x(n), is given by:

$$y(n) = x(Mn)$$

and is made up of every M^{th} sample value of x(n); the intervening (M-1) sample values of x(n) are dropped. This amounts to

$$y(0) = x(0), y(1) = x(M), y(2) = x(2M), y(3) = x(3M), \dots$$

The time between samples of y(.) is M times that between samples of x(.), or the sampling frequency of y(.) is reduced by a factor of M from that of x(.). The block diagram of a down sampler is shown below.



Example 5.3.1 As an example, if $x(n) = a^n u(n)$, a < 1, is the sequence:

$$n = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8$$

$$x(n) = \{1 \quad a \quad a^2 \quad a^3 \quad a^4 \quad a^5 \quad a^6 \quad a^7 \quad a^8 \quad \cdot \quad \cdot \}$$

then y(n) = x(2n), with M = 2, is its 2-fold down-sampled version and is obtained by keeping every other sample of x(n) and dropping the samples in between:

$$n = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5$$

$$y(n) = \{1 \quad a^2 \quad a^4 \quad a^6 \quad a^8 \quad a^{10} \quad . \quad .\}$$

In this example it is understood that the time between samples of y(n) is twice that between samples of x(n), or, the sampling rate of y(n) is one-half of that of x(n).

Example 7.3.2 Test the system y(n) = x(Mn), where *M* is a constant, for **time-invariance**.



Solution See also Unit I. For the input x(n) the output is

$$y(n) = T[x(n)] = x(Mn)$$

Delay this output by n_0 to get

$$y(n-n_0) = x(M(n-n_0)) = x(Mn-Mn_0) < (A)$$

Next, for the delayed input $x(n-n_0)$ the output is

$$y(n, n_0) = T[x(n-n_0)] = x(Mn) = x(Mn-n_0) < (B)$$

 \rightarrow e see that (A) \neq (B), that is, $y(n-n_0)$ and $y(n, n_0)$ are not equal. Delaying the input is not equivalent to delaying the output. So the system is *not* time-invariant. In other words the down-sampling operation is a time-varying system.

Spectrum of a down-sampled signal Given the signal x(n) whose spectrum is $X(\omega)$ or $X(e^{j\omega})$ we want to find the spectrum of y(n), the down-sampled version of x(n), denoted by $y(n) - Y(\omega)$.

Consider the periodic train of impulses, p(n), with period M

$$p(n) = 1, \quad n = 0, \pm M, \pm 2M, \dots$$

0, otherwise



The discrete Fourier series representation (see Example 1 in Unit II) of p(n) is

$$p(n) = \sum_{k=0}^{M-1} P_k e^{j2\pi k n/M}$$
, $0 \le n \le M-1$

The Fourier coefficients are given by M^{-1}

$$P_k = \frac{1}{M} \sum_{m=0}^{\infty} p(n) e^{-j 2\pi k n / M} = \frac{1}{M}, \qquad 0 \le k \le M - 1$$

Thus the DFS for p(n) is

$$p(n) = \frac{1}{M} \sum_{k=0}^{M-1} e^{j 2\pi k n / M}, \qquad 0 \le n \le M-1$$

Define the signal x'(n)

$$x'(n) = x(n)p(n) = x(n), \quad n = 0, \pm M, \pm 2M, \dots$$

0, otherwise

The sequence x'(n) consists of values of x(n) whenever $n = 0, \pm M, \pm 2M, \dots$, and zeros in between those points.



Define the down-sampled version y(n)

$$y(n) = x'(Mn) = x(Mn) p(Mn) = x(Mn)$$

The signal y(n) consists of values of x(Mn) at $n = 0, \pm 1, \pm 2, ...,$ but no zeros in between.

With y(n) = x'(Mn) = x(Mn) our objective is to find the spectrum $Y(\omega)$. Keep in mind that $X(\omega)$ periodic in ω since x(n) is a discrete-time sequence; and the same is true of $Y(\omega)$. Now the *z*-transform of y(n) is

$$Y(z) = \sum_{n=-\infty}^{\infty} y(n) z^{-n} = \sum_{n=-\infty}^{\infty} x'(Mn) z^{-n}$$

Set Mn = k: then n = k/M and the summation limits $n = \{-\infty \text{ to } \infty\}$ become $k = \{-\infty \text{ to } \infty\}$. Thus

$$Y(z) = \sum_{k=-\infty}^{\infty} x'(k) z^{-k/M} = \sum_{n=-\infty}^{\infty} x'(n) z^{-n/M} k$$

Here x'(n) = 0 except when *n* is a multiple of *M*. Substituting x(n) p(n) for x'(n) in the above equation,

$$Y(z) = \sum_{n=-\infty}^{\infty} x(n) p(n) z^{-n/M}$$

Substituting $\frac{1}{M} \sum_{k=0}^{M-1} e^{j2\pi k n/M}$ for $p(n)$ (from the DFS) in the above equation,
$$Y(z) = \sum_{n=-\infty}^{\infty} x(n) \left[\frac{1}{M} \sum_{k=0}^{M-1} e^{j2\pi k n/M} \right]_{-n/M} = \frac{1}{M} \sum_{k=0}^{M-1} \sum_{n=-\infty}^{\infty} x(n) e_{j2\pi kn/M} z_{-n/M}$$
$$= \frac{1}{M} \sum_{k=0}^{\infty} \sum_{n=-\infty}^{\infty} x(n) \left(e^{-j2\pi k/M} z^{1/M} \right)^{-n}$$
$$= X \left(e^{-j2\pi k/M} z^{1/M} \right)$$

MATLAB. To demonstrate the stretching and shifting of $X(\omega)$ to $X((\omega-2\pi k)/M)$ for k = 1 and M = 2, that is, $X((\omega-2\pi)/2)$. This is done in 3 steps: (1) $X(\omega)$, (2) $X(\omega/2)$, and (3) $X((\omega-2\pi)/2)$

w = -2*pi: pi/256: 2*pi; subplot(3, 1, 1), plot(w, cos(w)); xlabel ('\omega, rad/sample'), ylabel('X(\omega)'); grid; title ('X(\omega)') % subplot(3, 1, 2), plot(w, cos(w/2)); xlabel ('\omega, rad/sample'), ylabel('X(\omega /2)'); grid; title ('Stretched by factor 2: X(\omega /2)') % subplot(3, 1, 3), plot(w, cos((w-2*pi)/2));

xlabel ('\omega, rad/sample'), ylabel('X(($\langle nega - 2 \rangle pi)/2$)'); grid; title ('And shifted by $2 \rangle pi$: X(($\langle nega - 2 \rangle pi)/2$)')



where, for simplicity, we have used the notation $X(\omega)$ instead of $X(e^{j\omega})$. This expression for $Y(e^{j\omega})$ is a sum of M terms. Note that the function $X(\omega - 2\pi k)$ is a shifted (by $2\pi k$) version of $X(\omega)$ and $X(\omega/M)$ is a stretched (by a factor M) version of $X(\omega)$. Thus $Y(e^{j\omega})$ is the sum of M uniformly shifted and stretched versions of $X(e^{j\omega})$ each scaled by the factor (1/M). The shifting

in multiples of 2π corresponds to the factor $(\omega - 2\pi k)$ in the argument of X(.), and the stretching by the factor M corresponds to the M in $(\omega - 2\pi k)/M$. Note that the amount of shift is also affected by the factor M, that is, the amount of shift doesn't stay at $2\pi k$ but ends up being $2\pi k/M$.

The expression for $Y(e^{j\omega})$ contains a total of M versions of $X(e^{j\omega})$, one original and (M - M)1) shifted replicas. Each of these is also stretched by a factor of M, so $X(e^{j\omega})$ should have been preshrunk, that is, band limited, to π/M before undertaking the down-sampling. Writing out the expression for $Y(e^{j\omega})$ in full, we have $2\pi k$

$$Y(e^{j\omega}) = \frac{1}{M} \sum_{k=0}^{\infty} X \begin{bmatrix} 0 \\ 0 \\ M \end{bmatrix}$$

$$1 \begin{bmatrix} 0 \\ 0 \\ M \end{bmatrix} (\omega - 2\pi 1) (\omega - 2\pi (M - 1))]$$

$$M \begin{bmatrix} 0 \\ M \end{bmatrix} X \begin{bmatrix} 0 \\ M \end{bmatrix} + X \begin{bmatrix} 0 \\ M \end{bmatrix} + \dots + \begin{bmatrix} 0 \\ X \end{bmatrix} \begin{bmatrix} 0 \\ M \end{bmatrix}$$

The first term that makes up $Y(e^{j\omega})$, that is, $X \begin{bmatrix} \Box \\ M \end{bmatrix}$, is shown in the figure below. The figure

implicitly uses M = 2. In general there will be (M-1) shifted replicas of this term.



This is also written in the form

 $Y(e^{j\omega}) = \frac{1}{2} \sum X(e^{j(\omega-2\pi k)/2}) = \frac{1}{2} \left[X(e^{j\omega/2}) + X(e^{j(\omega-2\pi)/2}) \right]$

$$=\frac{1}{2}\left[X\left(e^{j\omega/2}\right)+X\left(e^{j\omega/2-j\pi}\right)\right]=\frac{1}{2}\left[X\left(e^{j\omega/2}\right)+X\left(-e^{j\omega/2}\right)\right]$$

To recapitulate, before we decided to down sample $X(\omega)$ was originally band limited to π on the digital frequency scale (that is, $F_x/2$ Hz). We then band limited it to π/M (that is, $F_x/2M$ Hz) and down sampled by a factor of M.

Aliasing Down-sampling by a factor of M, in itself, is simply retaining every M^{th} sample while dropping all samples in between. If, therefore, prior to down-sampling, the signal x(n) is indeed band-limited to π/M then we generate the down-sampled version y(n) by simply taking every M^{th} sample of x(n). This process is shown below in block diagram fashion. If in this set-up x(n) is not band-limited as required then the spectrum of y(n) will contain overlapping spectral components

of x(n) due to stretching, i.e., $X(\omega/M)$ will overlap $X(\omega-2\pi/M)$, etc. This results in **aliasing**.



Band-limiting x(n) to π/M (if not done already) is done by an **anti-aliasing filter** (digital low pass filter) with a cut-off frequency of π/M . The general process of *decimation then consists* of filtering followed by down sampling shown in block diagram below.



Unlike an analog anti-aliasing filter associated with an ADC, the filter in this diagram is a *digital* anti-aliasing filter specified as

$$H(\omega) = \begin{array}{c} 1, \quad 0 \leq |\omega| < \pi/M \\ 0, \quad \pi/M \leq |\omega| \leq \pi \end{array}$$

Note that π corresponds to $F_x/2$ and π/M corresponds to $F_x/2M$ where F_x is the sampling frequency of x(n).

Typically, in order to avoid (delay) distortion, the filter H(z) is a linear phase FIR filter with (N+1) coefficients $\{h(r), r = 0 \text{ to } N\}$. The output, v(n), of the low pass filter is then given by convolution

$$v(n) = \sum_{r=0}^{N} h(r) x(n-r)$$

and the decimated signal is

$$y(n) = v(nM) = \sum_{r=0}^{N} h(r) x(nM - r)$$

In summary, in order to down sample a signal by a factor of M:

- The signal should have been originally over-sampled by a factor of M (that is • originally band limited to π/M and over-sampled). In this case the signal is downsampled straightaway; no pre-filter is needed. OR
- The signal, assumed originally band limited to π , should be band-limited to π/M • by a pre-filter; the signal is then down-sampled. In this case there will be some loss of information.

Example 5.3.3 Consider the signal $x(n) = a^n u(n)$, a < 1.

- a) Determine the spectrum $X(\omega)$
- b) If x(n) is applied to a decimator that reduces the sampling rate by a factor of 2 determine the output spectrum
- c) Show that the spectrum in part (b) is simply the Fourier transform of x(2n)
- d) Plot the spectra of x(n) and x(2n) for a = 0.905

Solution [See also Unit I]

a) The spectrum of x(n) is given by its DTFT

$$X(\omega) = \sum_{\substack{n = -\infty \\ n = -\infty}} x(n) e^{-j\omega n} = \sum_{\substack{n = 0 \\ n = 0}} a^n e^{-j\omega n} = \sum_{\substack{n = 0 \\ n = 0}} \left(a e^{-j\omega} \right)^n$$
$$= \frac{1}{1 - a e^{-j\omega}}, \qquad \left| a e^{-j\omega} \right|^{<1}$$

This spectrum is not band-limited but we may pretend it is. This may also be obtained as $X(\omega)$ $=X(z)\Big|_{z=e^{j\omega}}.$ $\mathbf{v}(n) = \mathbf{x}(2n)$ $\mathbf{r}(\mathbf{n})$

b) The spectrum of
$$y(n) = x(2n)$$
 is given by

$$Y(\omega) = \frac{1}{M} \sum_{k=0}^{\infty} X \left[\begin{array}{c} \square \\ M \end{array} \right]$$
which, with $M = 2$, becomes $\left(\begin{array}{c} \omega - 2\pi k \right) \quad 1 \left[\begin{array}{c} \omega \\ M \end{array} \right]$

$$Y(\omega) = \frac{1}{2} \sum_{k=0}^{\infty} X \left[\begin{array}{c} \square \\ 2 \end{array} \right] = \frac{1}{2} \left[\begin{array}{c} X \\ \square \\ 2 \end{array} \right] \left[\begin{array}{c} \square \\ 2 \end{array} \left[\begin{array}{c} \square \\ 2 \end{array} \right] \left[\begin{array}{c} \square$$

c) The Fourier transform of y(n) = x(2n) = a u(2n) = a u(n) is

$$Y(\omega) = \sum_{n=0}^{\infty} a^{2n} e^{-j\omega n} = \sum_{\infty} \left(a^2 e^{-j\omega} - \frac{1}{1 - a^2 e^{-j\omega}}, - \frac{1}{1 - a^2 e^{-j\omega}} \right)^n$$

d) The spectra.

b = [1]; % Numerator coefficient a1 = [1, -0.905]; a2 = [1, -0.819]; %Denominator coefficients w = -pi: pi/256: pi; %A total of 512 points [X1w] = freqz(b, a1, w); [X2w] = freqz(b, a2, w)subplot(2, 1, 1), plot(w, abs(X1w)); legend ('Spectrum of x(n)');





Up-sampling

Assume that x(n) is obtained from the continuous-time signal x(t) by sampling at F_x Hz. We now wish to generate a signal y(n) by up-sampling x(n) by a factor of L, that is, we are increasing the sampling rate by a factor of L. This amounts to

- 1. Converting x(n) to x(t) using a D/A converter.
- 2. Resampling x(t) at LF_x Hz to produce y(n).

