

Discrete time signals and systems

## Discrete time signals

- Discrete time: discrete in time only (sampled but not quantised)

Digital: discrete in time and amplitude (sampled and quantised)

- Discrete time signals  $x = (\dots, x_{-1}, x_0, x_1, \dots)$  typically have real/complex elements, and could be right-sided, left-sided, two-sided or finite-length

↳ Right-sided: starts at  $x_m$  and  $x_k = 0$  for  $k < m$  [start at  $x_0$ : semi-infinite]

↳ Left-sided: ends at  $x_m$  and  $x_k = 0$  for  $k > m$

↳ Two-sided: no highest/lowest index for non-zero values

↳ Finite-length: starts at  $x_m$  and ends at  $x_{m+l-1}$  (length  $l$ )

- A discrete time signal may be a sampled continuous signal, in which case, w/ reference to the original continuous sample  $x_t = x(kT)$ , we can denote the discrete time signal as

$$x = (x(0), x(kT), x(2kT), \dots)$$

where  $T$  is the sampling interval

- We could treat a discrete time signal as a sequence of no./ has a normalised sampling freq. of 1 Hz.

## Discrete time systems

- A discrete time system takes discrete time signal inputs and produces discrete time signal outputs.

- The definition of linear time-invariant (LTI) system  $L$  extends from continuous to discrete time

If for  $\{u_k\}$  we get  $\{y_k\}$  and for  $\{u'_k\}$  we get  $\{y'_k\}$ , then

↳ Linear:  $L(\alpha\{u_k\} + \beta\{u'_k\}) = \alpha L(\{u_k\}) + \beta L(\{u'_k\}) = \alpha\{y_k\} + \beta\{y'_k\}$

↳ Time-invariant:  $L(\{u_{k+\Delta}\}) = \{y_{k+\Delta}\}$

- Discrete LTI systems are governed by difference eqns w/ const. coeffs.

## Discrete convolution

- The convolution of the sequences  $\{a_k\}$  and  $\{b_k\}$  is the sequence  $\{y_k\}$ .

$$y_k = \sum_{n=0}^{\infty} a_n b_{k-n}$$

- The summation limits can be simplified for the special cases:

↳ Semi-infinite sequence  $\{a_k\}_{k \geq 0}$  and  $\{b_k\}_{k \geq 0}$

$$y_k = \sum_{n=0}^{K} a_n b_{k-n}$$

↳ Finite-length sequence  $\{a_k\}_{k=0}^{L-1}$

$$y_k = \sum_{n=0}^{L-1} a_n b_{k-n}$$

# For Personal Use Only -bkwk2

Convolution and polynomial multiplication

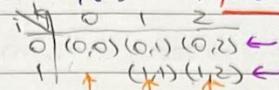
- We can map semi-infinite sequences to power series, i.e.  $\{u_k\} \rightarrow u(D)$

$$u(D) = a_0 + a_1 D^1 + a_2 D^2 + \dots$$

The convolution of semi-infinite sequences is equivalent to multiplication of power series.

$$\begin{aligned} a(D)b(D) &= \left( \sum_{i=0}^{\infty} a_i D^i \right) \left( \sum_{j=0}^{\infty} b_j D^j \right) = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} a_i b_j D^{i+j} = \sum_{k=0}^{\infty} \sum_{i+j=k} a_i b_{k-i} D^k \quad (\text{sub } k=i+j) \\ &= \sum_{k=0}^{\infty} \left( \sum_{i=0}^k a_i b_{k-i} \right) D^k = \sum_{k=0}^{\infty} Y_k D^k = y(D) \rightarrow a(D)b(D) = y(D) \end{aligned}$$

where we used  $\sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \leftrightarrow \sum_{k=0}^{\infty} \sum_{i+j=k}$



- $D$  is a dummy variable — we never plan to assign a value to the variable  $\rightarrow$  formal power series.
- As we do not assign values to the dummy variable, we do not need to worry about the radius of convergence (sufficient for mapping convolutions to multiplications)
- Z-transform is not a formal power series  $\&$  we do assign a complex value to the variable  $z$ .

## Z-transform

### Z-transform

- For a semi-infinite sequence  $\{a_k\}_{k \geq 0}$ , the one-sided Z-transform is defined as

$$Z[\{a_k\}] = A(z) = \sum_{k=0}^{\infty} a_k z^k$$

For a two-sided sequence  $\{a_k\}_{k=-\infty}^{\infty}$ , the two-sided Z-transform is defined as

$$Z[\{a_k\}] = A(z) = \sum_{k=-\infty}^{\infty} a_k z^k$$

- For the two-sided Z-transform, the sequence  $\{a_k\}$  is not uniquely determined by  $A(z)$ , so we must also specify the radius of convergence (not an issue for one-sided transforms)

### Properties of Z-transform

① Linearity :  $Z[\alpha \{a_k\} + \beta \{b_k\}] = \alpha Z[\{a_k\}] + \beta Z[\{b_k\}]$ ,  $Z[\frac{1}{k} \{a_k\}] = \frac{1}{z} Z[\{a_k\}]$

- The Z-transform only involves doing linear operations on the sequence values,

$$Z[\alpha \{a_k\} + \beta \{b_k\}] = \sum_{k=-\infty}^{\infty} (\alpha a_k + \beta b_k) z^k = \alpha \sum_{k=-\infty}^{\infty} a_k z^k + \beta \sum_{k=-\infty}^{\infty} b_k z^k = \alpha Z[\{a_k\}] + \beta Z[\{b_k\}]$$

It also follows that if the sequence defn involves my parameter  $\lambda$ , i.e.  $a_k = f_n(\lambda)$ ,

$$Z[\frac{1}{k} \{a_k\}] = \sum_{k=-\infty}^{\infty} \frac{1}{k} a_k z^k = \frac{1}{\lambda} \sum_{k=-\infty}^{\infty} a_k z^k = \frac{1}{\lambda} Z[\{a_k\}]$$

② Conjugate symmetry :  $A(\bar{z}) = \overline{A(z)}$

- The Z-transform has conjugate symmetry for real-valued signals  $\{a_k\}$ .

$$A(\bar{z}) = \sum_{k=0}^{\infty} a_k \bar{z}^k = \left( \sum_{k=0}^{\infty} a_k z^k \right)^* = \overline{A(z)}$$

- For the two-sided transform, there is the reverse property that for  $\{a_k\}$  s.t.  $a_{-k} = \bar{a}_k$ ,

the Z-transform of  $\{a_k\}$ ,  $A(z)$  is real on the unit circle.

$$\text{(e.g. } a_k = a + jb, a_{-k} = a - jb \text{ and } z = e^{j\theta}, A(e^{j\theta}) = a_0 + \sum_{k=1}^{\infty} (a+jb)e^{jk\theta} + (a-jb)e^{-jk\theta} = a_0 + 2 \sum_{k=1}^{\infty} a_k \cos(k\theta) + j \sin(k\theta))$$

# For Personal Use Only -bkwk2

$$\textcircled{3} \text{ Time delay / time advance: } z[\{a_{k+1}\}] = z^{-1}A(z) + a_0 ; z[\{a_{k+3}\}] = z^3A(z) - za_0$$

- Defining the time delay operation:  $\{a_k\} \rightarrow \{a_{k+1}\}$ .

$$z[\{a_{k+1}\}] = \sum_{k=0}^{\infty} a_{k+1} z^k = z \sum_{k=0}^{\infty} a_k z^{k+1} = z \sum_{k=0}^{\infty} a_k z^k + a_0 = z^{-1}A(z) + a_0$$

Assuming  $a_0 = 0$ , then  $z[\{a_{k+1}\}] = z^{-1}A(z)$ .