We may view y(n) as though it were generated by sampling an underlying analog signal y(t) (or x(t) for that matter) at a rate $F_y = LF_x$ Hz. As in the case of down-sampling we prefer to do this entirely in the digital domain.

Given the signal x(n) that was obtained at a certain sampling rate we can obtain a new signal y(n) from x(n) with a sampling rate that is *L* times that of x(n). The signal y(n), an upsampled version of x(n), is given by:

$$y(n) = x(n/L), \qquad n = 0, \pm L, \pm 2L, \dots$$

0, otherwise

and is constructed by placing (L-1) zeros between every pair of consecutive samples of x(n). The time between samples of y(n) is (1/L) of that between samples of x(n), or the sampling frequency of y(n) is increased by a factor of L from that of x(n). The block diagram of an up-sampler is shown below.



Example 5.4.1 As an example, if $x(n) = a^n u(n)$, a < 1, is the sequence:

 $n = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8$ $x(n) = \{1 \quad a \quad a^2 \quad a^3 \quad a^4 \quad a^5 \quad a^6 \quad a^7 \quad a^8 \quad \cdot \quad \cdot \}$

then y(n) = x(n/2), with L = 2, is its 2-fold up-sampled version and is obtained by inserting a 0 between each pair of consecutive values in x(n)

In this example it is understood that the time between samples of y(n) is one half of that between samples of x(n), or, the sampling rate of y(n) is twice that of x(n).

Example 5.4.2 Test the system y(n) = x(n/L), where *L* is a constant, for **time-invariance**.



Solution See also Unit I. The system y(n) = x(n/L) in itself only partially defines an up-sampler. But the following goes to show that the up-sampling operation is a time-varying system. For the input x(n) the output is

$$y(n) = T[x(n)] = x(n/L)$$

Delay this output by n_0 to get

$$y(n-n_0) = x((n-n_0)/L) = x((n/L)-(n_0/L)) < (A)$$

Next, for the delayed input $x(n-n_0)$ the output is

$$y(n, n_0) = T[x(n-n_0)] = x(n/L) = x((n/L)-n_0) < (B)$$

 \rightarrow e see that (A) \neq (B), that is, $y(n-n_0)$ and $y(n, n_0)$ are not equal. Delaying the input is not equivalent to delaying the output. So the system y(n) = x(n/L) is *not* time-invariant. Therefore the up sampler defined by

$$y(n) = x(n/L),$$
 $n = 0, \pm L, \pm 2L, \dots$
0, otherwise

is not time-invariant; it is a time-varying system.

Spectrum of an up-sampled signal Given the signal x(n) whose spectrum is $X(\omega)$ or $X(e^{j\omega})$ we want to find the spectrum of y(n), the up-sampled version of x(n), denoted by $y(n) - Y(\omega)$.

The signal y(n), with a sampling rate that is *L* times that of x(n), is given by:

$$y(n) = x(n/L), \qquad n = 0, \pm L, \pm 2L, \dots$$

0, otherwise

We obtain the *z*-transform and from it the spectrum:

$$Y(z) = \sum_{n = -\infty}^{\infty} y(n) z^{-n} = \sum_{n = 0, \pm L, \pm 2 L, \dots}^{\infty} x(n / L) z^{-n} + \sum_{\substack{n = \text{ other than } \pm kL \\ k = all \text{ integers}}}^{\infty} 0 z^{-n}$$
$$= \sum_{n = 0, \pm L, \pm 2 L, \dots}^{\infty} x(n / L) z^{-n}$$

Set n/L = k: this leads to n = kL, and the summation indices $n = \{0, \pm L, \pm 2L, \pm 3L, ...\}$ become $k = \{-\infty \text{ to } \infty\}$, so that

$$Y(z) = \sum_{k=-\infty}^{\infty} x(k) \ z^{-kL} = \sum_{k=-\infty}^{\infty} x(k) \left(z^{L} \right)^{k} = X \left(z^{L} \right)$$

Setting $z = e^{j\omega}$ gives us the spectrum

$$Y(e^{j\omega}) = Y(z)|_{z=e^{j\omega}} = X(e^{j\omega L})$$
 or $Y(\omega) = X(\omega L)$

Thus $Y(\omega)$ is an L-fold compressed version of $X(\omega)$; the value of X(.) that occurred at ωL occurs at ω , (that is, at $\omega L/L$) in the case of Y(.). In going from X to Y the frequency values are pushed in toward the origin by the factor L. For example, the frequency ωL is pushed to $\omega L/L$, the frequency π is pushed to π/L , 2π is pushed to $2\pi/L$, etc.

Shown below are the spectra $X(\omega)$ and $Y(\omega)$ for 2-fold up-sampling, that is, L = 2. Note that $X(\omega)$ is periodic to start with so that the frequency content of interest is in the base range $(-\pi \le \omega \le \pi)$ with replicas of this displaced by multiples of 2π from the origin on either side. Due to

up-sampling the frequency content of $X(\omega)$ in the range $(-\pi \le \omega \le \pi)$ is compressed into the range $(-\pi/L \le \omega \le \pi/L)$ of $Y(\omega)$, that is, into $(-\pi/2 \le \omega \le \pi/2)$, centered at $\omega = 0$. The first replica of $X(\omega)$ in the range $(\pi \le \omega \le 3\pi)$, centered at 2π , is compressed to the range $(\pi/2 \le \omega \le 3\pi/2)$ of $Y(\omega)$, centered at π ; its counterpart, in $(-3\pi \le \omega \le -\pi)$, centered at -2π , is compressed to $(-3\pi/2 \le \omega \le -\pi/2)$, centered at $-\pi$. If, for the purpose of discussion, we consider the range $(0, 2\pi)$ as one fundamental period then the replica in the range $(\pi/2, 3\pi/2)$ of *Y* is an **image (spectrum)** and needs to be filtered out with a low pass filter (**anti-imaging filter**) of band-width $\pi/2$. With L = 2 this is the only image in $(0, 2\pi)$.

Furthermore, while the spectrum $X(\omega)$ is periodic with a period = 2π , the spectrum $Y(\omega)$, on account of the image, is a 2-fold periodic repetition of the base spectrum in $(-\pi/2 \le \omega \le \pi/2)$; the image spectrum is actually spurious/unwanted; further the periodicity of $Y(\omega)$ is still 2π .



These observations can be extended to larger values of *L*. For L = 3, for instance, there will be two image spectra (a 3-fold periodic repetition of the base spectrum in $(-\pi/3 \le \omega \le \pi/3)$, and the anti-imaging filter band width will be $\pi/3$.

In general, up-sampling of x(n) by a factor of *L* involves

- Inserting L-1 zeros between successive pairs of sample values of x(n).
- The spectrum $Y(\omega)$ of the up-sampled signal is an *L*-fold compressed version of $X(\omega)$. As a result $Y(\omega)$ contains *L*-1 images and is an *L*-fold periodic repetition of the base spectrum in $(-\pi/L \le \omega \le \pi/L)$.
- The anti-imaging filter band width is π/L .

The over-all scheme of up-sampling is shown in block diagram below. Unlike an analog anti-imaging filter associated with a DAC, the filter in this diagram is a *digital anti-imaging filter*.



In this diagram the pass band gain of the anti-imaging filter is shown as 1. This gain is actually chosen equal to *L* to compensate for the fact that the average value of y(n) is 1/L times the average value of x(n) due to the presence of the inserted zeros.

$$H(\omega) = L, \quad 0 \le |\omega| < \pi/L$$
$$0, \quad \pi/L \le |\omega| \le \pi$$

Note that π corresponds to $F_x/2$ and π/L corresponds to $F_x/2L$ where F_x is the sampling frequency of x(n).

The output of the low pass filter is given by the convolution sum

$$y(n) = \sum_{r = -\infty}^{\infty} h(n-r) v(r)$$

where its input is

$$v(r) = x(r/L), \qquad r = 0, \pm L, \pm 2L, \dots$$

0, otherwise

Now v(r) = 0 except at r = kL, where k is all integers from $-\infty$ to ∞ . Thus we have

$$v(kL) = x(kL/L) = x(k)$$

The convolution sum may be written as

$$\sum_{r=-\infty}^{\infty} h(n-r) v(r) = \sum_{k=-\infty}^{\infty} h(n-kL) v(kL) = \sum_{k=-\infty}^{\infty} h(n-kL) x(k)$$

so that the interpolated signal is

$$y(n) = \sum_{k=-\infty} h(n - kL) x(k)$$

Illustration Given the signal $x(n) = \{1, a, a^2, a^3, a^4, a^5, a^6, a^7, a^8, a^9, a^{10}, ...\}$, its 2-fold upsampled version is obtained by inserting a 0 between each pair of consecutive samples in x(n):

 $y(n) = \{1, 0, a, 0, a^2, 0, a^3, 0, a^4, 0, a^5, 0, a^6, 0, a^7, 0, a^8, 0, a^9, 0, a^{10}, ...\}$
Intuitively, even visually, y(n) contains higher (or, more) frequencies than x(n) because of the inserted zeros. For instance, consider the first two or three samples in each sequence. In the case of x(n) the changes from 1 to a to a^2 are smoother than the fluctuations in y(n) from 1 to 0 to a to 0 to a^2 ; these latter fluctuations are the higher frequencies not originally contained in x(n). It is these higher frequencies that are represented by the image in the spectrum of y(n) prior to anti-imaging filtering. The anti-imaging filter removes or smoothes out the higher frequency fluctuations from the up-sampled version; *this smoothing is manifested in the form of the interpolated zeros being replaced by nonzero values*.

Cascading sampling rate converters

Given a discrete-time signal x(n) we may want to convert its sampling rate by a non integer factor, in particular, by a rational number. For instance, we may be interested in x(3n/2). This involves a 2-fold up-sampling and a 3-fold down-sampling for a net down sampling by a factor of 1.5 (= 3/2). The sequence x(3n/5) involves a 5-fold up-sampling and a 3-fold down-sampling for a net up-sampling by a factor of 1.67 (= 5/3).

In general, in a cascade of an *M*-fold down-sampler and an *L*-fold up-sampler the positions of the two samplers are inter-changeable with no difference in the input-output behavior *if and only if M and L are co-prime (relatively prime*, that is, *M* and *L* do not have a common factor). The sequence x(3n/2) may be generated by cascading the up-sampler and the down-sampler in either order, that is, down followed by up or vice versa. However, a cascade of a 6-fold down-sampler (M = 6) followed by a 4-fold up-sampler (L = 4) is not the same as a cascade of a 4-fold up-sampler followed by a 6-fold down-sampler even though in both cases M/L = 6/4. This is because *M* and *L* have a common factor, that is, the rational number M/L is not in its reduced form. The ratio M/L should be reduced to 3/2; then the 3-fold down-sampler and the 2-fold up-sampler are interchangeable in position.

Example 5.5.1 Given $x(n) = e^{-n/2} u(n)$, find x(5n/3).

Answer We borrow this from an earlier section. Our objective is to present the earlier solution and then reformulate it in the context of cascading up- and down-samplers. The sequence

$$x(n) = e^{-n/2} u(n) = (e^{-1/2})^n u(n) = (0.606)^n u(n) = a^n u(n)$$

where $a = e^{-l/2} = 0.606$, is sketched below:



With y(n) = x(5n/3), we evaluate y(.) for several values of *n* (we have assumed here that x(n) is zero if *n* is not an integer):

 $\begin{array}{l} y(0) = x(5 \ . \ 0 \ / \ 3) = x(0) = e^{-0/2} = 1 \\ y(1) = x(5 \ . \ 1 \ / \ 3) = x(5 \ / \ 3) = 0 \\ y(2) = x(5 \ . \ 2 \ / \ 3) = x(10 \ / \ 3) = 0 \\ y(3) = x(5 \ . \ 3 \ / \ 3) = x(5) = e^{-5/2} = a^5 \\ \cdots \\ y(6) = x(5 \ . \ 6 \ / \ 3) = x(10) = e^{-10/2} = a^{10} \end{array}$

The general expression for y(n) can be written as

$$y(n) = x(5n/3) = e^{-(5n/3)/2}$$
, *n* as specified below
= $e^{-5n/6}$, $n = 0, 3, 6, ...$
0, otherwise

The sequence is sketched below:



We shall recast this problem in terms of cascading the up- and down-samplers. In the expression y(n) = x(5n/3) there is a 3-fold up-sampling and a 5-fold down-sampling. Since the numerator 5 is greater than the denominator 3 there is a *net down-sampling by a factor of 1.67* (= 5/3). Let us first do a 3-fold up-sampling of x(n) followed by a 5-fold down-sampling of the resulting sequence. That is, given the sequence x(n)

$$x(n) \qquad \qquad \uparrow 3 \qquad y_u(n) = x(n/3) \qquad \qquad \downarrow 5 \qquad y(n) = x(5n/3)$$

$$n = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9 \quad 10 \quad . \quad .$$

$$x(n) = \{1 \quad a \quad a^2 \quad a^3 \quad a^4 \quad a^5 \quad a^6 \quad a^7 \quad a^8 \quad a^9 \quad a^{10} \quad . \quad .\}$$

we define $y_u(n) = x(n/3)$, and then $y(n) = y_u(5n) = x(5n/3)$. The sequences $y_u(n)$ and y(n) are given below.

$$y_u(n) = x(n/3) = e^{-n/6} u(n/3) = a^{n/3}$$
 $n = 0, 3, 6, ...$
0, otherwise

$$y(n) = y_u(5n) = x(5n/3) = e^{-5n/6} u(5n/3) = a^{5n/3}$$
 $n = 0, 3, 6, ...$
0, otherwise

Alternatively, we may first do a 5-fold down sampling followed by a 3-fold up-sampling:

$$\begin{array}{c} x(n) & y_d(n) = x(5n) \\ \downarrow 5 & \uparrow 3 \end{array} \begin{array}{c} y(n) = x(5n/3) \\ \uparrow 3 & \downarrow 5 \end{array}$$

$$y_d(n) = x(5n) = \{1, a^5, a^{10}, a^{15}, a^{20}, \dots \}$$

$$y(n) = y_d(n/3) = x(5n/3) = \{1, 0, 0, a^5, 0, 0, a^{10}, 0, 0, a^{15}, 0, 0, a^{20}, \dots \}$$

The net effect is that between the first two terms (1 and a^5) of the final output y(.) we have dropped four original terms and inserted two zeros.