- Defining the time advance operation:  $\{a_k\} \rightarrow \{a_{k+3}\}$

$$z[\{a_{k+3}\}] = \sum_{k=0}^{\infty} a_{k+3} z^k = z^3 \sum_{k=0}^{\infty} a_{k+1} z^{k+1} = z^3 \sum_{k=0}^{\infty} a_k z^k - za_0 = z^3A(z) - za_0$$

\* The five-shift properties can be applied repeatedly for longer shifts.

$$\textcircled{4} \text{ Convolution: } z[\{a_k\} * \{b_k\}] = A(z)B(z)$$

- Consider the  $z$ -transform of the convolution between  $\{a_k\}, \{b_k\}$ .

$$z[\{a_k\} * \{b_k\}] = \sum_{k=0}^{\infty} \left( \sum_{i=0}^k a_i b_{k-i} \right) z^k = \sum_{i=0}^{\infty} \sum_{k=i}^{\infty} a_i b_{k-i} z^k = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} a_i b_j z^{i+j} = \left( \sum_{i=0}^{\infty} a_i z^i \right) \left( \sum_{j=0}^{\infty} b_j z^j \right) = A(z)B(z)$$

where we used  $\sum_{k=0}^{\infty} \sum_{i=0}^k \leftrightarrow \sum_{i=0}^{\infty} \sum_{k=i}^{\infty}$  as before

## Common $z$ -transforms

$$\textcircled{1} \text{ Geometric: } z[q^k] = \frac{1}{1-qz} ; z[1] = \frac{1}{1-z} ; z[e^{ak}] = \frac{1}{1-e^{az}}$$

- Consider the geometric sequence  $a_k = q^k$ .

$$z[q^k] = \sum_{k=0}^{\infty} q^k z^k = \sum_{k=0}^{\infty} (qz)^k = \frac{1}{1-qz}$$

$q=1 \rightarrow$  step sequence,  $a_k=1$

$$z[1] = \frac{1}{1-z}$$

$q = e^{ak} \rightarrow$  exponential sequence,  $a_k = e^{ak}$

$$z[e^{ak}] = \frac{1}{1-e^{az}}$$

$$\textcircled{2} \text{ Sine/cosine: } z[\sin(ak)] = \frac{z \sin a}{1-2z \cos a + z^2} ; z[\cos(ak)] = \frac{1-z \cos a}{1-2z \cos a + z^2}$$

- Applying linearity, we can get the  $z$ -transform for the sine/cosine sequences.

$$a_k = \sin(ak) : z[\sin(ak)] = z\left[\frac{1}{2j}(e^{jk} - e^{-jk})\right] = \frac{1}{2j}\left(\frac{1}{1-je^{-az}} - \frac{1}{1-je^{az}}\right) = \frac{z \sin a}{1-2z \cos a + z^2}$$

$$a_k = \cos(ak) : z[\cos(ak)] = z\left[\frac{1}{2}(e^{jk} + e^{-jk})\right] = \frac{1}{2}\left(\frac{1}{1-je^{-az}} + \frac{1}{1-je^{az}}\right) = \frac{1-2z \cos a}{1-2z \cos a + z^2}$$

$$\textcircled{3} \text{ Scale by } r^k : z[r^k a_k] = G(r^k z)$$

- Scaling a sequence  $\{a_k\}$  w/ a geometric sequence  $r^k$ , i.e.  $a_k = r^k g_k$ .

$$z[r^k g_k] = \sum_{k=0}^{\infty} r^k g_k z^k = \sum_{k=0}^{\infty} r^k (r^k z)^k = G(r^k z)$$

$$\textcircled{4} \text{ Ramp: } z[k] = \frac{z^{-1}}{(1-z)^2}$$

- Differentiating the geometric sequence,

$$z[kq^{k-1}] = z\left[\frac{d}{dt} z^t\right] = \frac{1}{z} \frac{d}{dt} z^t = \frac{1}{z} \frac{1}{1-z} = \frac{z^{-1}}{(1-z)^2}$$

$q=1 \rightarrow$  ramp sequence,  $a_k = k$

$$z[k] = \frac{z^{-1}}{(1-z)^2}$$

# For Personal Use Only -bkwk2

Initial value theorem (INT) and final value theorem (FVT)

- The INT states that

$$\lim_{z \rightarrow \infty} A(z) = a_0$$

$$\Rightarrow \lim_{z \rightarrow \infty} A(z) = \lim_{z \rightarrow \infty} \left( \sum_{k=0}^{\infty} a_k z^{-k} \right) = \lim_{z \rightarrow \infty} (a_0 + a_1 z^{-1} + a_2 z^{-2} + \dots) = a_0$$

- The FVT states if all the poles of  $(z-1)A(z)$  lie strictly within the unit circle,

$$\lim_{z \rightarrow 1^-} (z-1)A(z) = \lim_{k \rightarrow \infty} a_k$$

$\hookrightarrow$  All poles of  $A(z)$  are in the unit circle, except possibly a pole at  $z=1$ .

$$\text{Assuming distinct poles, } A(z) = \frac{a_0}{1-z^{-1}} + \sum_i \frac{a_i}{1-p_i z^{-1}} \quad \text{where } |p_i| < 1$$

$$\text{By inverse z-transform, } a_k = a_0 + \sum_i a_i p_i^k$$

$$\text{Taking limits } k \rightarrow \infty, \quad \lim_{k \rightarrow \infty} a_k = \lim_{k \rightarrow \infty} (a_0 + \sum_i a_i p_i^k) = a_0$$

$$\text{Consider } (z-1)A(z) = z(1-z^{-1})A(z) = z a_0 + (z-1) \sum_i \frac{a_i}{1-p_i z^{-1}}$$

$$\text{Taking limits } z \rightarrow 1, \quad \lim_{z \rightarrow 1^-} (z-1)A(z) = \lim_{z \rightarrow 1^-} (z a_0 + (z-1) \sum_i \frac{a_i}{1-p_i z^{-1}}) = a_0$$

$$\therefore a_0 = \lim_{k \rightarrow \infty} a_k = \lim_{z \rightarrow 1^-} (z-1)A(z)$$

## Z-transform and other transforms

### Z-transform and Laplace transform

- For continuous time signals  $x(t)$ , the Laplace transform is defined as

$$\tilde{x}(s) = \int_0^\infty x(t) e^{-st} dt$$

Analogously, we can define a transform for a discrete time signal  $\{x(kT)\}_{k \geq 0}$ , as

$$\sum_{k=0}^{\infty} x(kT) e^{-skT}$$

If we define  $x_k = x(kT)$  and  $z = e^{sT}$ , we get the Z-transform

$$X(z) = \sum_{k=0}^{\infty} x_k z^{-k}$$

- Alternatively, we can obtain the Z-transform by taking the Laplace transform of a continuous time representation of a discrete time signal,  $x^*(t)$ .

$$x^*(t) = \sum_{k=0}^{\infty} x(kT) \delta(t-kT)$$

$$\therefore L[x^*(t)] = \int_0^\infty \sum_{k=0}^{\infty} x(kT) \delta(t-kT) e^{-st} dt = \sum_{k=0}^{\infty} x(kT) e^{-skT} = \sum_{k=0}^{\infty} x_k z^{-k} = Z[\{x_k\}]$$

# For Personal Use Only -bkwk2

## Z-transform and DTFT

- consider a band-limited signal  $x(t)$  where  $x(t)=0$  for  $t < 0$ ,

the DTFT of the sampled signal is defined as

$$X(\theta) = \sum_{k=0}^{\infty} x(kT) e^{j\theta k}$$

where  $\theta/kT$  is the normalised frequency.

If we define  $X_k = x(kT)$  and  $z = e^{j\theta}$ , we get the Z-transform

$$x(z) = \sum_{k=0}^{\infty} X_k z^k$$

$\therefore X(z = e^{j\theta}) = X(\theta)$ , i.e. Z-transform on the unit circle is the DTFT.

- Since the Z-transform on the unit circle is the DTFT, we can express the IDTFT as

$$x_k = x(kT) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\theta) e^{jk\theta} d\theta = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\theta}) e^{jk\theta} d\theta$$

i.e. we can do an inverse Z-transform using  $X_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\theta}) e^{jk\theta} d\theta$

- The Z-transform on the unit circle shows the normalised Fourier spectrum

(i.e. normalised frequency  $\theta \in [-\pi, \pi]$  and actual frequency  $w \in [-f_s/2, f_s/2]$ ).

- For real-valued signals, the conjugate symmetric property states  $X(e^{j\theta}) = \overline{X(e^{-j\theta})}$

(i)  $|X(e^{j\theta})| = |X(e^{-j\theta})|$  and (ii)  $L(X(e^{j\theta})) = -L(X(e^{-j\theta}))$

## Inverse Z-transform

### Inverse Z-transform

- The inverse Z-transform can be found using

↳ Partial fractions

↳ Polynomial long division

↳ Taylor expansion.

↳ IDTFT formula.

#### ① Partial fractions

- Decompose the rational function into partial fractions and use standard Z-transforms

- e.g.:  $A(z) = \frac{z^3}{z^2(z-1)(z-2)}$  find  $\{a_k\}_{k \geq 0}$ .

$$A(z) = \frac{z^3}{z^2(z-1)(z-2)} = z^2 \left[ \frac{1}{z-1} - \frac{1}{z-2} \right] = z^3 \left[ \frac{1}{z-1} - \frac{1}{z-2} \right] \rightarrow a_k = \begin{cases} 0 & k < 2 \\ 1 & k=3 \\ -1 & k=2 \\ 2 & k \geq 3 \end{cases}$$

#### ② Polynomial long division

- Explicitly compute the coefficients of  $z^0, z^1, \dots \rightarrow$  first few terms of  $\{x_k\}$ . [not analytic]

- e.g.:  $A(z) = \frac{z^3}{z^2(z-1)(z-2)}$  find  $\{a_k\}_{k \geq 0}$ .

$$A(z) = \frac{z^3}{z^2(z-1)(z-2)} = \frac{1}{z^2} \frac{z^3}{z^3 - 3z^2 + 2}$$

$$\frac{z^3 - 3z^2 + 2}{z^3} \overline{z^3 - 3z^2 + 2}$$

$$- \frac{z^3 + 2z^2 - 2z^3 - 6z^4 + \dots}{-2z^4} \rightarrow A(z) = z^{-2}(z^1 - 2z^3 - 6z^4 + \dots)$$

$$- \frac{-2z^4 + 6z^2 + 4z^3}{-6z^4 + 4z^3} = z^3(1 - 2z^2 - 6z^3 + \dots)$$

$$\therefore \{a_k\}_{k \geq 0} = \{0, 0, 0, 1, 0, -2, -6, \dots\}$$

# For Personal Use Only -bkwk2

## ③ Taylor expansion

- Use the Taylor expansion to get a power series  $\rightarrow$  read off the coefficients.

- e.g.  $A(z) = \ln(1+z^2)$  Find  $\{a_k\}_{k \geq 0}$ .

$$A(z) = \ln(1+z^2) = z^{-1} - \frac{1}{2}z^2 + \frac{1}{3}z^3 - \frac{1}{4}z^4 + \dots \rightarrow \{a_k\}_{k \geq 0} = \{0, 1, -\frac{1}{2}, \frac{1}{3}, -\frac{1}{4}, \dots\}$$

## ④ IDTFT formula.

- Given the Z-transform on the unit circle, we can use  $a_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} A(e^{j\theta}) e^{j k \theta} d\theta$ .

- e.g.  $A(z=e^{j\theta}) = \begin{cases} 1 & \text{if } \theta \in [-\pi/4, \pi/4] \\ 0 & \text{otherwise.} \end{cases}$

$$a_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} A(e^{j\theta}) e^{j k \theta} d\theta = \frac{1}{2\pi} \int_{-\pi/4}^{\pi/4} 1 \cdot e^{jk\theta} d\theta = \frac{1}{\pi k} \frac{1}{2j} (e^{jk\pi/4} - e^{-jk\pi/4}) = \frac{1}{\pi k} \frac{\sin(k\pi/4)}{2j} = \frac{1}{\pi k} \sin(\frac{k\pi}{4}).$$

System transfer function

System transfer function

Pulse response of LTI systems

- We define the unit pulse (Kronecker delta function)  $\delta_k$  in discrete time as

$$\delta_k = \begin{cases} 1 & k=0 \\ 0 & k \neq 0 \end{cases}$$

- If the i/p to a system is the unit pulse input  $\{x_k\}_{k \geq 0} = \delta_k = (1, 0, 0, \dots)$ , then the o/p to a system  $\{h_k\}_{k \geq 0} = (h_0, h_1, h_2, \dots)$  is the pulse response of the system.

- For a LTI system w/ pulse response  $\{h_k\}_{k \geq 0}$ , we can write the general i/p  $\{x_k\}_{k \geq 0}$  in the form

$$x_k = \sum_{n=0}^k x_n \delta_{k-n}$$

Using the superposition principle, the o/p of the system is thus

$$y_k = \sum_{n=0}^k x_n h_{k-n}$$

$Z$ -transfer function.

- The system described by the following difference eqn, subject to zero initial conditions is LTI.

$$y_k + a_1 y_{k-1} + \dots + a_n y_{k-n} = b_0 x_k + b_1 x_{k-1} + \dots + b_m x_{k-m}, \quad x_k = y_k = 0 \text{ for } k < 0.$$

Taking the  $Z$ -transform,

$$Y(z)[1 + a_1 z^{-1} + \dots + a_n z^{-n}] = X(z)[b_0 + b_1 z^{-1} + \dots + b_m z^{-m}]$$

We define the transfer function  $G(z)$  to be

$$G(z) = \frac{Y(z)}{X(z)} = \frac{b_0 + b_1 z^{-1} + b_2 z^{-2} + \dots + b_m z^{-m}}{1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_n z^{-n}}$$

Pulse response and  $Z$ -transfer function.

- The o/p  $\{y_k\}_{k \geq 0}$  of a LTI system can be found by convolving the i/p  $\{x_k\}_{k \geq 0}$  w/ the pulse response  $\{h_k\}_{k \geq 0}$  of the system

$$y_k = \sum_{n=0}^k x_n h_{k-n} \quad \text{or} \quad \{y_k\} = \{x_k\} * \{h_k\}$$

Taking the  $Z$ -transform gives

$$Y(z) = H(z)X(z)$$

where  $H(z) = \frac{Y(z)}{X(z)}$  is the system transfer.

- The system transfer function  $H(z)$  is the  $Z$ -transform of the system pulse response  $\{h_k\}_{k \geq 0}$ .

$$H(z) = Z[\{h_k\}]$$

Solution of linear difference equations.

- Take the  $Z$ -transform on the linear difference eqn, applying the linearity and time shift property.
- Substitute in any initial conditions and/or inputs
- Solve for  $Y(z)$  and take the inverse  $Z$ -transform to find  $\{y_k\}_{k \geq 0}$ .

# For Personal Use Only -bkwk2

Finite impulse response (FIR) and infinite impulse response (IIR)

- FIR filters have pulse responses that terminate after a finite no. of time steps.

$$\{h_k\}_{k \geq 0} = \{h_0, h_1, \dots, h_m, 0, 0, \dots\}.$$

IIR filters have pulse responses that do not terminate after a finite no. of time steps

$$\{h_k\}_{k \geq 0} = \{h_0, h_1, h_2, \dots\}.$$

- FIR systems have transfer functions of the special form

$$G(z) = b_0 + b_1 z^{-1} + \dots + b_m z^{-m} = \frac{b_0 z^m + b_1 z^{m-1} + \dots + b_m}{z^m}$$

where  $a_k = 0$  (no feedback), and all the poles are at  $z=0$ .

IIR systems have transfer functions in the general form

$$G(z) = \frac{b_0 + b_1 z^{-1} + \dots + b_m z^{-m}}{1 + a_1 z^{-1} + \dots + a_n z^{-n}}$$

where the poles are at arbitrary locations.

## Causality

- discrete-time systems whose pulse response  $\{h_k\}$ , where  $h_k = 0$  for  $k < 0$  are causal.
- consider a non-causal pulse response  $\{g_k\} = \{\dots, g_{-2}, g_{-1}, g_0, g_1, \dots\}$ .

By convolution, the o/p of the system is

$$y_k = \sum_{n=-\infty}^{\infty} g_n x_{k-n} = \dots \underbrace{g_{-2} x_{k+2}}_{\text{future}} + \underbrace{g_{-1} x_{k+1}}_{\text{future}} + \underbrace{g_0 x_k}_{\text{past}} + \underbrace{g_1 x_{k-1}}_{\text{past}} + \dots$$

i.e. the current value of the o/p  $y_k$  depends on future values of the i/p  $x_{k+1}, x_{k+2}, \dots$

- In causal systems, only the past of an i/p sequence  $x_k, x_{k-1}, \dots$  can influence the current o/p  $y_k$ .
- The transfer function of causal systems have the form

$$H(z) = b_0 + b_1 z^{-1} + b_2 z^{-2} + \dots$$

for a transfer function w/ the following form to be causal

$$H(z) = \frac{b_0 z^m + b_1 z^{m-1} + \dots + b_m}{z^n + a_1 z^{n-1} + \dots + a_n}$$

we require  $m \leq n$ , otherwise we can rewrite the transfer function as

$$H(z) = c_1 z^{m-n} + \dots + c_m z^0 + \frac{b_0 z^m + b_1 z^{m-1} + \dots + b_m}{z^n + a_1 z^{n-1} + \dots + a_n}$$

which includes non-causal terms  $z, z^2, \dots, z^{m-n}$

\* Note for causal systems,  $\lim_{z \rightarrow \infty} G(z) = g_0 < \infty$ .