Example 5.5.2 Given $x(n) = e^{-n/2} u(n)$, find x(3n/5). Here there is a 5-fold up-sampling and a 3-fold down sampling. Since the denominator is bigger there is a *net up-sampling by a factor of* **1.67**.

$$x(n) = \{1, a, a^2, a^3, a^4, a^5, a^6, a^7, a^8, a^9, a^{10}, \ldots\}$$

Method A Up-sampling followed by down sampling is given below. The 5-fold up-sampled signal, $y_u(n)$, is obtained by inserting 4 zeros shown in bold face between every pair of consecutive samples in x(n)

$$y_u(n) = x(n/5)$$

= {1, 0, 0, 0, 0, a, 0, 0, 0, a², 0, 0, 0, 0, a³, 0, 0, 0, 0, a⁴, 0, 0, 0, 0, a⁵, 0, 0, 0, 0, a⁶, 0, 0, 0, 0, a⁷, 0, 0, 0, 0, a⁸, 0, 0, 0, 0, a⁹, 0, 0, 0, 0, a¹⁰, ...}

The 3-fold down-sampled signal, $y_1(n)$, is obtained by keeping every third sample in $y_u(n)$ and discarding the rest (shown underlined)

Method B Down-sampling followed by up sampling is given below. The 3-fold down-sampled signal, $y_d(n)$, is obtained by keeping every third sample in x(n) and discarding the rest (shown umderlined)

$$x(n) = \{1, \underline{a}, \underline{a^2}, a^3, \underline{a^4}, \underline{a^5}, a^6, \underline{a^7}, \underline{a^8}, a^9, \underline{a^{10}}, \underline{a^{11}}, \dots \}$$
$$y_d(n) = x(3n) = \{1, a^3, a^6, a^9, a^{12}, \dots \}$$

 $v_2(n) = v_3(n/5) = x(3n/5)$

The 5-fold up-sampled signal, $y_2(n)$, is obtained by inserting 4 zeros shown in bold face between every pair of samples in $y_d(n)$

$$y_2(n) = \{1, \mathbf{0}, \mathbf{0}, \mathbf{0}, \mathbf{0}, a^3, \mathbf{0}, \mathbf{0}, \mathbf{0}, \mathbf{0}, a^6, \mathbf{0}, \mathbf{0}, \mathbf{0}, \mathbf{0}, a^9, \mathbf{0}, \mathbf{0}, \mathbf{0}, \mathbf{0}, a^{12}, \ldots\} = \{1, 0, 0, 0, 0, a^3, 0, 0, 0, 0, a^6, 0, 0, 0, 0, a^9, 0, 0, 0, 0, a^{12}, \ldots\}$$

It is seen that $y_1(n) = y_2(n)$.

Example 5.5.3 Given that $x(n) = \{1, a, a^2, a^3, a^4, a^5, a^6, a^7, a^8, a^9, a^{10}, ...\}$ is the input,

- 1. Find the output $y_{I}(n)$ of a cascade of a 2-fold up-sampler followed by a 4-fold down sampler.
- 2. Find the output $y_2(n)$ of a cascade of a 4-fold down sampler followed by a 2-fold up-sampler.

Solution Note that the down-sampling factor M = 4 and the up-sampling factor L = 2 are not coprime since they have a factor in common. The ratio M/L = 4/2, as given, is not in its reduced form. As a result we do not expect that $y_1(n)$ and $y_2(n)$ will be equal. Specifically, in the first case we have

$$x(n) = \{1, a, a^2, a^3, a^4, a^5, a^6, a^7, a^8, a^9, a^{10}, \ldots\}$$

Up-sample by inserting a zero (shown bold face) between consecutive samples of x(n) resulting in $y_u(n)$

$$y_u(n) = x(n/2) = \{1, 0, a, 0, a^2, 0, a^3, 0, a^4, 0, a^5, 0, a^6, 0, a^7, 0, a^8, 0, a^9, 0, a^{10}, \ldots\}$$

Down-sample by keeping every fourth sample of $y_u(n)$ and discarding the three samples in between resulting in $y_l(n)$

$$y_1(n) = y_u(4n) = x(4n/2) = \{1, a^2, a^4, a^6, a^8, a^{10}, \ldots\}$$

In the second case

$$x(n) = \{1, a, a^{2}, a^{3}, a^{4}, a^{5}, a^{6}, a^{7}, a^{8}, a^{9}, a^{10}, \ldots\}$$
$$y_{d}(n) = x(4n) = \{1, a^{4}, a^{8}, a^{12}, \ldots\}$$

$$y_2(n) = y_d(n/2) = x(4n/2) = \{1, 0, a^4, 0, a^8, 0, a^{12}, \ldots\}$$

It is seen that $y_1(n) \neq y_2(n)$.

Sampling Rate Conversion by a Rational Factor L/M Here the sampling rate is being converted by a non-integral factor such as 0.6 or 1.5. That is, given x(n) with a sampling rate of F_x we want to obtain y(n) with a sampling rate of F_y of, say, $0.6F_x$ (decimation) or $1.5F_x$ (interpolation).

Take, for instance L/M = 3/5. Here the basic approach is to first interpolate (up-sample) by a factor of L = 3 and then decimate (down-sample) by a factor of M = 5. The net effect of the cascade of interpolation followed by decimation is to change the sampling rate by a rational factor L/M, that is,

$$F_{y} = \begin{pmatrix} L \\ \Box \\ M \end{pmatrix} F_{x} = \begin{pmatrix} 3 \\ -5 \end{pmatrix}_{x} = 0.6 F_{x}$$

The corresponding signal is given by y(n) = x(5n/3), ignoring the filters involved. (This can also be done by first down-sampling and then up-sampling).

The block diagram of the scheme where the interpolator precedes the decimator is shown below.

Rational sampling rate conversion



In general, if L < M we have a rational decimator and if L > M we have a rational interpolator. In this set-up interpolation is done before decimation in order to work at the higher sampling rate so as to preserve the original spectral characteristics of x(n). Recall that unless x(n) was originally over-sampled, decimation in itself or decimation prior to interpolation will modify the spectrum of x(n) irrecoverably.

The above configuration has an added benefit that the two filters $H_u(z)$ and $H_d(z)$ in series (which operate at the same sampling rate) can be combined into a single equivalent low pass filter with a frequency response of $H(\omega) = H_u(\omega)H_d(\omega)$. The simplified configuration is shown below.

Rational sampling rate conversion



The bandwidth of the anti-imaging filter $H_u(z)$ is π/L rad., and that of the anti-aliasing filter $H_d(z)$ is π/M rad., so that the bandwidth of the composite anti-imaging and anti-aliasing filter $H(\omega)$ is

$$\omega_c = \min\left(\frac{\pi, \pi}{L, M}\right)$$

and the frequency response is given by

$$H(\omega) = L, \quad 0 \le |\omega| < \omega_c$$
$$0, \quad \omega_c \le |\omega| \le \pi$$

In the time domain, the output of the up-sampler, v(n), is given by

$$v(n) = x(n/L), \qquad n = 0, \pm L, \pm 2L, \dots$$

0, otherwise

and the output of the linear time-invariant filter H(z) is

$$w(n) = \sum_{k = -\infty} h(n-k) v(k)$$

Since v(k) = 0 except at k = rL, where *r* is an integer between $-\infty$ to ∞ , we set k = rL. As *k* goes from $-\infty$ to ∞ , *r* goes from $-\infty$ to ∞ , and v(rL) = x(rL/L) = x(r).

$$w(n) = \sum_{r=-\infty}^{\infty} h(n - rL) \ v(rL) = \sum_{r=-\infty}^{\infty} h(n - rL) \ x(r) = \sum_{k=-\infty}^{\infty} h(n - kL) \ x(k)$$

Finally the output of the down sampler is

$$y(n) = w(Mn) = \sum_{k = -\infty} h(Mn - kL) x(k)$$

In summary, sampling rate conversion by the factor L/M can be achieved by first increasing the sampling rate by L, accomplished by inserting L-1 zeros between successive samples of the input x(n), followed by linear filtering of the resulting sequence to eliminate unwanted images of $X(\omega)$ and, finally, by down-sampling the filtered signal by the factor M to get the output y(n). The sampling rates are related by $F_y = (L/M)F_x$. If $F_y > F_x$, that is, L > M, the low pass filter acts as an anti-imaging post-filter to the up-sampler. If $F_y < F_x$, that is, L < M, the low pass filter acts as an anti-aliasing pre-filter to the down-sampler.

Example 5.5.4 The signal $x(t) = \cos 2\pi 2t + 0.8 \sin 2\pi 4t$ is sampled at 40 Hz to generate x(n).

- *a*) Give an expression for x(n)
- b) Design a sampling rate converter to change the sampling frequency of x(n) by a factor of 2.5. Give an expression for y(n).
- *c)* Design a sampling rate converter to change the sampling frequency of x(n) by a factor of 0.4. Give an expression for y(n).

Solution

a) $x(n) = \cos 2\pi 2nT + 0.8 \sin 2\pi 4nT = \cos \pi n/10 + 0.8 \sin \pi n/5$

b) The rate conversion factor is L/M = 2.5 = 5/2. We do this by first up-sampling by a factor of 5, then down-sampling by a factor of 2. Roughly speaking, y(n) = x(2n/5).

Rational sampling rate conversion



The up-sampling requires an anti-imaging LP filter of bandwidth of $\pi/L = \pi/5$ rad., and a gain of 5. The down-sampling requires an anti-aliasing LP filter of band width $\pi/M = \pi/2$ rad. When the up-sampler precedes the down-sampler we have the following configuration where the filter has the smaller of the two band widths, that is, $\pi/5$ rad. The gain is 5 and is always determined by the up-sampler.

Write equations for v(n), w(n), and y(n) based on the above diagram and assuming H(z) is an FIR filter of N coefficients.

c) The rate conversion factor is L/M = 0.4 = 2/5. We do this by first up-sampling by a factor of 2, then down-sampling by a factor of 5. Roughly speaking, y(n) = x(5n/2). Band width of filter = $\pi/5$ rad., and gain = 2, once again determined by the up-sampler.



Write equations for v(n), w(n), and y(n) based on the above diagram and assuming H(z) is an FIR filter N coefficients.

Multistage conversion The composite anti-imaging and anti-aliasing LP filter band width is π/L or π/M whichever is smaller. That is, the filter band width is determined by the larger number of *L* and *M*. However, if either *L* or *M* is a very large number the filter bandwidth is very narrow.

Narrowband FIR (linear phase) filters can require a very large number of coefficients (see Unit VI, FIR Filters, Example 6). This can pose problems in

- 1. Increased storage space for coefficients,
- 2. Long computation time, and
- 3. Detrimental finite word length effects

The latter drawback is minimized by using a multistage sampling rate converter where the conversion ratio L/M is factored into the product of several ratios each of which has its own smaller L and M valu/es. If for in stance, the ratio L/M is split into the product of two ratios

$$\frac{L}{M} = \left\| \begin{array}{c} L_{1} \\ L_{1} \\ M \end{array} \right\| \left\| \begin{array}{c} L_{2} \\ M \\ M_{1} \\ M \end{array} \right\|$$

where the L's and M's on the right hand side are smaller, we may implement the rate conversion in two stages as shown below



Example 7.5.5 [Rational sampling rate converter][CD, DAT] Digital audio tape (DAT) used in sound recording studios has a sampling rate of 48 kHz, while a compact disc (CD) is recorded at a sampling rate of 44.1 kHz. Design a sampling rate converter that will convert the DAT signal x(n) to a signal y(n) for CD recording.

Over-sampling analog-to-digital converter (ADC) [Ref. SKMitra, Sec 4.8.4] A practical difficulty with analog to digital conversion is the need for a low pass analog anti-aliasing prefilter to band limit the signal to less than half of the sampling rate. High-order analog filters are expensive, and they are also difficult to keep in calibration. The combination of over-sampling followed by down-sampling can be used to transfer some of the anti-aliasing burden from the analog into the digital domain, and thereby use a simpler low order analog filter.

As an example, a typical compact disk encoding system may employ an *over-sampling sigma-delta* A/D *converter* which over-samples at 3175.2 kHz which is then brought down to the CD sampling rate of 44.1 kHz. This amounts to over-sampling by a factor of 72 (= 3175.2/44.1).

Suppose the range of frequencies of interest in the signal x(t) is $0 \le |F| \le F_M$. Normally we would band-limit x(t) to the maximum frequency F_M with a sharp cut-off analog low pass filter and sample it at a rate of $F_s = 2F_M$ at least (strictly, $F_s \ge 2F_M$). Suppose that we instead over-sample x(t) by an integer factor M, at $F_s = M$ ($2F_M$). This significantly reduces the requirements for the anti-aliasing filter which may be specified more leniently as

$$H_a(s) = 1, \quad 0 \le |F| \le F_M$$

$$0, \quad MF_M \le |F| < \infty$$

Though this is still an ideal low pass filter in the pass band and stop band, its transition band width is no longer zero and it may be approximated with an inexpensive first- or second-order Butterworth filter as shown below.



We are paying the price of a *higher sampling rate* for the benefit of a *cheaper analog anti-aliasing filter*. The result is a discrete-time signal that is sampled at a much higher rate than $2F_M$. Following the sampling operation, we can reduce this sampling rate to the minimum value using a decimator. The resulting structure of the *over-sampling ADC* is shown in the block diagram below. There are *two* anti-aliasing filters, a low-order analog filter $H_a(s)$ with cut-off frequency F_M rad/sec., and a high-order digital filter $H_d(z)$, with a cut-off frequency (π/M) rad/sample.



A second benefit of using the over-sampling ADC is the *reduction in quantization noise*. If q is the quantization step size (precision), then the quantization noise in x(n) is $\sigma_x^2 = q^2 \frac{1}{2}$, and the noise appearing in the output y(n) in the above scheme is $\sigma_y^2 = \sigma_x^2 M = q^2 \frac{1}{2} M$, a reduction by a factor of M.

Identities

A sampling rate converter (the M or L operation) is a linear time-varying system. On the other hand, the filters $H_d(z)$ and $H_u(z)$ are linear time-invariant systems. In general, the order of a sampling rate converter and a linear time-invariant system *cannot* be interchanged. We derive below several identities, two of which are known as **noble identities** (viz., identities 3 and 6, all the others being special cases), which help to swap the position of a filter with that of a down-sampler or up-sampler by *properly modifying the filter*.