- Systems that operate in function of time and in real time are always causal.

- Examples of non-causal "real world systems" include

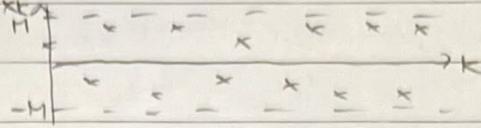
↳ Image processing (function of location)

↳ signal processing system w/ stored/buffered data (not in real time).

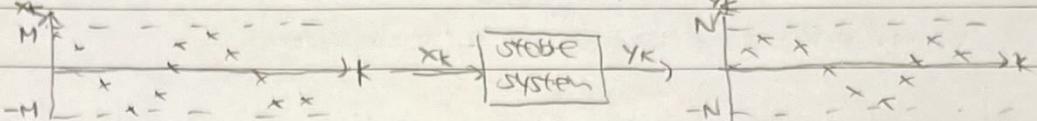
# For Personal Use Only -bkwk2

Stability.

- A signal  $\{x_k\}$  is bounded if there exists a const.  $M$  s.t.  $|x_k| < M \forall k$



- A discrete time system is stable if bounded i/p's give bounded o/p's (BIBO stability)



- Let  $G$  be a discrete time system w/ a rational transfer function

closed system

$$G(z) = \frac{b(z)}{a(z)} = \frac{b_0 z^m + \dots + b_m}{z^n + a_1 z^{n-1} + \dots + a_n}$$

where  $m \leq n$  and there are no common factors b/w  $a(z)$  and  $b(z)$ .

- Let the pulse response of  $G$  be  $\{y_k\}_{k>0}$ , then the following are equivalent: (i)  $G$  is stable,

(ii) all of the roots of  $a(z)$ ,  $p_i$ , satisfy  $|p_i| < 1$ , and (iii)  $\sum_{k=0}^{\infty} |y_k|$  is finite.

↪ Proof of (i)  $\rightarrow$  (ii) (Using contrapositive,  $\neg(iii) \rightarrow \neg(ii)$ , i.e.  $\neg(|p_i| < 1 \forall i) \rightarrow G \text{ not stable}$ )

• Suppose  $p_1, p_2, \dots$  are distinct. Then we can decompose  $G$  using partial fractions

$$G(z) = \frac{\alpha_1}{1-p_1 z^{-1}} + \dots + \frac{\alpha_n}{1-p_n z^{-1}} \rightarrow y_k = \alpha_1 p_1^k + \dots + \alpha_n p_n^k$$

• Suppose  $|p_i| > 1$  for some  $i$ , then  $y_k$  is unbounded.

∴ A bounded i/p (bounded) gives an unbounded o/p, so  $G$  is unstable.

• Suppose we have a pole on the unit circle, i.e.  $p_i = e^{j\theta_i}$  for some  $i$ .  $G$  can rewritten as

$$G(z) = \frac{\alpha_i}{1-e^{j\theta_i} z^{-1}} + \tilde{G}(z) \rightarrow y_k = \alpha_i e^{j\theta_i k} + \dots$$

The term that corresponds to the pole on the unit circle  $y_k = \alpha_i e^{j\theta_i k} \rightarrow$  bounded

• Now consider the i/p  $U(z) = \frac{a_i z^{-1}}{1-e^{j\theta_i} z^{-1}}$ , which corresponds to the sequence  $u_k \rightarrow$  bounded.

The corresponding o/p is  $y(z) = G(z)U(z) = \left( \frac{\alpha_i}{1-e^{j\theta_i} z^{-1}} + \tilde{G}(z) \right) \left( \frac{a_i z^{-1}}{1-e^{j\theta_i} z^{-1}} \right) = \frac{\alpha_i^2 z^{-1}}{(1-e^{j\theta_i} z^{-1})^2} + \tilde{G}(z)U(z)$

where the first term corresponds to the sequence  $y_k^i = \alpha_i^2 k e^{j\theta_i (k-1)}$  → unbounded  $\neg(y_k) \text{ grows w/k}$ .

∴ A bounded i/p gives an unbounded o/p, so  $G$  is unstable

$\Rightarrow$  If  $G$  is stable, then all of the  $a(z)$ ,  $p_i$ , satisfy  $|p_i| < 1$ .

↪ Proof of (ii)  $\rightarrow$  (iii) (i.e.  $|p_i| < 1 \forall i \rightarrow \sum_{k=0}^{\infty} |y_k| < \infty$ )

• Suppose  $p_1, p_2, \dots$  are distinct. Then we can decompose  $G$  using partial fractions.

$$G(z) = \frac{\alpha_1}{1-p_1 z^{-1}} + \dots + \frac{\alpha_n}{1-p_n z^{-1}} \rightarrow y_k = \alpha_1 p_1^k + \dots + \alpha_n p_n^k$$

∴  $\sum_{k=0}^{\infty} |y_k| \leq \alpha_1 \sum_{k=0}^{\infty} |p_1|^k + \dots + \alpha_n \sum_{k=0}^{\infty} |p_n|^k$  where RHS is finite if  $|p_i| < 1 \forall i$

↪ Proof of (iii)  $\rightarrow$  (i) (i.e.  $\sum_{k=0}^{\infty} |y_k| < \infty \rightarrow G$  is stable)

• Let  $\{x_k\}$  be an arbitrary bounded i/p, i.e.  $|x_k| \leq M$  for  $k \geq 0$ . Then the o/p is given by

$$y_k = \sum_{n=0}^{\infty} g_n x_{kn} \rightarrow |y_k| = \sum_{n=0}^{\infty} |g_n x_{kn}| \leq \sum_{n=0}^{\infty} |g_n| |x_{kn}| < \sum_{n=0}^{\infty} |g_n| M = M \sum_{n=0}^{\infty} |g_n|$$

∴ Given  $\sum_{n=0}^{\infty} |g_n| < \infty$ , we have  $|y_k| = M \sum_{n=0}^{\infty} |g_n| < \infty$  (bounded o/p).

(We can also prove (i)  $\rightarrow$  (iii) using contrapositive, w/  $u_{N+k} = \text{sign}(y_k) \rightarrow y_N = \sum_{k=0}^N g_k u_{N+k} = \sum_{k=0}^N |g_k| > M$ )

# For Personal Use Only -bkwk2

## Transient and steady state response

Transient and steady state for discrete-time systems.

- The general sol'n of a linear difference eqn. is

$$\{Y_k\}_{k \geq 0} = \{y_k\}_{k \geq 0} + \{s_k\}_{k \geq 0}$$

where  $\{y_k\}$  is the transient response and  $\{s_k\}$  is the steady state response.

- The transient response  $\{y_k\}$  decays for stable systems.

## Pulse response and poles

- Consider a LTI system  $G(z)$  subjected to a pulse imp  $\{x_k\} = \{1, 0, 0, \dots\}$ , so  $\Delta(z) = 1$ .

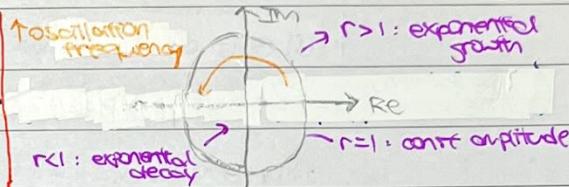
$$\text{The resp } \{y_k\} \text{ is given by } y_k = z^k [Y(z)] = z^k [G(z)\Delta(z)] = z^k [G(z)]$$

i.e. the poles of  $G(z)$  define the response to a pulse.

- In general, for a pole at  $z = re^{j\theta}$ ,

↳ the magnitude  $r$  gives the amplitude decay rate

↳ the argument  $\theta$  gives the oscillation rate



- Recall all the poles inside the unit circle  $|p| < 1 + i$  → stable system.