Recall that the input-output description of a down-sampler is

$$y(n) = x(Mn) - Y(z) = \frac{1}{M} \sum_{k=0}^{M-1} X\left(e^{-j2\pi k/M} z^{1/M}\right) - Y(e^{j\omega}) = \frac{1}{M} \sum_{k=0}^{M-1} X\left(e^{j(\omega-2\pi k)/M}\right)$$

and the same for an up-sampler is

$$y(n) = x(n/L), n = 0, \pm L, \pm 2L, \dots - Y(z) = X(z^L) - Y(e^{j\omega}) = X(e^{j\omega L})$$

0, otherwise

which we use in the following development.

Example 7.6.1 Show that the following systems are equivalent.



In the structure on the left the process of generating y(.) consists of multiplying every input sample x(.) by a and then, in the down-sampling process, dropping (M-1) of these products for every M^{th} one we keep. The structure on the right is more efficient computationally: the (M-1) samples are dropped first and every M^{th} is multiplied by a, that is, only the samples that are retained are multiplied. The number of multiplications is reduced by $\frac{100(M-1)}{M}$ %.

Example 5.6.2 Show that the following systems are equivalent.



In the structure on the left the process of generating y(.) consists of first up-sampling, that is, inserting zeros between consecutive points of x(n) and then multiplying by a. In the process the (L-1) zeros are also multiplied. The structure on the right is more efficient computationally: the sequence x(n) is first multiplied and then the zeros inserted. The number of multiplications is reduced by $\frac{100(L-1)}{r}$ %.

Identity #1 If we use the notation $M{.}$ to mean the down-sampling of the signal in braces, then we have

$$\$M\{ax_1(n) + bx_2(n)\} = \$M\{ax_1(n)\} + \$M\{bx_2(n)\}$$

$$= a\$M\{x_1(n)\} + b\$M\{x_2(n)\}$$

$$< (2)$$

In words, the result of down-sampling the weighted sum of signals equals the weighted sum of the down-sampled signals. In other words, the two block diagrams below are equivalent.



In the diagram on the left the weighted sum of two inputs is down-sampled:

$$w(n) = a x_1(n) + b x_2(n)$$
 — $W(e^{j\omega}) = a X_1(e^{j\omega}) + b X_2(e^{j\omega})$

The output and its spectrum are then given by

$$y(n) = w(Mn)_{M-1}$$

$$Y(e^{j\omega}) = \frac{1}{M} \sum_{\substack{k \\ M - 1}} W\left(e^{j(\omega - 2\pi k)/M}\right)$$

$$= \frac{a}{M} \sum_{\substack{k \\ M - 2}} X_1\left(e^{j(\omega - 2\pi k)/M}\right) + \frac{b}{M} \sum_{\substack{k \\ M - 2}} X_2\left(e^{j(\omega - 2\pi k)/M}\right) \quad < (A)$$

In the diagram on the right the weighted inputs are down-sampled and added to form the output. M^{-1}

$$y_{1}(n) = a x_{1}(Mn) \qquad - \qquad Y_{1}(e^{j\omega}) = \frac{a}{M} \sum_{\substack{k=0 \ b \ M-1}}^{m-1} X_{1}\left(e^{j(\omega-2\pi k)/M}\right)$$
$$y_{2}(n) = b x_{2}(Mn) \qquad - \qquad Y_{2}(e^{j(\omega-2\pi k)/M}) = \sum_{\substack{m=0 \ M \ k=0}}^{m-1} X_{2}\left(e^{j(\omega-2\pi k)/M}\right)$$

The output and its spectrum are given by

$$y(n) = y_{1}(n) + y_{2}(n)$$

$$Y(e^{j\omega}) = Y_{1}(e^{j\omega}) + Y_{2}(e^{j\omega}) = \frac{a}{M} \sum_{k=0}^{M-1} \left(e_{j(\omega-2\pi k)/M} \right) + \frac{b}{M} \sum_{k=0}^{M-1} X_{2} \left(e^{j(\omega-2\pi k)/M} \right) < (B)$$

Equations (A) and (B) are identical. QED.

Eventually, by virtue of Example 6.1, the down-samplers are moved to the upstream side of the multipliers which would correspond to Eq. (2).

Identity #2 A delay of *M* sample periods before an *M*-fold down-sampler is the same as a delay of one sample period after the down-sampler.

$$x(n) \xrightarrow{ z^{-M} y_1(n) } \downarrow M \xrightarrow{ y(n) } x(n) \xrightarrow{ x(n) } \downarrow M \xrightarrow{ y_2(n) } z^{-1} \xrightarrow{ y(n) }$$

In the first case (diagram on the left) we have

$$y_l(n) = x(n-M)$$
 — $Y_l(z) = z^{-M} X(z)$ < (1)

and

$$y(n) = y_{I}(Mn) \qquad - \qquad Y(z) = \frac{1}{M} \sum_{k=0}^{M-1} Y_{1} \left(e^{j 2 - \pi_{k/M}} z^{1/M} \right) \qquad < (2)$$

Note from (1) that

$$Y_{1}\left(e^{-j2\pi k/M} z^{1/M}\right) = Y_{1}(z)\Big|_{z=e^{-j2\pi k/M} z^{1/M}} = z^{-M}X(z)\Big|_{z=e^{-j2\pi k/M} z^{1/M}}$$

Substituting this in (2) $= \underbrace{\left(e^{-j2\pi k/M} z^{1/M}\right)^{-M} X\left(e^{-j2\pi k/M} z^{1/M}\right)}_{M-1} X\left(e^{-j2\pi k/M} z^{1/M}\right)$ $Y(z) = \frac{1}{M} \sum_{k=0}^{M} \left(e^{-j2\pi k/M} z^{1/M}\right)^{-M} X\left(e^{-j2\pi k/M} z^{1/M}\right)$ $= \frac{1}{M} \sum_{k=0}^{M-1} \left(e^{j2\pi kM/M} z^{-M/M}\right) X\left(e^{-j2\pi k/M} z^{1/M}\right) = \frac{1}{M} \sum_{k=0}^{M-1} \left(1 z^{-1}\right) X\left(e^{-j2\pi k/M} z^{1/M}\right)$ $= z^{-1} \frac{1}{M} \sum_{k=0}^{M-1} X\left(e^{-j2\pi k/M} z^{1/M}\right) < (A)$

In the second case (diagram on the right) we have

$$y_{2}(n) = x(nM) \qquad - \qquad Y_{2}(z) = \frac{1}{M} \sum_{k=0}^{M-1} X\left(e^{-2\pi k/M} z^{1/M}\right) \qquad < (3)$$
$$y(n) = y_{2}(n-1) \qquad - \qquad Y(z) = z^{-1} Y_{2}(z) \qquad < (4)$$

Substituting from (3) into (4) we have

$$Y(z) = z^{-1} Y_2(z) = z^{-1} \frac{1}{M} \sum_{k=0}^{M-1} X\left(e^{-2\pi k/M} z^{1/M}\right) < (B)$$

Equations (A) and (B) are identical. QED.

Identity #3 (Noble identity) An *M*-fold down-sampler followed by a linear time invariant filter H(z) is equivalent to a linear time invariant filter $H(z^M)$ followed by an *M*-fold down-sampler. Note that the second identity is a special case of this identity with $H(z) = z^{-1}$ and $H(z^M) = z^{-M}$.



For the system consisting of the filter followed by the down-sampler we have

$$Y_1(z) = H(z^M) X(z)$$

$$Y(z) = \frac{1}{M} \sum_{k} Y^{k} \left(e^{-j2\pi k/M} z^{1/M} \right)$$

Note that

$$Y_{1}\left(e^{-j2\pi k/M} z^{1/M}\right) = Y(z) \Big|_{z=e^{-j2\pi k/M} z^{1/M}}$$

= $H\left(z^{M}\right) X(z) \Big|_{z=e^{-j2\pi k/M} z^{1/M}}$
= $H\left(\left(e^{-j2\pi k/M} z^{1/M}\right)^{M}\right) X\left(e^{-j2\pi k/M} z^{1/M}\right)$
= $H\left(e^{-j2\pi kM/M} z^{M/M}\right) X\left(e^{-j2\pi k/M} z^{1/M}\right)$
= $H(1z) X\left(e^{-j2\pi k/M} z^{1/M}\right) = H(z) X\left(e^{-j2\pi k/M} z^{1/M}\right)$

Thus

$$Y(z) = \frac{1}{M} \sum_{k \in Q}^{M-1} H(z) X \left(e^{-j 2\pi k/M} z^{1/M} \right)$$

= $H(z) \frac{1}{M} \sum_{k=0}^{M-1} X \left(e^{-j 2\pi k/M} z^{1/M} \right)$ < (A)

For the system consisting of the down sampler followed by the filter we have $y_2(n) = x(nM)$

$$Y_{2}(z) = \frac{1}{M} \sum_{m=1}^{M-1} X\left(e^{-2\pi k/M} z^{1/M}\right)$$

$$Y(z) = H(z) Y_{2}(z) = H(z) \frac{1}{M} \sum_{m=1}^{M-1} X\left(e^{-2\pi k/M} z^{1/M}\right) \quad < (B)$$

Equations (A) and (B) are identical. Thus the two systems are equivalent.

Identity #4 (This identity contains no summing junction as does identity #1.) Eventually, by virtue of Example 6.2, the up-samplers are moved to the downstream side of the multipliers.



Identity #5 A delay of one sample period before an *L*-fold up-sampler is the same as a delay of *L* sample periods after the up-sampler.



For the system consisting of the up-sampler preceded by z^{-1} we have

$$Y_{l}(z) = z^{-l} X(z) Y(z) = Y_{l}(z^{L}) = z^{-L} X(z^{L}) \prec (A)$$

For the system consisting of the up-sampler followed by z^{-L} we have

$$Y_{2}(z) = X(z^{L})$$

$$Y(z) = z^{-L} Y_{2}(z) = z^{-L} X(z^{L}) \prec (B)$$

Since equations (A) and (B) are identical the two systems are equivalent.

Identity #6 (Noble identity) An *L*-fold up-sampler preceded by a linear time invariant filter H(z) is equivalent to a linear time invariant filter $H(z^L)$ preceded by an *L*-fold up-sampler. Note that the fifth identity is a special case of this identity with $H(z) = z^{-1}$ and $H(z^L) = z^{-L}$.



For the system consisting of the filter followed by the up-sampler we have

$$Y_{l}(z) = H(z) X(z) Y(z) = Y_{l}(z^{L}) = H(z^{L}) X(z^{L}) < (A)$$

For the system consisting of the up-sampler followed by the filter we have

$$Y_{2}(z) = X(z^{L}) Y(z) = H(z^{L}) Y_{2}(z) = H(z^{L}) X(z^{L}) < (B)$$

Since equations (A) and (B) are identical the two systems are equivalent.

FIR implementation of sampling rate conversion

The anti-aliasing filter in a decimator and the anti-imaging filter in an interpolator may each be either an FIR or an IIR filter, the former being preferred since it offers linear phase. We give here the FIR implementation.

Implementation of Decimator The process of decimation consists of an anti-aliasing (low pass) filter followed by a down-sampler. We repeat below the block diagram developed earlier.



Taking the filter H(z) to be an FIR filter, the decimator is implemented as shown below, using a direct form structure for H(z). Note that the coefficients $\{b_i, i = 0 \text{ to } (N-1)\}$, used in earlier formulations, are the same as $\{h(i), i = 0 \text{ to } (N-1)\}$ used in this diagram. Further, the FIR filter here is implemented with *N* coefficients rather than (N+1) coefficients.



The implementation equations which correspond to the above structure are N^{-1}

$$v(n) = \sum_{r=0}^{\infty} h(r) x(n-r)$$
 and $y(n) = v(Mn) = \sum_{r=0}^{\infty} h(r) x(Mn-r)$

We first compute v(n) for all values of n. Then y(n) is obtained by retaining every M^{th} value of v(.), dropping the intervening (M-1) values. In other words there are (M-1) computations of v(.) that could be avoided.

We may use identity #1 to move the down-sampler to the left of the adders and the result of Example 6.1 to move it to the upstream side of the multipliers as shown below. As a result the number of multiplications is reduced by $\frac{100(M-1)}{M}$ %.



Implementation of Interpolator The process of interpolation consists of an up-sampler followed by an anti-imaging (low pass) filter. We repeat below the block diagram developed earlier.



Taking the filter H(z) to be an FIR filter, the interpolator is implemented as shown below, using a direct form structure for H(z). Note that the coefficients $\{b_i, i = 0 \text{ to } (N-1)\}$, used in earlier formulations, are the same as $\{h(i), i = 0 \text{ to } (N-1)\}$ used in this diagram. Further, the FIR filter here is implemented with *N* coefficients rather than (N+1) coefficients.

We shall find it more convenient to use the transposed form of the FIR filter rather than the structure actually shown here, so here follows a digression on the transposed structure.



Start of Digression

Transposed Structure According to the transposition theorem the transposed form of a filter has the same transfer function as the filter. The transposed form of a given filter structure is found as follows:

- 1. Construct the signal flow graph of the filter.
- 2. Reverse the direction of arrow on every branch.
- 3. Interchange the inputs and outputs.
- 4. Reverse the roles of all nodes: an adder becomes a pick-off point and a pick-off point becomes an adder.

If we apply this procedure to the FIR structure in the above interpolator the result is the transpose shown below. The intermediate steps are omitted.