- consider the following pole locations

↳ Real poles  $[p = \lambda]$  :  $y_k = \lambda^k$

$$G(z) = \frac{1}{1-\lambda z^{-1}} \rightarrow y_k = \lambda^k$$

↳ Conjugate symmetric poles  $[p_1 = re^{j\theta}, p_2 = re^{-j\theta}]$  :  $y_k = 2\lambda^k \cos(k\theta)$

$$G(z) = \frac{1}{1-re^{j\theta}z^{-1}} + \frac{1}{1-re^{-j\theta}z^{-1}} \rightarrow y_k = \lambda^k (e^{jk\theta} + e^{-jk\theta}) = 2\lambda^k \cos(k\theta)$$

+ Real-valued impulse responses have complex conjugate pair poles.

↳ Repeated poles  $[p_1 = \lambda, p_2 = \lambda]$  :  $y_k = (1+k) \lambda^k$

$$G(z) = \frac{1}{(1-\lambda z^{-1})^2} = \frac{z^{-1}}{(1-\lambda z^{-1})^2} \cdot z \rightarrow y_k = (1+k) \lambda^k$$

## Stable filters and steady state response.

- In the  $z$ -domain, any linear filter can be written as

$$A(z) Y(z) = B(z) X(z) + C(z, y_i)$$

where  $C(z, y_i)$  takes into account initial conditions of the system

- Rearranging for  $Y(z)$  and taking the inverse  $z$ -transform,

$$Y_k = z^k \left[ \frac{B(z)}{A(z)} X(z) + \frac{C(z, y_i)}{A(z)} \right]$$

$$\text{As } k \rightarrow \infty, \quad \lim_{k \rightarrow \infty} Y_k = \lim_{k \rightarrow \infty} z^k \left[ \frac{B(z)}{A(z)} X(z) + \frac{C(z, y_i)}{A(z)} \right] = \lim_{k \rightarrow \infty} z^k \left[ \frac{B(z)}{A(z)} X(z) \right]$$

- For a stable system, stable roots in  $A(z)$  enforce the exponential decay of  $z^k \left[ \frac{C(z, y_i)}{A(z)} \right]$ ,

whereas  $\frac{B(z)}{A(z)} X(z)$  may not necessarily decay exponentially → stable filter 'forgets' initial conditions.

# For Personal Use Only -bkwk2

Steady state response to a step input

- provided all the poles of  $(z-1)Y(z) = (z-1)G(z)U(z)$  are strictly within the unit circle, we can use FUT,

$$\lim_{k \rightarrow \infty} Y_k = \lim_{z \rightarrow 1} (z-1)G(z)$$

- For a step i/p,  $U(z) = \frac{1}{1-z}$ , so the steady-state value of the op is given by  $\lim_{k \rightarrow \infty} Y_k = G(1)$

$$\lim_{k \rightarrow \infty} Y_k = \lim_{z \rightarrow 1} (z-1)Y(z) = \lim_{z \rightarrow 1} (z-1)G(z)U(z) = \lim_{z \rightarrow 1} (z-1)G(z) \frac{1}{1-z} = \lim_{z \rightarrow 1} zG(z) = G(1)$$

(For the continuous case, FUT:  $\lim_{t \rightarrow \infty} y(t) = \lim_{s \rightarrow 0} sG(s)$ ,  $U(s) = \frac{1}{s}$ , so  $\lim_{t \rightarrow \infty} y(t) = G(0)$  for step i/p.)

Steady state response to a sinusoidal input

- consider the complex sinusoidal  $y_p(k) = e^{j\theta k}$  to a system  $G(z) = \frac{b(z)}{(z-p_1)(z-p_2)\dots(z-p_n)}$

$$Y(z) = G(z)X(z) = \frac{b(z)}{(z-p_1)(z-p_2)\dots(z-p_n)} \frac{1}{1-e^{j\theta z}} = \frac{G(e^{j\theta})}{(1-e^{j\theta z})} + \text{terms of the form } \frac{B_i}{1-p_i z}$$

(In general,  $Y(z) = \frac{F_0}{1-e^{j\theta z}} + \text{terms of the form } \frac{B_i}{1-p_i z}$ , consider  $\left[(1-e^{j\theta z})Y(z)\right] \Big|_{z=e^{j\theta}}$ ,

$$\left[\left((e^{j\theta} z^{-1})Y(z)\right)\right] \Big|_{z=e^{j\theta}} = B_0, \left[\left((e^{j\theta} z^{-1})G(z)\right)Y(z)\right] \Big|_{z=e^{j\theta}} = G(e^{j\theta}) \rightarrow F_0 = G(e^{j\theta})$$

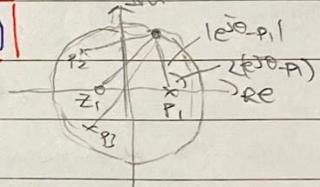
- Taking the inverse z-transform,

$$Y_k = G(e^{j\theta})e^{j\theta k} + \text{terms decaying to 0 as } k \rightarrow \infty.$$

$$\therefore \text{The stationary op is } Y_k = G(e^{j\theta})e^{j\theta k} = |G(e^{j\theta})| e^{j(\theta k + \angle G(e^{j\theta}))}$$

- Now consider the real sinusoidal i/p  $x_r = \cos(\theta k) = \frac{1}{2}(e^{j\theta k} + e^{-j\theta k})$ . Applying linearity,

$$Y_k = \frac{|G(e^{j\theta})|}{2} \left( e^{j(\theta k + \angle G(e^{j\theta}))} + e^{-j(\theta k + \angle G(e^{j\theta}))} \right) = |G(e^{j\theta})| \cos(\theta k + \angle G(e^{j\theta}))$$



Bode diagrams

- A bode diagram plots the filter magnitude  $|G(e^{j\theta})|$  and phase shift  $\angle G(e^{j\theta})$  at each frequency  $\theta$ .

- Consider a transfer function of the form  $G(z) = C \frac{\prod_{k=1}^m (z-p_k)}{\prod_{k=1}^n (z-z_k)}$

$$\hookrightarrow |G(e^{j\theta})| = |C| \frac{\prod_{k=1}^m |e^{j\theta} - z_k|}{\prod_{k=1}^n |e^{j\theta} - p_k|} \quad \hookrightarrow |G(e^{j\theta})|_{dB} = 20\log(|G(e^{j\theta})|) = 20\log(|C|) + \sum_{k=1}^m \log|e^{j\theta} - z_k| - \sum_{k=1}^n \log|e^{j\theta} - p_k|$$

i.e.,  $|G(e^{j\theta})| = \text{product of distances from zeros to } e^{j\theta} \text{ divided by product of distances from poles to } e^{j\theta}$

(in dB, the products/divisions become sums/subtractions)

$$\hookrightarrow \angle G(e^{j\theta}) = \angle C + \sum_{k=1}^m \angle(e^{j\theta} - z_k) - \sum_{k=1}^n \angle(e^{j\theta} - p_k)$$

The  $\angle G(e^{j\theta})$  = sum of angles from zeros to  $e^{j\theta}$  minus sum of angles from poles to  $e^{j\theta}$ .

- When  $e^{j\theta}$  is close to a pole, the magnitude of the response rises (resonance)

when  $e^{j\theta}$  is close to a zero, the magnitude of the response falls (null)

- since  $e^{j\theta} = e^{j(\theta k)}$ , going more than a complete revolution is redundant

- For real-valued signals, the conjugate symmetric property of z-transform means  $G(p e^{j\theta}) = \overline{G(p e^{j\theta})}$ ,

i.e.  $|G(p e^{j\theta})| = |G(p e^{j\theta})|$  and  $\angle G(p e^{j\theta}) = -\angle G(p e^{j\theta})$ , so we frequencies mirror the frequencies.

- therefore, bode diagrams are generally limited to  $\theta \in [0, \pi]$  and  $\omega \in [0, \omega_2]$ .

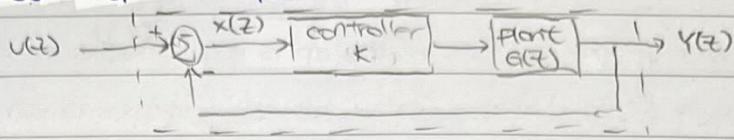
- The gain and phase margins can be read off the bode diagram as in the continuous time case

# For Personal Use Only -bkwk2

## Nyquist stability criterion

### Closed-loop control systems

- Consider the closed-loop control system below.



where we will assume the open loop system is causal (but not necessarily assume it is stable)

$$x(z) = U(z) - Y(z); \quad Y(z) = KG(z)x(z) \rightarrow Y(z) = \frac{U(z) \cdot KG(z)}{1 + KG(z)}$$

$$\therefore H(z) = \frac{Y(z)}{U(z)} = \frac{KG(z)}{1 + KG(z)}$$

- \* same form as for continuous time signals.