Start of Aside

As an aside note that the FIR structure is simple enough that the following algebraic manipulation can be used to proceed from the original FIR structure to the transpose structure. Let the system function be

$$\frac{V(z)}{U(z)} = H(z) = h(0) + h(1)z^{-1} + h(2)z^{-2} + h(3)z^{-3} + h(4)z^{-4} + \dots + h(N-1)z^{-(N-1)}$$

This may be rearranged as

V

$$(z) = H(z)U(z) = h(0)U(z) + h(1)z^{-1}U(z) + h(2)z^{-2}U(z) + h(3)z^{-3}U(z) + h(4)z^{-4}U(z) + \dots + h(N-1)z^{-(N-1)}U(z) = h(0)U(z) + z^{-1}(h(1)U(z) + z^{-1}(h(2)U(z) + z^{-1}(h(3)U(z) + \dots - \dots + z^{-1}(h(N-2)U(z) + z^{-1}h(N-1)U(z))))$$

This last equation, proceeding from right to left, performs the following in the time domain:

- Multiply u(n) by h(N-1), giving h(N-1)u(n)
- Delay by 1 unit, giving h(N-1)u(n-1)
- Add to h(N-2)u(n), giving h(N-2)u(n)+h(N-1)u(n-1)
- Delay by 1 unit, giving h(N-2)u(n-1)+h(N-1)u(n-2)
- Add to h(N-3)u(n), giving h(N-3)u(n)+h(N-2)u(n-1)+h(N-1)u(n-2)
- Delayby1 unit, giving h(N-3)u(n-1)+h(N-2)u(n-2)+h(N-1)u(n-3)
- ...
- Add to h(0)u(n)

all of which yields the implementation of the difference equation

 $v(n) = h(0)u(n) + h(1)u(n-1) + h(2)u(n-2) + \dots$

...+
$$h(N-2)u(n-N-2)+h(N-1)u(n-N-1)$$

Further, the equation

$$V(z) = h(0)U(z) + z^{-1}(h(1)U(z) + z^{-1}(h(2)U(z) + z^{-1}(h(3)U(z) + \dots))$$

... +
$$z^{-1}(h(N-2)U(z)+z^{-1}h(N-1)U(z))))$$

also suggests the transpose structure previously developed according to the rules. End of Aside End of Digression **Resumption of Implementation of Interpolator** Using the transposed form of the FIR filter the structure of the interpolator appears as below:



We may use identity #4 and the result of Example 6.2 to move the up-sampler to the right of the multipliers as shown below. As a result the number of multiplications is reduced by

$$\frac{100(L-1)}{L}\%.$$



Polyphase structures

The polyphase structure for FIR Filters was developed for the efficient implementation of sampling rate converters; however, it can be used in other applications. Further, the polyphase structure can be developed for any filter, FIR or IIR. We give below an introduction.

Polyphase Structure for FIR Filters The impulse response of the FIR filter h(n) is of finite length, *N*. The system function with *N* coefficients is

$$H(z) = \sum_{n=0}^{\infty} h(n) z^{-n} = h(0) + h(1) z^{-1} + h(2) z^{-2} + h(3) z^{-3} + h(4) z^{-4} + \dots + h(N-1) z^{-(N-1)}$$

We shall use another parameter *M*: we shall divide the number of coefficients into *M* groups (or branches or phases), modulo *M*. In other words the *N* terms in H(z) are arranged into *M* branches with each branch containing at most (Int (N-1/M)+1) terms.

Type 1 polyphase decomposition For illustration, let N = 11 and M = 2 so that one group contains 6 coefficients and the other 5 as developed below:

$$H(z) = \sum_{n=0}^{10} h(n) \ z^{-n} = h(0) + h(1) \ z^{-1} + h(2) \ z^{-2} + h(3) \ z^{-3} + h(4) \ z^{-4} + \dots + h(10) \ z^{-10}$$

= $h(0) + h(2)z^{-2} + h(4)z^{-4} + h(6)z^{-6} + h(8)z^{-8} + h(10)z^{-10} > 1^{\text{st}} \text{ group}$
+ $h(1)z^{-1} + h(3)z^{-3} + h(5)z^{-5} + h(7)z^{-7} + h(9)z^{-9} > 2^{\text{nd}} \text{ group}$
= $h(0) + h(2) \ z^{-2} + h(4) \ z^{-4} + h(6) \ z^{-6} + h(8) \ z^{-8} + h(10) \ z^{-10}$
+ $z^{-1} \{ h(1) + h(3) \ z^{-2} + h(5) \ z^{-4} + h(7) \ z^{-6} + h(9) \ z^{-8} \}$

Define

$$Int\left(\overline{N-1/M}\right) = \text{ integer part of the argument}$$

$$P_{0}(z) = h(0) + h(2)z^{-1} + h(4)z^{-2} + h(6)z^{-3} + h(8)z^{-4} + h(10)z^{-5} = \sum_{n=0}^{lnt(1-1/2)} \sum_{n=0}^{lnt(1-1/2)} h(2n+0)z^{-n}$$
(More generally, $P_{0}(z) = \sum_{n=0}^{n=0} h(Mn+0)z^{-n}$)
$$P_{1}(z) = h(1) + h(3)z^{-1} + h(5)z^{-2} + h(7)z^{-3} + h(9)z^{-4} = \sum_{n=0}^{lnt(1-1/2)} h(2n+1)z^{-n} \quad (h(11) = 0)$$
(More generally, $P_{1}(z) = \sum_{n=0}^{lnt(N-1/M)} h(Mn+1)z^{-n}$)

1

In this specific case we have

$$\frac{Y(z)}{X(z)} = H(z) = P_0(z^2) + z^{-1}P(z^2) = \sum_{n=0} z^{-m} P(z^2)$$

This decomposition of H(z) is known as **type 1 polyphase decomposition**. The corresponding structure is shown below. The functions $P_0(z^2)$ and $P_1(z^2)$ can each be implemented as a direct form.



By observing the expressions for $P_0(z)$ and $P_1(z)$ we can further generalize the functions $P_m(z)$ for any *m* as

$$P_m(z) = \sum_{n=0}^{lnt(N-1/M)} h(Mn+m)z^{-n}, \quad m = 0 \text{ to } (M-1)$$

Further, the system function becomes

$$\frac{Y(z)}{X(z)} = H(z) = P_0(z^M) + z^{-1}P_1(z^M) + \dots + z^{-(M-1)}P_{M-1}(z^M) = \sum_{m=0}^{M-1} z^{-m}P_1(z^M)$$

Example 5.8.1 As another example, let N = 11 and M = 3.

$$H(z) = \sum_{n=0}^{\infty} h(n) \ z^{-n} = h(0) + h(1) \ z^{-1} + h(2) \ z^{-2} + h(3) \ z^{-3} + h(4) \ z^{-4} + \dots + h(10) \ z^{-10}$$

= $h(0) + h(3)z^{-3} + h(6)z^{-6} + h(9)z^{-9} - 1^{\text{st}} \operatorname{group}$
+ $h(1)z^{-1} + h(4) \ z^{-4} + h(7) \ z^{-7} + h(10)z^{-10} - 2^{\text{nd}} \ \text{group}$
+ $h(2) \ z^{-2} + h(5) \ z^{-5} + h(8) \ z^{-8} - 3^{\text{rd}} \ \text{group}$
 $H(z) = h(0) + h(3) \ z^{-3} + h(6) \ z^{-6} + h(9) \ z^{-9}$
+ $z^{-1} \{ \ h(1) + h(4) \ z^{-3} + h(7) \ z^{-6} + h(10) \ z^{-9} \}$
+ $z^{-2} \{ \ h(2) + h(5) \ z^{-3} + h(8) \ z^{-6} \}$

Define

$$Int\left(\overline{N^{-1}/M}\right) = \text{ integer part of the argument}$$

$$P_{0}(z) = h(0) + h(3)z^{-1} + h(6)z^{-2} + h(9)z^{-3} = \sum_{n=0}^{\infty} h(3n+0)z^{-n}$$

$$Int\left(\overline{N^{-1}/M}\right)$$
(More generally, $P_{0}(z) = \sum_{n=0}^{\infty} h(Mn+0)z^{-n}$)
$$P_{1}(z) = h(1) + h(4)z^{-1} + h(7)z^{-2} + h(10)z^{-3} = \sum_{n=0}^{\ln t} (11-1/3)$$
(More generally, $P_{1}(z) = \sum_{n=0}^{\infty} h(Mn+1)z^{-n}$)

$$P_{2}(z) = h(2) + h(5)z^{-1} + h(8)z^{-2} = \sum_{\substack{n=0 \\ h(3n+1)z^{-n} \\ h(m=1)M}}^{Int (m=1/3)} (h(11) = 0)$$
(More generally, $P_{2}(z) = \sum_{\substack{n=0 \\ h(m=1)M}}^{Int (m=1)M} h(Mn+1)z^{-n})$
In this specific case we have
$$Y(z) = H(z) = P(z^{3}) + z^{-1}P(z^{3}) + z^{-2}P(z^{3}) = \sum_{\substack{n=0 \\ m=0}}^{2} z^{-m} P(z^{3})$$



Generalization As mentioned earlier, in the general case for an arbitrary $M (\leq N)$ we have $I_{M}(N = 1/M)$

$$P_m(z) = \sum_{n=0}^{\infty} h(Mn+m) z^{-n}, \quad m = 0 \text{ to } (M-1)$$

and

$$\frac{Y(z)}{X(z)} = H(z) = P_0(z^M) + z^{-1}P_1(z^M) + \dots + z^{-(M-1)}P_{M-1}(z^M) = \sum_{m=0}^{M-1} z^{-m}P_m(z^M)$$

This overall operation is known as polyphase filtering.



Type 2 polyphase decomposition Given

$$\frac{Y(z)}{X(z)} = H(z) = P_0(z^M) + z^{-1}P_1(z^M) + \dots + z^{-(M-1)}P_{M-1}(z^M) = \sum_{m=0}^{M-1} z^{-m}P_n(z^M)$$

set m = M - 1 - k: as *m* goes from 0 to M - 1, *k* goes from M - 1 to 0. Thus

$$H(z) = \sum_{k=M-1}^{0} z^{-(M-1-k)} P_{M-1-k}(z)$$

Let

$$Q_k(z^M) = z^{-(M-1-k)} P_{M-1-k}(z^M)$$

This gives us the type 2 polyphase decomposition

$$\sum_{k=M-1}^{0} \sum_{k=0}^{M-1} Q_{k} (z^{M}) = Q_{k}(z^{M}) + Q_{1}(z^{M}) + \dots + Q_{M-1}(z^{M})$$

Type 3 polyphase decomposition Given

$$\frac{Y(z)}{X(z)} = H(z) = P_0(z^M) + z^{-1}P_1(z^M) + \dots + z^{-(M-1)}P_{M-1}(z^M) = \sum_{m=0}^{M-1} z^{-m}P_m(z^M)$$

set m = -k: as *m* goes from 0 to *M*-1, *k* goes from 0 to -(M-1) to 0. Thus

$$H(z) = \sum_{k=0}^{-(M-1)} z^{k} P_{-k}(z^{M})$$

Let

$$R_{k}(z^{M}) = z^{k} P_{-k}(z^{M})$$

Example 5.8.2 For the system

$$H(z) = \sum_{n=0}^{\infty} h(n) z^{-n} = h(0) + h(1)z^{-1} + h(2)z^{-2} + h(3)z^{-3} + h(4)z^{-4} + h(5)z^{-5} + h(6)z^{-6}$$

= { $h_0 + h_2 z^{-2} + h_4 z^{-4} + h_6 z$ }⁶ + { $h_1 z^{-1} + h_3 z^{-3} + h_5 z$ }⁻⁵
= { $h_0 + h_2 z^{-2} + h_4 z^{-4} + h_6 z$ }⁶ + z^{-1} { $h_1 + h_3 z^{-2} + h_5 z^{-4}$ }
 $P_0(z) = h_0 + h_2 z^{-1} + h_4 z^{-2} + h_6 z^{-3}$

 $P_{I}(z) = h_{1} + h_{3} z^{-1} + h_{5} z^{-2}$ The structure is shown below (left). As mentioned earlier $P_{0}(z^{2})$ and $P_{1}(z^{2})$ can each be implemented as a direct form.

Polyphase realization

Canonical polyphase realization



The structure shown on the right is obtained by moving the delay element z^{-1} to the right of $P_1(z^2)$ these two being in series in the second phase. In this latter case the two systems $P_1(z^2)$ and $P_1(z^2)$ can share the same delay elements (that is, storage locations) even though each has its own set of coefficients, thus resulting in a canonical polyphase realization, shown below.



Polyphase Structure for IIR Filters The anti-aliasing filter in a decimator and the anti-imaging filter in an interpolator may each be either an FIR or an IIR filter. The polyphase structure can be developed for any filter, FIR or IIR, and any finite value of M. We now proceed to the case where h(n) is an infinitely long sequence:

$$H(z) = \sum_{n = -\infty} h(n) z^{-n} = \dots + h(-M) z^{M} + \dots h(-2) z^{2} + h(-1) z^{1}$$

+ $h(0) + h(1) z^{-1} + h(2) z^{-2} + \dots + h(M-1) z^{-(M-1)}$
+ $h(M) z^{-M} + h(M+1) z^{-(M+1)} \dots$

Once again we arrange the terms into M groups or branches in a modulo-M fashion. Each branch contains infinitely many terms. The terms are arranged in tabular form below:

$H(z) = \sum_{n = -\infty}^{\infty} h(n) z^{-n}$					
1 st branch<		$+ h(-M) z^{M}$	+ h(0)	$+ h(M) z^{-M}$	
2 nd branch		$+h(-M+1) z^{M-1}$	$+ h(1) z^{-1}$	$+ h(M+1) z^{-(M+1)}$	•••
			$+ h(2) z^{-2}$	$+h(M+2) z^{-(M+2)}$	•••
					•••
i th branch<		$h(-M+i-1)z^{-(-M+i-1)}$	$h(i-1)z^{-(i-1)}$	$h(M+i-1)z^{-(M+i-1)}$	•••
		•••	•••		
		$+h(-2)z^{2}$	$+h(M-2)z^{-(M-2)}$	$+h(2M-2)z^{-(2M-2)}$	
M th branch<		$+h(-1) z^{1}$	$+h(M-1)z^{-(M-1)}$	$+h(2M-1)z^{-(2M-1)}$	

$$H(z) = [\dots + h(-M)z^{M} + h(0) + h(M)z^{-M} + h(2M)z^{-2M} + \dots] \qquad 1^{\text{st}} \text{ row} \\ + [\dots + h(-M+1)z^{M-1} + h(1)z^{-1} + h(M+1)z^{-(M+1)} + \dots] \qquad 2^{\text{nd}} \text{ row} \\ + [\dots] + \dots \\ + [\dots + h(-M+i-1)z^{-(-M+i-1)} + h(i-1)z^{-(i-1)} + h(M+i-1)z^{-(M+i-1)} \\ + h(2M+i-1)z^{-(2M+i-1)} + \dots] \\ \dots \\ + [\dots + h(-2)z^{2} + h(M-2)z^{-(M-2)} + h(2M-2)z^{-(2M-2)} + \dots]$$

+[...+
$$h(-1)z^{1}+h(M-1)z^{-(M-1)}+h(2M-1)z^{-(2M-1)}+...] Mth row$$

We factor out z^{0} from the first row (branch), z^{-1} from the 2^{nd} row, and, in general, $z^{-(i-1)}$ from the i^{th} row to get

$$\begin{split} H(z) &= [\dots + h(-M)z^{M} + h(0) + h(M)z^{-M} + h(2M)z^{-2M} + \dots] & 1^{\text{st}} \text{ row} \\ &+ z^{-1} [\dots + h(-M+1)z^{M} + h(1) + h(M+1)z^{-M} + \dots] & 2^{\text{nd}} \text{ row} \\ &+ z^{-2} [\dots] + \dots \\ &+ z^{-(i-1)} [\dots + h(-M+i-1)z^{M} + h(i-1) + h(M+i-1)z^{-M} + h(2M+i-1)z^{-2M} + \dots] \\ &\dots \\ &+ z^{-(M-2)} [\dots + h(-2) z^{M} + h(M-2) + h(2M-2)z^{-M} + \dots] \\ &+ z^{-(M-1)} [\dots + h(-1) z^{M} + h(M-1) + h(2M-1) z^{-M} + \dots] \end{split}$$

Define

$$P_0(z^M) = [\dots + h(-M)z^M + h(0) + h(M)z^{-M} + h(2M)z^{-2M} + \dots] = \sum_{n = -\infty} h(nM)z^{-nM}$$

œ

so that

$$P_0(z) = [\dots + h(-M)z + h(0) + h(M)z^{-1} + h(2M)z^{-2} + \dots] = \sum_{n = -\infty}^{\infty} h(nM) z^{-n}$$

Similarly,

$$P^{l}(z) = [\dots + h(-M+1)z + h(1) + h(M+1)z^{-1} + h(2M+1)z^{-2} + \dots] = \sum_{\substack{n = -\infty \\ n = -\infty}}^{\infty} h(nM+1) z^{-n}$$

In general

$$P_m(z) = \sum_{n = -\infty}^{\infty} h(nM + m) z^{-n}, \qquad m = 0 \text{ to } (M - 1)$$

And the system function can now be written M^{-1}

$$H(z) = \sum_{m=0}^{M-1} z^{-m} P_m(z^{M})$$

This is called the *M*-component polyphase decomposition of H(z). The *M* functions $P_m(z)$ are the polyphase components of H(z). This overall operation is known as polyphase filtering. As an example, for M = 3, we have

example, for
$$M = 3$$
, we have

$$\frac{Y(z)}{X(z)} = H(z) = \sum_{m=0}^{2} z^{-m} P_{m}(z^{3}) = P_{0}(z^{3}) + z^{-1} P(z^{3}) + z^{-2} P(z^{3})$$

$$Y(z) = P_{0}(z^{3}) X(z) + z^{-1} P(z^{3}) X(z) + z^{-2} P(z^{3}) X(z)$$

This last equation leads to the structure below (left):



We may also rearrange the output equation as

$$Y(z) = P_{0}(z^{3}) X(z) + z^{-1} P(z^{3}) X(z) + z^{-2} P(z^{3}) X(z)$$

= $P_{0}(z^{3}) X(z) + z^{-1} \{P(z^{3}) X(z) + z^{-1} P(z^{3}) X(z)\}$

This last equation leads to the polyphase structure shown above (right), known as the transpose polyphase structure because it is similar to the transpose FIR filter realization.

Polyphase structure for a decimator

The decimator block diagram is shown below: it consists of an anti-aliasing filter, H(z), which could be an FIR or an IIR filter, followed by an *M*-fold down sampler.



We replace the filter H(z) by its *M*-component polyphase decomposition

$$H(z) = \sum_{m=0}^{M-1} z^{-m} P_m(z^{M})$$

The sub filters $P_0(\frac{M}{z})$, $P_1(\frac{M}{z})$, ..., $P_{M-1}(\frac{M}{z})$ could be FIR or IIR depending on H(z). The block diagram then appears as below.



We may use identity #1 to move the down-sampler to the immediate right of $P_m(z^M)$ in each branch, and then use identity #3 to move the down-sampler from the immediate right to the immediate left of $P_m(.)$ while at the same time changing $P_m(\frac{M}{z})$ to $P_m(z)$. The result appears as below. It can be seen from this diagram that in this structure the number of multiplications is reduced by a factor of M.



**** Continuing with the development, comparing $P \sum_{n=1}^{\infty} z_n = {}^{\infty} h(rM + m) z_n$

with the defining equation of the z-transform

$$P_{m}(z) = \sum_{r=-\infty} p_{m}(r) z_{-r}$$

we identify

$$p_m(r) = h(rM + m)$$

Note that M is a constant and m is a parameter. Upon substituting for $P_m(z)$ in the system function

$$H(z) = \sum_{m=0}^{M-1} z^{-m} P_m(z^{M})$$

we have

$$H(z) = \sum_{m=0}^{M-1} z^{-m} \sum_{r=-\infty}^{\infty} h(rM+m) (z^{M})^{-r} = \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} h(rM+m) z^{-(rM+m)}$$
$$= \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} p_{m}(r) z^{-(rM+m)}$$

The output transform is

$$Y(z) = H(z)X(z) = \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} p_m \text{ (r) } X(z) z^{-(rM+m)}$$

The output is

$$y(n) = \mathbf{z}^{-1} \left\{ \sum_{\substack{m=0 \ r=-\infty \\ m=0 \ r=-\infty}}^{M-1} \sum_{\substack{m=0 \ r=-\infty \\ m=0 \ r=-\infty}}^{\infty} p_m(r) X(z) z_{-(rM+m)} \right\} = \sum_{\substack{m=0 \ r=-\infty \\ m=0 \ r=-\infty}}^{M-1} \sum_{\substack{m=0 \ r=-\infty}}^{\infty} p_m(r) Z_{-1} \left\{ X(z) z_{-(rM+m)} \right\}$$

Define

$$x_m(r) = x(rM - m)$$

Note that M is a constant and m is a parameter.

$$x_m(-r) = x(-rM - m)$$

Shifting the sequence by n units we get

$$x_m(n-r) = x(n-rM-m)$$

The output now may be written

$$y(n) = \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} p_m(r) x_m(n-r)$$

Define the operation of polyphase convolution as

$$y(m) = p_m(n)^* x_m(n) = \sum_{r=-\infty} p_m(r) x_m(n-r)$$

Then

$$y(n) = \sum_{m=0}^{M-1} p_m(n) * x_m(n) = \sum_{m=0}^{M-1} y_m(n)$$

This overall operation is known as polyphase filtering.

FINITE WORD LENGTH EFFECTS

Finite Wordlength Effects

All the signals and systems are digital in DSP. The digital implementation has finite accuracy. When numbers are represented in digital form, errors are introduced due to their finite accuracy. These errors generate finite precision effects or finite wordlength effects.

.....(1)

Let us consider an example if the first order IIR filter to illustrate how errors are encountered in discretization. Such filter can be described as,

$$y(n) = \alpha y(n-1) + x(n)$$

The z-transform of above equation gives

 $z[y(n)] = \alpha z[y(n-1)+x(n)]$

 $X(z) = \alpha_{Z^{-1}} Y(z) + X(z)$

Hence the transfer function

$$H(z) = \frac{1}{1 - \alpha z^{-1}} = \frac{z}{z - \alpha}$$

Here observe that ' α ' is the filter coefficient when this filter is implemented on some DSP processor or software, ' α ' can have only discrete values. Let the discrete values of α be represented by α . Hence the actual transfer function which is implemented is given as,

$$H\left(z\right)=\frac{z}{z-\alpha}$$

The transfer function given by above equation is slightly different from H(z). Hence the actual frequency response will be different from desired response.

The input x(n) is obtained by sampling the analog input signal. Since the quantizer takes only fixed(discrete) values of x(n), error is introduced. The actual input can be denoted by x(n).

x(n) = x(n) + e(n)

Here e(n) is the error introduced during A/D conversion process due to finite wordlength of the quantizer. Similarly error is introduced in the multiplication of α and y(n-1) in equation(1). This is because the product α y(n-1) has to be quantized to one of the available discrete values. This introduces error. These errors generate finite wordlength effects.

Finite Wordlength Effects in IIR Digital Filters

When an IIR filter is implemented in a small system, such as an 8-bit microcomputer, errors arise in representing the filter coefficients and in performing the arithmetic operations indicated by the difference equation. These errors degrade the performance of the filter and in extreme cases lead to instability.

Before implementing an IIR filter, it is important to ascertain the extent to which its performance will be degraded by finite wordlength effects and to find a remedy if the degradation is not acceptable. The effects of these errors can be reduced to acceptable levels by using more bits but this may be at the expense of increased cost.

The main errors in digital IIR filters are:

i ADC Quantization Noise:

This noise is caused by representing the samples of the input data ,by only a small number of bits.

ii. Coefficient quantization errors:

These errors are caused by representing the IIR filter coefficients by a finite number of bits.

iii. Overflow errors

These errors are caused by the additions or accumulation of partial results in a limited register length.

iv. Product round-offerrors

These errors are caused when the output, and results of internal arithmetic operations are rounded to the permissible wordlength.

Finite Wordlength Effects in FFT Filters

As in most DSP algorithms, the main errors arising from implementing FFT algorithms using fixed point arithmetic are

i. Round off errors

These errors are produced when the product $W^k B$ is truncated or rounded to the system wordlength.

ii. Overflow errors

These errors result when the output of a butterfly exceeds the permissible wordlength.

iii. Coefficient quantization errors

These errors result from representing the twiddle factors using a limited number of bits.

LIMIT CYCLES

5.2.1 Overflow oscillations

Limit Cycle: The finite wordlength effects are analyzed using the linear model of the digital systems. But nonlinearities are introduced because of quantization of arithmetic operations. Because of these nonlinearities, the stable digital filter under infinite precision may become unstable under finite preision.Because of this instability, oscillating periodic output is generated.Such output is called limit cycle.The limit cycle occur in IIR filters due to feedback paths.

Types of Limit Cycles

There are two types of limit cycles.

- (1) Granular and
- (2) Overflow.

1. Granular Limit Cycles

The granular limit cycles are of low amplitude. These cycles occur in digital filters when the input signal levels are very low. The granular limit cycles are of two types. They are

- i. Inaccessible limit cycles
- ii. Accessible limit cycles.

2. Overflow Limit Cycles

Overflow limit cycles occur because of overflow due to addition in digital filters implemented with finite precision. The amplitudes of overflow limit cycles are very large and it can cover complete dynamic range of the register. This further leads to overflow causing cumulative effect. Hence overflow limit cycles are more serious than granular limit cycles.

Transfer Characteristics and Example



Figure: Transfer characteristics of adder having overflow limit cycles

Because of overflow limit cycle oscillations the output fluctuates between minimum and maximum values. The above figure shows the transfer characteristic of an adder that exihibit overflow limit cycle oscillations. Here f(v) indicates addition operation. Consider the addition of following two numbers in sign magnitude form.

$$x_{1} = 0.111 \text{ i.e., } \frac{7}{8}$$

$$x_{2} = 0.101 \text{ i.e., } \frac{5}{8}$$

$$x_{2} = 0.101 \text{ i.e., } \frac{5}{8}$$

Then

1 - 2Here overflow has occured in addition due to finite precision and the digit before decimal point makes the number negative.

Signal Scaling:

Need for Scaling: Limit cycle oscillations can be avoided by using the nonlinear transfer characteristic. But it introduces distortion in the output. Hence it is better to perform signal scaling such that overflow or underflow does not occur and hence limit cycle oscillations can be avoided.

Implementation of Signal Scaling

Figure shows the direct form-II structure of IIR filter. Let the input x(n) be scaled by a factor s_0 before the summing node to prevent overflow. With the scaling, the transfer function will be,



Figure: Direct form-II realization for second order IIR filter

. .

Let us define the transfer function

$$H'(z) = \frac{W(z)}{X(z)}$$

From figure we can write above transfer function as,

$$H'(z) = \frac{s_0}{1 + a_1 z^{-1} + a_2 z^{-2}}$$

 $A(z) = 1 + a_1 z^{-1} + a_2 z^{-2}$

Since

 $H'(z) = \frac{S_0}{A(z)}$ $H'(z) = \frac{W(z)}{W(z)}$

Since,

$$\frac{W(z)}{X(z)} = \frac{s_0}{A(z)}$$
$$W(z) = \frac{s_0 X(z)}{A(z)}$$

Let

(or)

 $S(z) = \frac{1}{A(z)^{2}}, \text{ then above equation becomes,}$ $W(z) = s_0 S(z) X(z)$

Evaluating z-transform on unit circle, we put $z = e^{j\omega}$ in above equation,

$$W(e^{j\omega}) = s \delta(e^{j\omega}) X(e^{j\omega})$$

Taking inverse Fourier transform of above equation,

$$\omega(\mathbf{n}) = \frac{1}{2\pi} \int W(\mathbf{e}^{j\omega}) \mathbf{e}^{j\omega \mathbf{n}} d\omega \qquad = \frac{1}{2\pi} \int S S_0(\mathbf{e}^{j\omega}) X(\mathbf{e}^{j\omega}) \mathbf{e}^{j\omega \mathbf{n}} d\omega$$

Hence $\omega^2(\mathbf{n}) = \frac{s_0^2}{4\pi^2} \left| \int S(e^{j\omega}) X(e^{j\omega}) e^{j\omega |\mathbf{n}|} d\omega \right|^2$

Schwartz inequality states that,

$$\left|\int \mathbf{x}_{1}(t)\mathbf{x}_{2}(t)dt\right|^{2} \leq \left|\mathbf{x}_{1}(t)^{2}dt\int |\mathbf{x}_{2}(t)|^{2}dt\right|$$

Using this relation we can write above equation as, $\frac{1}{2} = \frac{1}{100} \frac{1}{2} \frac{1}{100} \frac{1}{2} \frac{1}{100} \frac{1}{2} \frac{1}{100} \frac{1}{100$

$$\omega^{2}(n) \leq \frac{s}{4\pi^{2}} |S(e)| d\omega \int X(e) d|_{1}$$

since $e^{j\omega n} = |1$

$$\frac{|a|}{2} \int_{2}^{j} \frac{1}{2} \int_{1}^{j} \frac{1}{2} \int_{1}^{j} \frac{1}{2\pi} \int_$$

Parseval's theorem states that $\frac{1}{2\pi} \int |X(e_{j\omega})|^2 = \sum_{n=0}^{\infty} x^2(n)$. Then above equation can be written as,

$$\omega^{2}(\mathbf{n}) \leq s_{0}^{2} \left[\frac{1}{2\pi} \int \left| S(e^{j\omega})^{2} d\omega \sum_{\mathbf{n}=0}^{\infty} x^{2}(\mathbf{n}) \right| \right]$$
We have put $z = e^{j\omega}$. Hence $dz = je^{j\omega}d\omega$ or

$$d\omega = \frac{dz}{z} = \frac{dz}{z}$$
, since
 $je^{j\omega}$ jz $e^{j\omega} = z$