## Nyquist diagrams

- The Nyquist diagram of  $G$  is  $G(e^{j\theta})$  plotted on the complex plane for  $\theta \in [-\pi, \pi]$ . To draw the Nyquist plot,
  - ↳ Read from the Bode diagram, obtain the real/imaginary pair using the gain/phase.
  - ↳ Substitute various value of  $\theta$  into  $G(e^{j\theta})$  and plot directly (try values like  $\theta=0, \theta=\pi/2, \theta=\pi$ ).
- If there are open-loop poles on the unit circle, the Nyquist trajectory,  $G(e^{j\theta})$  will go to infinity as it passes through the pole on the unit circle.
- For real-valued systems, the conjugate symmetry property of z-transforms gives  $G(e^{j\theta}) = \overline{G(e^{j\theta})}$ , thus to find a complete Nyquist diagram, we only need to plot  $G(e^{j\theta})$  for  $0 \leq \theta \leq \pi$ , then draw in the complex conjugate locus.

## Nyquist stability criterion.

- The Nyquist stability criterion allows us to predict the stability of the closed loop  $H(z) = \frac{KG(z)}{1 + KG(z)}$  from the Nyquist diagram of the open loop  $G$
- The Nyquist stability criterion states that the closed loop system is stable iff the no. of encirclements of  $-1/k$  by  $G(e^{j\theta})$  as  $\theta$  increases from  $-\pi$  to  $\pi$  equals the no. of open loop unstable poles  $N_{op,0}$ .
- Our goal is to vary  $k$  to shift the pt  $-1/k$  on the real axis s.t. we get a stable closed loop system.
- The interpretation is the same as the continuous-time Nyquist diagram - gain and phase margins can be read off in the same way and have the same meanings.

can encircle

# For Personal Use Only -bkwk2

Encirclements of the origin by  $F(e^{j\theta})$ .

- consider the rational function  $F(z) = A \frac{\prod_{k=1}^M (z-z_k)}{\prod_{k=1}^N (z-p_k)}$ . The phase of  $F(e^{j\theta})$  is given by.

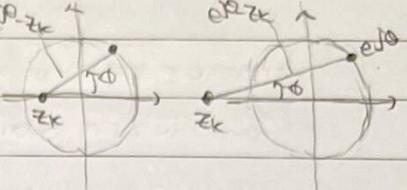
$$\angle F(e^{j\theta}) = \angle A + \sum_{k=1}^M \angle(e^{j\theta} - z_k) - \sum_{k=1}^N \angle(e^{j\theta} - p_k)$$

- As  $\theta$  increases from 0 to  $2\pi$  (or equivalently  $-\pi$  to  $\pi$ ),  $e^{j\theta} - z_k$

(i)  $|z_k| < 1$ , then the increase in  $\angle(e^{j\theta} - z_k)$  is  $2\pi$

(ii)  $|z_k| > 1$ , then the increase in  $\angle(e^{j\theta} - p_k)$  is 0.

\* similar for poles but neg effect



- Let  $N_{zi}/N_{pi}$  be the no. of zeros/poles of  $F(z)$  inside the unit circle

- The overall increase  $\alpha$  of  $\angle F(e^{j\theta})$  as  $\theta$  increases from 0 to  $2\pi$  is

$$\alpha = 2\pi(N_{zi} - N_{pi})$$

The no. of encirclements  $C$  of the origin by the Nyquist trajectory  $F(e^{j\theta})$  as  $\theta$  increased from 0 to  $2\pi$  is

$$C = N_{zi} - N_{pi}$$

Encirclements of the origin by  $1+KG(z)$ .

- The open loop transfer function is  $H_o(z) = G(z) = \frac{b(z)}{a(z)}$   $\rightarrow$  open loop poles are solns to  $a(z)=0$

The closed loop transfer function is  $H_c(z) = \frac{KG(z)}{1+KG(z)}$   $\rightarrow$  closed loop poles are solns to  $1+KG(z) = \frac{a(z)+kb(z)}{a(z)} = 0$ .

- Consider  $1+KG(z) = \frac{a(z)+kb(z)}{a(z)}$

$\hookrightarrow$  The zeros of  $1+KG(z)$  (soln to  $a(z)+kb(z)=0$ ) are the closed loop poles

$\hookrightarrow$  The poles of  $1+KG(z)$  (soln to  $a(z)=0$ ) are the open loop poles.

- The no. of zeroes/poles of  $1+KG(z)$  inside the unit circle,  $N_{zi}/N_{pi}$  is equivalent to the no. of

closed/open loop poles inside the unit circle  $N_{cp,i}/N_{op,i}$ , i.e.

$$N_{zi} = N_{cp,i}$$

$$N_{pi} = N_{op,i}$$

- The no. of encirclements  $C$  of the origin by the Nyquist trajectory  $1+KG(e^{j\theta})$  as  $\theta$  increases from 0 to  $2\pi$  is

$$C = N_{zi} - N_{pi} = N_{cp,i} - N_{op,i}$$

Pole arithmetic

- Let  $m/n$  be the degree of  $b(z)/a(z)$ . Since the open-loop system  $G(z) = \frac{b(z)}{a(z)}$  is causal,  $m \leq n$ .

Note that  $\deg(a(z) + kb(z)) = \max(\deg(a(z)), \deg(kb(z))) = \max(n, m) = n$ .

Therefore,  $1+KG(z) = \frac{a(z)+kb(z)}{a(z)}$  is a fraction of two polynomials of degree  $n \rightarrow$  both have  $n$  roots  $\rightarrow$  no. of zeros = no. of poles =  $n$ .

- Let the subscript cp/op denote closed/open loop poles, z/p denote zeroes/poles and

i/o denote inside/outside the unit circle. [As before,  $N_{op} = N_z$ ,  $N_{cp} = N_p$ ]

$$n = \underline{N_{pi}} + \underline{N_{cp,o}} = \underline{N_{op,i}} + \underline{N_{cp,o}}$$

$\therefore$  The no. of encirclements  $C$  satisfies  $C = N_{cp,i} - N_{op,i} = (n - N_{cp,o}) - (n - N_{op,o}) = N_{op,o} - N_{cp,o}$ .

# For Personal Use Only -bkwk2

Encirclements of the point  $-1/k$  by  $G(e^{j\theta})$ .

- When  $1+KG(e^{j\theta})$  encircles the origin, equivalently,
- $\hookrightarrow KG(e^{j\theta})$  encircles the pt.  $-1$
- $\hookrightarrow G(e^{j\theta})$  encircles the pt.  $-1/k$ .

- The no. of encirclements  $C$  of the pt.  $-1/k$  by the Nyquist trajectory  $G(e^{j\theta})$  as  $\theta$  increases from 0 to  $2\pi$  is

$$C = N_{op,0} - N_{p,0}$$

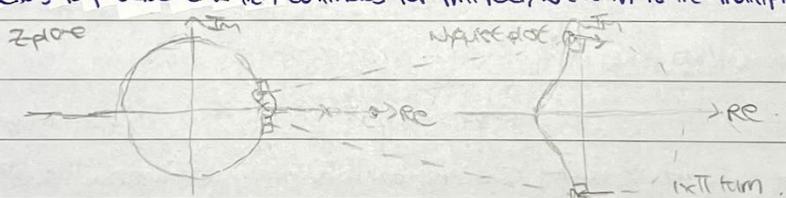
- The closed loop system is stable if  $N_{op,0} = 0$ , i.e. no closed loop poles outside the unit circle.
- therefore, for closed loop stability, we req.

$$\boxed{C = N_{op,0} = N_{p,0}}$$

as stated in the Nyquist stability criterion.

Open loop poles on the unit circle

- The Nyquist trajectory  $G(e^{j\theta})$  goes to infinity as it passes an open-loop pole on the unit circle, so to correctly count the no. of encirclements, we need to detach a closed curve for the Nyquist locus.
- To close the Nyquist trajectory, we indent the path of  $z$  around the poles on the unit circle w/ a small semi-circular excursion outside the unit circle.
- This means open loop poles on the unit circle are counted as stable in the stability criterion.
- i.e.  $N_{op,0} = N_{p,0}$  is the no. of open loop poles strictly outside the unit circle.
- The "right turn" in the  $z$ -plane then gives a "right turn" in the Nyquist plot, and a circular arc of large radius is produced which continues for  $m\pi$  rad, where  $m$  is the multiplicity of the unit circle pole.