Putting these values in equation

$$\omega^{2}(\mathbf{n}) \leq s_{0}^{2} \left| \frac{1}{2\pi 1} \int_{1}^{2} |S(e^{j\omega})|^{2} \cdot \frac{dz}{\sqrt{3}} \sum_{n=0}^{\infty} x^{2}(\mathbf{n}) \right|$$

$$\leq s^{2} \leq x^{2}(\mathbf{n}) \cdot \frac{|S(z)_{2}|^{2}}{\sqrt{3}} \int_{1}^{\infty} \frac{|S(z)_{2}|^{2}}{\sqrt{3}} \frac{|S(z)$$

Here $|S(z)|^2 = S(z)S(z^{-1})$. Then we have,

$$\omega^{2}(\mathbf{n}) \leq s_{0}^{2} \sum_{n=0}^{\infty} x^{2}(\mathbf{n}) \frac{1}{2\pi j} S(z) S(z^{-1}) z^{-1} dz$$

 $\omega(x) \leq \sum_{n} (x^{2} n s^{2} \int_{0}^{\infty} (x^{2} x^{2} - s^{2} \int_{0}^{\infty} (x^{2} - s^{2} f^{2}) \int_{0}^{\infty} (x^{2} - s^{2} f^{2}) \int_{0}^{\infty} (x^{2} - s^{2}) \int_{0}^{\infty} (x^{2} - s^{2$

Here the integration is executed over a closed contour i.e. $\omega^2(n) \le s^{2z} \sum x^2(n) \cdot \frac{1}{2\pi j^{\circ}} S(z) \cdot S(z^{-1}) z^{-1} dz$

(or)

Here $\omega^2(n)$ represents instantaneous energy of signal after first summing node. And $x^2(n)$ represents instantaneous energy of input signal. Overflow will not occur if

$$\omega^{2}(\mathbf{n}) \leq \sum_{n=0}^{\infty} \mathbf{x}^{2}(\mathbf{n})$$

For this equation to be true we get following condition from equation
$$1 \quad \mathbf{S}(\mathbf{z})\mathbf{S}(\mathbf{z}^{-1})\mathbf{z}^{-1}d\mathbf{z} = 1$$
$$s_{0}^{2} 2\pi \mathbf{j} \int_{\mathbf{C}}^{\mathbf{C}} \mathbf{E} = \frac{1}{\mathbf{A}(\mathbf{z})}$$

Earlier we have defined
$$\mathbf{S}_{0}^{\mathbf{Z}} = \frac{1}{\mathbf{A}(\mathbf{z})}$$
. Hence above condition becomes,
$$\frac{2s_{0}^{1}}{2}\pi \mathbf{j}_{\mathbf{C}}^{1} \mathbf{C} \frac{\mathbf{C}_{0}^{\mathbf{Z}^{-1}d\mathbf{z}}}{\mathbf{C}\mathbf{A}\mathbf{z}\cdot\mathbf{A}\mathbf{z}} = 1$$
$$s_{0}^{2} = \frac{1}{2}\pi \mathbf{j}_{\mathbf{C}}^{1} \mathbf{C} \frac{\mathbf{C}_{0}^{\mathbf{Z}^{-1}d\mathbf{z}}}{2\pi \mathbf{j}} = 1$$

(or)

Above equation gives the value of scaling factor s_0 to avoid overflow

5.3 ROUND OFF NOISE IN IIR DIGITAL FILTERS

Statistical Model for Analysis of Round-off Error Multiplication:

We perform arithmetic operations like addition and multiplication some errors will be occured. Those errors are called arithmetic errors. The results of arithmetic operations are required to be quantized so that they can occupy one of the finite set of digital levels. Such operation can be visualized as multiplier (or other arithmetic operation) with quantizer at its output.



Figure: Quantization of multiplication or product

The above process can be represented by a statistical model for error analysis. The output v (n) and error $e_{\alpha}(n)$ in product quantization process.i.e.,

$$\upsilon(n) = \upsilon(n) + e_{\alpha}(n)$$

The statistical model is shown below.



Figure: Statistical model for analysis of round-off error multiplication

For the analysis purpose following assumptions are made.

-) The error sequence $\{e_{\alpha_{(n)}}\}$ is the sample sequence of a stationary white noise process.
- ii $e_{\alpha}(n)$ is having uniform distribution over the range of quantization error.
- iii) The sequence $\{e_{\alpha}(n)\}$ is uncorrelated with the sequence $\upsilon(n)$ and input sequence x(n).

Computational output round off noise

Product Round-off Errors and its Reduction:

The results of product or multiplication operations are quantized to fit into the finite wordlength, when the digital filters are implemented using fixed point arithmetic. Hence errors generated in such operation are called product round off errors.

The effect of product Round-off errors can be analyzed using the statistical model of the quantization process. The noise due to product round-off errors reduces the signal to noise ratio at the output of the filter. Some times this ratio may be reduces below acceptable levels. Hence it is necessary to reduce the effects of product round-off errors.

There are two solutions available to reduce product round-off errors.

a) Error feedback structures and

b)State space structure.

The error feedback structures use the difference between the unquantized and quantized signal to reduce the round-off noise. The difference between unquantized and quantized signal is fed back to the digital filter structure in such a way that output noise power due to round-off errors is reduced.

First Order Error-feedback Structure to reduce Round-off Error:

The results of product or multiplication operations are quantized to fit into the finite wordlength, when the digital filters are implemented using fixed point arithmetic. Hence errors generated in such operation are called product round off errors.

The effect of product round-off errors can be analyzed using the statistical model of the quantization process.

Let the quantization error signal be given as the difference between unquantized signal y(n) and quantized signal v(n) .i.e.,

e(n) = y(n) - v(n)



Figure: First order error feedback structure to reduce product round-off error

This error signal is fed back in the structure such that round-off noise is reduced. Such structure for first order digital filter is shown in figure.

The incorporation of quantization error feedback as shown in figure helps in reducing the noise power at the output . This statement can be proved mathematically.

Round-off Errors in FFT Algorithms:

FFT is used in large number of applications. Hence it is necessary to analyze the effects due to finite wordlengths in FFT algorithms. The most critical error in FFT computation occurs due to arithmetic round-off errors.

The DFT is used in large number of applications such as filtering, correlation, spectrum analysis etc. In such applications DFT is computed with the help of FFT algorithms. Therefore it is important to study the quantization errors in FFT algorithms. These quantization effects mainly take place because of round-off errors. These errors are introduced when multiplications are performed in fixed point arithmetic.

FFT algorithms require less number of multiplications compared to direct computation of DFT. But it does not mean that quantization errors are also reduced in FFT algorithms.

Let σ_{x}^{2} represents the variance of output DFT coefficients i.e., |X(k)|.

For N-point DFT σ_x is given as,

For direct computation of DFT, the variance of quantization errors in multiplications is given as,

$$\sigma_q^2 = \frac{N}{3} \cdot \Delta^2 \qquad \dots \dots (2)$$

Here σ_{a}^{2} is variance of quantization errors and Δ is step size which is given as,

$$\Delta = 2^{-b}$$
(3)

And b is the number of bits to represent one level. Hence equation (2) becomes,

$$\sigma_q^2 = \frac{N}{3} \cdot 2^{-2b}$$
(4)

The signal to noise power ratio at the output (i.e., DFT coefficients) can be considered as the measure of quantization errors. This ratio is the ratio of variance of DFT coefficients (σ^2) to the variance of quantization errors (σ^2) i.e.,

Signal to noise ratio in direct computation of $DFT = \begin{pmatrix} \sigma^2 \\ \sigma^2 \\ \sigma^2 \\ \sigma^2 \\ Direct DFT \end{pmatrix}$

From equation (1) and equation (2) we have,

31 T

$$\begin{pmatrix} \sigma^2 \end{pmatrix} = \underbrace{\underline{SN}}_{2b} \qquad 2b \\ \underbrace{-}_{2 \atop q} \\ \left[\begin{array}{c} \sigma_q \end{array} \right]_{\text{Direct DFT}} \\ \mathbf{N}_{3} \\ \cdot 2^{-2b} \\ \cdot 2^{-2b$$

When DFT is computed using FFT algorithms, the variance of the signal remains same i.e., $\frac{1}{2}$

$$\sigma_x = \frac{1}{3N}$$
 from equation (1)(6)

But with algorithms the variance of the quantization errors is given as,

$$\sigma_q^2 = \frac{2}{3} \cdot 2^{-2b} \qquad \dots \dots (7)$$

Hence signal to noise ration in FFT algorithms is,

$$\begin{array}{c} \left(\sigma^{2}\right) = \frac{1}{3N} \\ \left(\gamma^{2}\right) = \frac{1}{3N} \\ \left(\gamma^{2}\right) = \frac{1}{3 \cdot 2^{2b}} 2N \end{array}$$

In the above expression, the signal to noise ration is inversely proportional to N. whereas in direct DFT computation the signal to noise ratio is inversely proportional to N^2 as given by equation (5). This means quantization errors increase fast with increase in 'N' in direct computation of DFT. But in FFT algorithms the quantization errors increase slowly with increase in 'N'.

Product of Round-off Errors in IIR Digital Filters:

The results of product or multiplication operations are quantized to fit into the finite wordlength, when the digital filters are implemented using fixed point arithmetic. Hence errors generated in such operation are called product round off errors.

Product round-off error analysis is an extensive topic.Our presentation here will be brief and aims to make you aware of the nature of the errors, their effects and how to reduce them if necessary.

The basic operations in IIR filtering are defined by the familiar second- order difference equation:

$$y(n) = \sum_{k=0}^{2} b_k x(n-k) - \sum_{k=1}^{2} a_k y(n-k)$$

Where x(n-k) and y(n-k) are the input and output data samples, and b_k and a_k are the filter coefficients. In practice these variables are often represented as fixed point numbers. Typically, each of the products b_k x(n-k)

and $a_k y(n-k)$ would require more bits to represent than any of the operands. For example, the product of a B-bit data and a B-bit coefficient is 2B bits long.

Truncation or rounding is used to quantize the products back to the permissible wordlength. Quantizing the products leads to errors, popularly known as round-off errors, in the output data and hence a reduction in the SNR. These errors can also lead to small-scale oscillations in the output of the digital filter, even when there is no input to the filter.



Figure: Representation of the product quantization error: (a) a block diagram representation of the quantization process; (b) a linear model of the quantization process

The figure(a) represents a block diagram of the product quantization process, and figure (b) represents a linear model of the effect of product quantization. The model consists of an ideal multiplier, witk infinite precision, in series with an adder fed by a noise sample, e(n), representing the error in the quantized product ,where we have assumed, for simplicity, that x(n), y(n), and K are each represented by B bits. Thus

 $\mathbf{y}(\mathbf{n}) = \mathbf{K}\mathbf{x}(\mathbf{n}) + \mathbf{e}(\mathbf{n})$

The noise power, due to each product quantization, is given by

$$\sigma_r^2 = \frac{q_2}{12}$$

Where r symbolizes the round-off error and q is the quantization step defined by the wordlength to which product is quantized. The round-off noise is assumed to be a random variable with zero mean and constant variance. Although this assumption may not always be valid, it is useful in assessing the performance of the filter.

Product of Round-off Errors on Filter Performance:

The effects of round-off errors on filter performance depend on the type of filter structure used and the point at which the results are quantized.

The above figure represents the quantization noise model for the direct form building block. It is assumed in the figure that the input data,x(n),output data,y(n),and the filter coefficients are represented as B-bit numbers (including the sign bit). The products are quantized back to B bits after multiplication by rounding (or truncation).



Figure: Product quantization noise model for the direct form filter section. All the noise sources in (a) have

been combined in (b) as they feed to the same point

Since all five noise sources, e_1 to e_5 in figure(a), feed to the same point (that is into the middle adder), the total output noise power is the sum of the individual noise powers(figure(b)).



Figure: Product quantization noise model for the canonic filter section. The noise sources feeding the same point in (a) have been combined in (b)

$$\sigma_{\rm or}^{2} = \frac{5q}{12} \left| \frac{1}{2 \prod j} \int_{c}^{c} F(z)F(z^{-1}) \frac{dz}{z} \right|_{s_{1}}^{2} = \frac{5q}{12} \left[\sum_{k=0}^{c} f(k) s \right]_{1}^{2} = \frac{5q^{2}}{12} \left\| F(z) \right\|_{2}^{2} s_{1}^{2}$$

Where F(z) = $\frac{1}{1 + a z^{-1} + a z^{-2}}$

 $f(k) = Z^{-1}[F(z)]$ is the inverse z-transform of F(z), which is also the impulse response from each noise source to the filter output, $\frac{2}{2}$ is the L_2 norm squared and $\frac{q^2}{12}$ is the intrinsic product round-off noise power. The total noise power at the filter output is the sum of the product round-off noise and the ADC quantization noise.

$$\sigma_{0}^{2} = \sigma_{0A^{2} + \sigma_{0}r^{2}}$$

$$= \frac{q^{2} \left[\sum_{k=0}^{\infty} h^{2}(k) + \frac{2^{\infty}}{5s_{1}\sum_{k=0}^{2} f^{2}(k)}\right]}{12 \left[\sum_{k=0}^{2} h^{2}(k) + \frac{2^{\infty}}{5s_{1}\sum_{k=0}^{2} f^{2}(k)}\right]} = \frac{q^{2}}{12} \left[\frac{H(z)}{2} + \frac{5s}{12}\right]$$

For canonic section, figure(a), the noise model again includes a scale factor as this generates a round-off error of its own. The noise sources $e_1(n)$ to $e_3(n)$ all feed to the left adder, whilst the noise sources $e_4(n)$ to $e_6(n)$ feed directly into the filter output. Combination of the noise sources feeding to the same point leads to the noise model of figure(b). Assuming uncorrelated noise sources, the total noise contribution is simply the sum of

the individual noise contribution
$$g = F(z) + 1$$

$$\sigma_{or}^{2} = \sum_{k=0}^{2 \sigma} f(k) + \frac{2}{2} = F(z) + 1$$

Where f(k) is the impulse response from the noise source e_1 to the filter output, and F(z) the corresponding transfer function given by

$$F(z) = s_1 \frac{\overset{0}{\overset{1}{_1}} + b z^{-1} + b z^{-2}}{\overset{1}{_1} + a z^{-1} + a z^{-2}} = s_1 H(z)$$

The total noise (ADC+round-off noises) at the filter output is given by

$$\sigma_{0}^{2} = \sigma_{0A}^{2} + \sigma_{or}^{2} + s + h + (k) = h + (k) = h + (k) = \frac{3}{12} \left\{ 3 \left[1 + s^{2} + H(z) + H(z) + \frac{1}{2} + \sum_{k=0}^{\infty} 2^{2} + \sum_{k=0}^{\infty} 2^{k} +$$

METHODS TO PREVENT OVERFLOW

Prevent Overflow Limit Cycle Oscillations: The overflow limit cycles occur because of overflow due to addition in digital filters implemented with finite precision. The amplitudes of overflow limit cycles are very large and it can cover complete dynamic range of the register.

The specific design of filter coefficients do not assure prevention of overflow limit cycle oscillations. The transfer characteristic can be modified to avoid overflow limit cycle oscillations.



Figure: Prevention of overflow limit cycle oscillations

As shown in figure when an overflow or underflow is sensed, the output of the adder is set to its full scale value of ± 1 . This prevents oscillatory output. This nonlinearity of the characteristic causes very small distortion in the output because overflow/underflow occurs rarely.

Scaling is also used to prevent overflow limit cycle oscillations. Limit cycle free structures are normally used to avoid the effects of limit cycles.

Characteristics of a Limit Cycle Oscillation with respect to the System by the following Difference equation

$$y(n) = 0.95 y(n-1) + x(n)$$
.

Let y(n) be the output of the system after the product term 0.95 y(n-1) is quantized after rounding. i.e.,

$$y_{r}(n) = Q_{r}[0.95y(n-1)] + x(n)$$
$$x_{n} = \begin{cases} 0.75 & \text{for } n = 0 \\ 0 & \text{for } n \neq 0 \end{cases}$$

Let b = 4 bits are used to represent the quantized product excluding sign bit.

With n=0

Let

y_r(n)= Q_r[0.95y_r(n −1)]+ x(n) ∴ y (0)= Q [0.95y₁(−1)]+ x(0) = Q [0.95×0]+ 0.75 Since y_r(−1)= 0 = 0.75

$$[0.75]_{10} = [0.11]_2$$

: 4-bits rounded value of $[0.11]_2$ will be $[0.1100]_2$ i.e., 0.75 only.

 \therefore y (0)= 0.75 after 4 bits rounding

With n = 1

$$y_{r}(1) = Q_{r}[0.95y_{r}(0)] + x(1) = Q_{r}[0.95 \times 0.75] + 0 = Q_{r}[0.7125]$$
$$[0.7125]_{0} = [0.1\ 0\ 1\ 1\ 0\ 1\ 1\ 0\ 0\ 1\ 1\ 0\ 0\ \cdots]_{2}$$

Note that 0.7125 requires infinite binary digits for its representation. Let us round it to 4 bits.

$$\therefore Q_r [0.7125]_{10} = [0.1 \ 0 \ 1 \ 1]_2$$
 upto 4 bits

But decimal equivalent of $[0.1 \ 0 \ 1 \ 1]_2$ is 0.6875.

$$\therefore y_r(1) = 0.6875$$

This means the actual value of $y_r(1) = 0.7125$ is changes to 0.6875 due to 4-bits quantization.

With n = 2

Here $\delta = \frac{1}{2^{b}} = \frac{1}{2^{4}} = 0.0625$ $y(n-1) \le \frac{0.0625/2}{1-0.95} \le 0.625$

Dead band = [-0.625, +0.625]

Signal Scaling to Prevent Limit Cycle Oscillations: This is zero input condition. Following table lists the values of y(n) before and after quantization. Here the values are rounded to nearest integer value.

n	y(n) before quantization	y(n) after quantization
-1	12	12
0	10.8	11
1	9.72	10
2	8.748	9
3	7.8732	8
4	7.08588	7
5	6.377292	6
6	5.7395628	6
7	5.1656065	5
8	4.6490459	5

From table observe that if $y(-1) \le 5$, y(n) = y(-1) for $n \ge 0$ for zero input. Hence the dead band will be [-5,5].

Since the values are rounded to nearest integer after quantization, the step size will be $\delta = 1$. Hence dead band can also be calculated as follows:

$$y(n-1) = \frac{\delta/2}{1-\alpha}$$
, Here $\alpha = 0.9$, $y(n-1) = \frac{1/2}{1-0.9} = 5$

Thus the dead band is [-5,5].

Dynamic Range Scaling to Prevent the Effects of Overflow: The overflow can take place at some internal nodes when the digital filters are implemented by using fixed point arithmetic. Such nodes can be inputs/outputs of address or multipliers. This overflow can take place even if the inputs are scaled. Because of such overflow at intermediate points, produces totally undesired output or oscillations. The overflow can be avoided by scaling the internal signal levels with the help of scaling multipliers. These scaling multipliers are inserted at the apprpriate points in the filter structure to avoid possibilities of overflow. Sometimes these scaling multipliers are absorbed with the existing multipliers in the structure to reduce the total number and complexity.

At which node the overflow will take place is not known in advance. This is because the overflow depends upon type of input signal samples. Hence whenever overflow takes place at some node, the scaling should be done dynamically. Hence dynamic range scaling in the filter structure can avoid the effects of overflow.

Let $u_r(n)$ be the signal sample at r^{th} node in the structure. Then the scaling should ensure that,

 $|u_r(n)| \le 1$ for all r and n.

TRADE OFF BETWEEN ROUND-OFF AND OVERFLOW NOISE, MEASUREMENT OF COEFFICIENT QUANTIZATION EFFECTS THROUGH POLE-ZERO MOVEMENT

Errors in Rounding and Truncation Operations : The computations like multiplication or addition are performed the result is truncated or rounded to nearest available digital level. This operation introduces an error. Hence the performance of the system is changed from expected value.

Truncation Error: This error is introduced whenever the number is represented by reduced number of bits.

Let $Q_t(x)$ be the value after truncation , then truncation error will be,

 $\varepsilon_r = Q_t(x) - (x)$

Here x is the original value of the number.

Rounding Error : This error is introduced whenever the number is rounded off to the nearest digital level. The number of bits used to represent the rounded number are generally less than the number of bits required for actual number.

Let $Q_r(x)$ be the value after rounding. Then rounding error will be,

$$\varepsilon_r = Q_r(x) - x$$

Here x is the original value of a number.

Tradeoff between roundoff and overflow noise:

Scaling operation

Scaling is a process of readjusting certain internal gain parameters in order to constrain internal signals to a range appropriate to the hardware with the constraint that the transfer function from input to output should not be changes.

The filter in figure with unscaled node x has the transfer function

$$H(z) = D(z) + F(z)G(z) \qquad \dots \dots (1)$$

To scale the node x, we divide F(z) by some number β and multiply G(z) by the same number as in figure. Although the transfer function does not change by this operation, the signal level at node x has been changes. The scaling parameter β can be chosen to meet any specific scaling rule such as

$$l_{1} \text{scaling:} \beta = \sum_{i=0}^{\infty} f(i) \rfloor \qquad \dots \dots (2)$$
$$l_{2} \text{scaling:} \beta = \delta \sqrt{\sum_{i=0}^{\infty} \lfloor f^{2}(i) \rfloor} \qquad \dots \dots (3)$$

Where f(i) is the unit-sample response from input to the node x and the parameter δ can be interpreted to represent the number of standard deviations





Figure: A filter with unscaled node x and (b) A filter with scaled node x'

Representable in the register at node x if the input is unit-variance white noise. If the input is bound by $|u(n)| \le 1$, then,

$$\left| \mathbf{x}(\mathbf{n}) \right| = \left| \sum_{i=0}^{\infty} f(i) \mathbf{u}(\mathbf{n}-i) \right| \le \sum_{i=0}^{\infty} f(i) \left| \dots \dots (4) \right|$$

Equation represents the true bound on the range of x and overflow is completely avoided by l_1 scaling in (2), which is the most stringent scaling policy.

In many cases, input can be assumed to be white noise. Although we cannot compute the variance at node x. for unit-variance white noise input,

$$\begin{bmatrix} 2 \\ E x (n) \end{bmatrix} = \sum_{i=0}^{\infty} f^{2}(i) \qquad \dots \dots (5)$$

Since most input signals can be assumed to be white noise, l_2 scaling is commonly used. In addition, (5) can be easily computed. Since (5) is the variance (not a strict bound), there is a possibility of overflow, which can be reduced by increasing δ in (3). For large values of δ , the internal variables are scaled conservatively so that no overflow occurs. However, there is a trade –off between overflow and roundoff noise, since increasing δ deteriorates the output SNR (signal to noise ratio).



Figure: Model of roundoff error

Roundoff Noise: If two W-bit fixed point fraction numbers are multiplies together, the product is (2W-1) bit long. This product must eventually be quantized to W-bits by rounding or truncation. For example, consider the 1st-order IIR filter shown in figure. Assume that the input wordlength is W=8bits. If the multiplier coefficient wordlength is also the same, then to maintain full precision in the output we need to increase the output wordlength by 8 bits per iterations. This is clearly infeasible. The alternative is to roundoff (or truncate) the output to its nearest 8-bit representation.



Figure: Error probability distribution

The result of such quantization introduces roundoff noise e(n). For mathematical ease a system with roundoff can be modeled as an infinite precision system with an external error input. For example in the previous case (shown in figure) we round off the output of the multiply add operation and an equivalent model is shown in figure.

Although rounding is not a linear operation, its effect at the output can be analyzed using linear system theory with the following assumptions about e(n):

- 1. E(n) is uniformly distributed white noise.
- 2. E(n) is a wide –sense stationary random process, i.e., mean and covariance of e(n) are independent of the time index n.
- 3. E(n) is uncorrelated to all other signals such as input and other noise signals.

Let the wordlength of the output be W-bits, then the roundoff error e(n) can be given by

$$\frac{-2^{-(w-1)}}{2} \le e(n) \le \frac{2^{-(w-1)}}{2} \qquad \dots \dots (6)$$

Since the error is assumed to be uniformly distributed over the interval given in (6), the corresponding probability distribution is shown in figure, where Δ is the length of the interval (i.e., $2^{-(w-1)}$).

Note that since mean is zero, variance is simply $E[e^2(n)]$

Figure: Signal flow graph

In other words (8) can be rewritten as

$$\sigma_{e}^{2} = \frac{2^{-2w}}{3} \qquad \dots \dots (9)$$

Where σ_e^2 is the variance of the roundoff error in a finite precision, W-bit wordlength system. Since the

variance is proportional to 2^{-2w} , increase in wordlength by 1 bit decreases the error by a factor of 4.

The purpose of analyzing roundoff noise is to determine its effect at the output signal. If the noise variance at the output is not negligible in comparison to the output signal level, the worlength should be increase or some low noise structures should be used. Therefore, we need to compute SNR at the output, not just the noise gain to the output. In the noise analysis, we use a double length accumulator model, which means rounding is performed after two (2w-1)-bit products are added. Also, notice that multipliers are the sources for roundoff noise.

Effects of Coefficient Quantization in FIR filters: Let us consider the effects of coefficient quantization in FIR filters. Consider the transfer function of the FIR filter of length M,

$$H(z) = \sum_{n=0}^{M-1} h(n) z^{-n}$$

The quantization of h(n) takes place during implementation of filter. Let the quantized coefficients be denoted by h(n) and e(n) be the error in quantization. Then we can write,

$$h(n) = h(n) + e(n)$$

And the new filter transfer function becomes,

n=0

$$H(z) = \sum_{n=0}^{M-1} h(n)z^{-n} = \sum_{n=0}^{M-1} [h(n) + e(n)]z^{-n} = \sum_{n=0}^{M-1} (n)z^{-n} + \sum_{n=0}^{M-1} e(n)z^{-n} = H(z) + E(z)$$

Where, $E(z) = \sum_{n=0}^{M-1} e(n)z^{-n}$



Figure: Model of FIR filter with quantizer coefficients

Here observe that the FIR filter with quantized coefficients can be modelled as parallel connection of two FIR filters H(z) and E(z).

In the figure, H(z) is the FIR filter with unquantized coefficients, and E(z) is the FIR filter representing coefficient quantization error.

The frequency response of FIR filter with quantized coefficients as,

 $H(\omega) = H(\omega) + E(\omega)$

Here $E(\omega)$ is the error in the desired frequency response which is given as,

$$E(\omega) = \sum_{n=0}^{M-1} e(n)e^{-j\omega n}$$

Consider the magnitude of error i.e.,

$$\left| E(\omega) \right| = \left| \sum_{n=0}^{M-1} e(n) e^{-j\omega n} \right| \implies E(\omega) \leq \sum_{n=0}^{M-1} e(n) e^{-j\omega n} \quad \left| \implies \leq \sum_{n=0}^{M-1} |e(n)| \quad (\therefore e^{-j\omega n} = |1) |e(n)| \leq \sum_{n=0}^{M-1} |e(n)| = |1|$$

The upper bound is reached if all the errors have same sign and have the maximum value in the range. If we consider e(n) to be statistically independent random variables, then more realistic bound is given by standard derivation of $E(\omega)$ i.e.;

 $\sigma_{E}(\omega)$ is the standard derivation of error in frequency response i.e., E(ω).

DEADBAND EFFECTS

Deadband and Deadband of First Order Filter: Dead band is the range of output amplitudes over which limit cycle oscillations takeplace

Dead band of first order filter

Consider the first order filter,

$$y(n) = \alpha y(n-1) + x(n)$$

Here $\alpha y(n-1)$ is the product term. After rounding it to 'b' bits we get,

2

$$y(n) = Q_{t} \left[\alpha y(n-1) \right] + x(n)$$

When limit cycle oscillations take place,

$$Q_{r}[\alpha y(n-1)] = \pm y(n-1) \qquad \dots \dots (1)$$

The error due to rounding is less than $\frac{1}{2}$. Hence,

$$|Q[\alpha y(n-1)-\alpha y(n-1)]| \leq \delta_2$$

From equation (1) above equation can be written as,

$$\begin{array}{l} \ddagger y(n-1) - \alpha y(n-1) \leq \frac{\delta}{2} \\ \vdots y(n-1) \left[1 + \alpha \right] \leq \frac{\delta}{2} \\ \vdots y(n-1) \leq \frac{\delta/2}{1-|\alpha|} \end{array}$$