If  $m=1$ , the Nyquist locus will be asymptotic to a straight line as it tends to infinity.

$\hookrightarrow$  suppose that  $G(z)$  has a pole at  $z=1$ , i.e.  $G(z) = \frac{F(z)}{z-1}$ , where  $F(z)$  has no poles/zeros at  $z=1$ .

Taylor expanding about  $z=1$ ,  $G(z) = \frac{F(z)}{z-1} = \frac{1}{z-1} [F(1) + F'(1)(z-1) + \dots] \approx \frac{F(1)}{z-1} + F'(1)$

Note that  $\frac{1}{z-1} = \frac{1}{e^{j\theta/2}(e^{j\theta/2}-e^{-j\theta/2})} = \frac{e^{j\theta/2}}{e^{j\theta/2}-e^{-j\theta/2}} = \frac{\cos(\theta/2)-j\sin(\theta/2)}{2j\sin(\theta/2)} = -\frac{1}{2} - \frac{j}{2\sin(\theta/2)}$

$\therefore G(e^{j\theta}) = F'(1) + \frac{F(1)}{e^{j\theta/2}} = F'(1) - \frac{1}{2}F(1) - \frac{jF(1)}{2\sin(\theta/2)} \rightarrow \boxed{G(e^{j\theta}) \rightarrow F'(1) - \frac{1}{2}F(1)} \text{ as } \theta \rightarrow 0$

$\hookrightarrow$  consider  $G(e^{j\theta})$  directly, then Taylor expand the  $e^{j\theta}$  terms,  $e^{j\theta} = 1+j\theta + \frac{(j\theta)^2}{2!} + \dots$

Take out a factor of  $j\theta$  from the  $(e^{j\theta}-1)$  term, then expand to get the form  $\frac{1}{j\theta} (a+bj\theta+cj\theta^2)$

The x-coord of the asymptote is  $b$ .

$\rightarrow$  May need to use the binomial expansion  $(1+k\theta)^{-1} = 1-k\theta + O(\theta^2)$

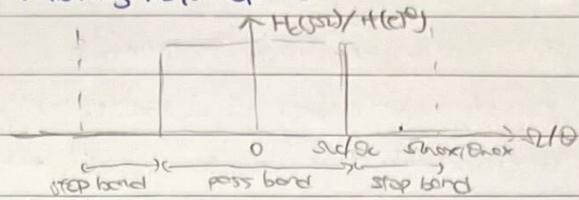
## Digital filters

# For Personal Use Only -bkwk2

### Design of filters

#### Ideal filters

- The ideal LPF has the following freq. response.



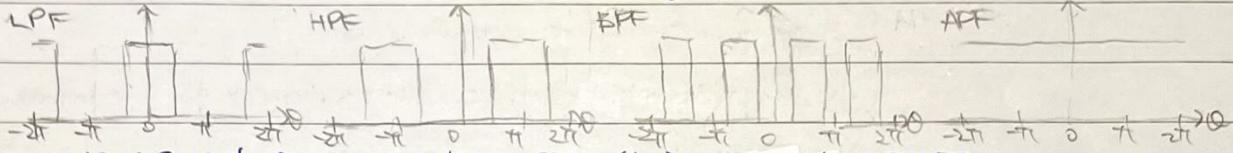
$$\text{Analog: } H(j\omega) = \begin{cases} 1 & \text{if } |\omega| \leq \omega_c \\ 0 & \text{if } |\omega| > \omega_c \end{cases}$$

$$\text{Digital: } H(e^{j\theta}) = \begin{cases} 1 & \text{if } |\theta| \leq \theta_c \\ 0 & \text{if } |\theta| > \theta_c \end{cases}$$

where  $\theta = \omega T$  is the normalised freq. and  $T = \frac{1}{f_s} = \frac{\pi}{\omega_s}$  is the sampling freq.

$$\text{and } \theta_c \leq \theta_{max} = \pi / \Delta\omega_c \leq \Delta\omega_{max} = \frac{\omega_s}{2}$$

- similarly, the ideal HPF/BPF/APF have the following freq. response.



- consider the ideal LPF. Its impulse response {h\_k} is found using the DTFT.

$$h_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} H(e^{j\theta}) e^{-jk\theta} d\theta = \frac{1}{2\pi} \int_{-\omega_c}^{\omega_c} e^{-jk\theta} d\theta = \frac{1}{2\pi} \frac{e^{-jk\omega_c} - e^{jk\omega_c}}{jk} = \frac{\sin(\theta_c k)}{\pi k} = \frac{\omega_c}{\pi} \sin(\theta_c k)$$

i.e. it is non-causal and has a two-sided infinite impulse response  $\rightarrow$  not realisable.

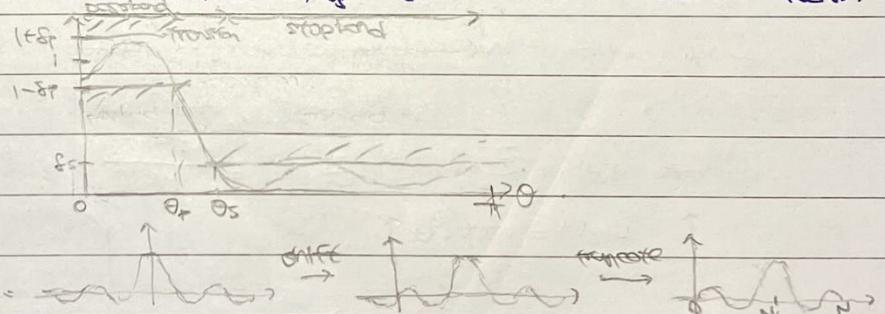
#### Realisable filter specifications

- A typical filter specification must specify the max. permissible deviations from the ideal, including

(i) band edge freq. ( $\omega_p, \omega_s$ ) (ii) max. passband ripple ( $1-f_p, 1+f_p$ ), (iii) min. stopband attenuation ( $\alpha, \delta_s$ ).

- Using the dB scale, we need to specify (i) passband ripple ( $20\log_{10}(1+f_p) \text{ dB}$ ), (ii) peak-to-peak passband ripple ( $20\log_{10}(1+2f_p) \text{ dB}$ ), (iii) min. stopband attenuation ( $-20\log_{10}(\delta_s) \text{ dB}$ ).

( $f_p = 0.06$  gives peak-to-peak passband ripple =  $2 \text{ dB}$ ,  $\delta_s = 0.01$  gives min. stopband attenuation =  $40 \text{ dB}$ ).



#### FIR filters.

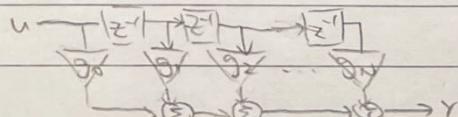
- The FIR filter G is derived by shifting + truncating at  $N+1$  samples of the ideal response H, i.e.

$$G(z) = \sum_{k=0}^N g_k z^{-k}, \text{ where } g_k = h_k \cdot w_k$$

and  $h_k$  is the non-causal impulse response of the ideal filter.

- FIR filters are feed-forward / non-recursive  $\rightarrow$  efficiently realised in hardware (FFT).

$$y_n = \sum_{k=0}^N g_k h_k n - k$$



- FIR filters are inherently stable.  $G(z) = \frac{\sum g_k z^{N-k}}{z^N} \rightarrow \text{all poles at } z=0 \text{ (inside unit circle)}$

# For Personal Use Only -bkwk2

## IIR filters

- IIR filters have a finite polynomial representation but infinite impulse response.

- IIR filter has feedback — the output "goes back" to the filter.

$$y_k = \underbrace{-a_1 y_{k-1} - \dots - a_n y_{k-n}}_{\text{feedback}} + b_0 x_k + b_1 x_{k-1} + \dots + b_n x_{k-n}$$

- IIR filters typically meet a given set of specifications w/ a much lower filter order than their FIR filter counterpart → used for digital implementation of analog filters.

## FIR vs IIR filters

- FIR filters: ✓ simple, stable, feedforward, fast implementation (FFT)

✓ can have exactly linear phase  $G(e^{j\omega}) = |G(e^{j\omega})| e^{j\phi}$

✓ design methods are generally linear

✗ higher filter order than IIR filter for some level of performance

✗ often greater delay than IIR filter w/ equal performance

- IIR filters: ✓ lower filter order than FIR filter for some level of performance

✓ shorter delay than FIR filter w/ equal performance.

✗ can be unstable, uses feedback.

## Design of FIR filters

Frequency distortion of truncation/windowing.

- To get the truncated unit sample response  $g_t$  from the shifted desired impulse, we effectively multiply by a rectangular window.

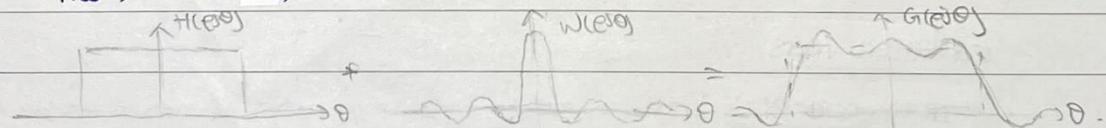
$$g_t = h_k \quad 0 \leq k \leq N \quad \Rightarrow \quad g_t = h_k w_k, \quad w_k = \begin{cases} 1 & 0 \leq k \leq N \\ 0 & \text{o.w.} \end{cases}$$

- Using duality of multiplication/convolution for Fourier transforms

$$G(e^{j\omega}) = \sum_{n=0}^N \int_{-\pi}^{\pi} H(e^{j\theta}) W(e^{j(\theta-\omega)}) d\theta$$

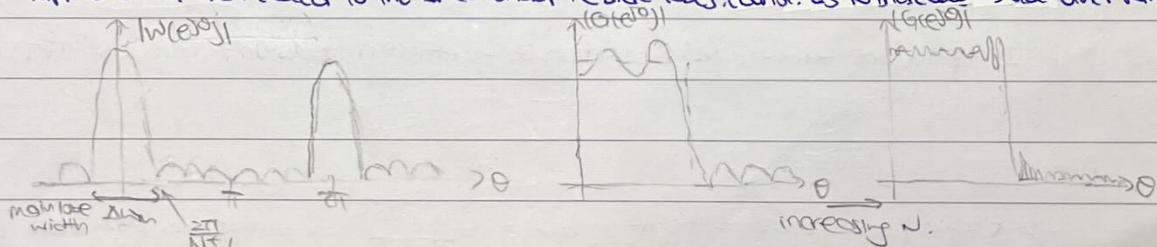
$$\text{where } W(e^{j\omega}) = \sum_{k=0}^N w_k e^{j\omega k} = \sum_{k=0}^N e^{j\omega k} = \frac{1 - e^{j\omega(N+1)}}{1 - e^{j\omega}} = e^{-j\omega \sum} \frac{\sin(\frac{\omega(N+1)}{2})}{\sin(\omega/2)}$$

$$H(e^{j\omega}) = \text{rect}\left(\frac{\omega}{2\pi}\right)$$



- The transition band is related to the main lobe. (Reduces as N → ∞)

Ripples of G are related to the area under the side lobes. (cont. as N increases → use other windows)



# For Personal Use Only -bkwk2

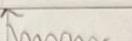
Design by window method.

- We can reduce freq. distortion by adopting diff. windows, common windows include.

↳ Rectangular

$$w_k = \begin{cases} 1 & 0 \\ 0 & \text{else} \end{cases}$$

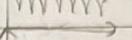
$$0 \leq k \leq N$$



↳ Triangular (Bartlett)

$$w_k = \begin{cases} \frac{2k/N}{N-2k/N} & 0 \leq k \leq N/2 \\ 0 & \text{else} \end{cases}$$

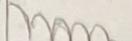
$$0 \leq k \leq N/2$$



↳ Hanning

$$w_k = \begin{cases} 0.5 + 0.5 \cos(\pi k/N) & 0 \leq k \leq N \\ 0 & \text{else} \end{cases}$$

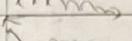
$$0 \leq k \leq N$$



↳ Hamming

$$w_k = \begin{cases} 0.54 - 0.46 \cos(\pi k/N) & 0 \leq k \leq N \\ 0 & \text{else} \end{cases}$$

$$0 \leq k \leq N$$



- The shape and duration of the window to control the resulting filter properties:

↳ Smaller side lobes yield better approx. of the ideal response

↳ Narrow transition bandwidth for increasing  $N$ , symmetric at cut-off  $\Omega_c$ .

↳ Some  $\delta$  for passband error and stopband approx.

	Peak side lobe amplitude (relative)	Aprox. of main lobe	Peak approx. error (dB)	$H(e^{j\theta})$	$G(e^{j\theta})$
Rectangular	-13	$\frac{4\pi}{N+1}$	-21		
Bartlett	-25	$\frac{4\pi}{N}$	-25		
Hanning	-31	$\frac{8\pi}{N}$	-44		
Hamming	-41	$\frac{8\pi}{N}$	-53		

$$G(e^{j\theta}) = \sum_{k=0}^N H(e^{j\theta}) w_k e^{-j\theta k} d\theta$$

- The window method can be used to quickly design filters to approx. a given target response.

- It does not impose amplitude response constraints (passband ripple, stopband attenuation) → need to be used iteratively to produce designs which meet such specs.

- To design a filter using the window method:

↳ 1) Select a suitable window function  $w_k$ , i.e. Hanning/Bartlett etc.

↳ 2) Specify an ideal freq. response  $H$ , i.e. LPF/BPF etc.

↳ 3) Compute the coeff. for the ideal filter,  $h_k$

↳ 4) Delay the ideal filter coeff  $h_k$  by  $N/2$ ,  $h_{k-N/2}$ , then multiply by the window function  $w_k$

↳ 5) Evaluate the freq. response of the resulting filter and repeat the above steps iteratively if needed.

## Low pass filter design (FIR)

-ej: Design a LPF w/ passband  $\theta_p = 0.2\pi$ , stopband  $\theta_s = 0.3\pi$  and approx. error  $\delta = 0.01$

① Max. allowed approx. error in dB =  $20 \log_{10} \delta = -40 \text{ dB} > -44 \text{ dB}$  → use Hanning window.

② cut-off freq.  $\Omega_c = \frac{\theta_p + \theta_s}{2} = 0.25\pi$  main lobe width  $\Delta\theta = \theta_p - \theta_s = 0.1\pi = \frac{\pi}{N} \rightarrow N = 80$

③ Ideal filter  $h_k = \frac{1}{2\pi} \int_{-\theta_c}^{\theta_c} H(e^{j\theta}) e^{j\theta k} d\theta = \frac{1}{2\pi} \int_{-\theta_c}^{\theta_c} e^{j\theta k} d\theta = \frac{1}{2\pi} \frac{e^{j\theta_c k} - e^{-j\theta_c k}}{j} = \frac{\sin(\theta_c k)}{\pi k} = \frac{\Omega_c}{\pi} \sin(\Omega_c k)$

④ Filter coeff  $g_k = h_k - \frac{1}{2} w_k$  =  $\frac{1}{4} \sin\left(\frac{\pi}{4}(k-40)\right) - \frac{1}{2}(1 - \cos(\frac{\pi k}{40}))$   $k \in [0, 80]$

# For Personal Use Only -bkwk2

Multi band FIR design .

- We can design HPF/BPF using composition of LPF

$$H_{\text{HPF}}[\theta, \pi](e^{j\theta}) = H_{\text{LPF}}[0, \pi] - H_{\text{LPF}}[0, \theta_1] \quad H_{\text{BPF}}[\theta_1, \theta_2](e^{j\theta}) = H_{\text{LPF}}[\theta_1, \theta_2] - H_{\text{LPF}}[0, \theta_1]$$

where  $H_{\text{LPF}}[0, \theta_1] = \begin{cases} 1 & 0 \leq \theta \leq \theta_1 \\ 0 & \theta > \theta_1 \end{cases}$  and its corresponding  $h_{\text{LPF}} = \frac{\theta_1}{\pi} \sin(\theta_1 k)$

- The ideal impulse response for HPF/BPF is found by linearity of the antitransform,

$$h_{\text{LPF}} = \sin(\pi k) - \frac{\theta_1}{\pi} \sin(\theta_1 k)$$

$$h_{\text{BPF}} = \frac{\theta_2}{\pi} \sin(\theta_2 k) - \frac{\theta_1}{\pi} \sin(\theta_1 k)$$

- The final filter is thus obtained as usual, using .

$$g_k = h_k w_k$$

Linear phase

- A linear phase filter has freq. response  $G(e^{j\theta})$  and satisfies

$$G(e^{j\theta}) = |G(e^{j\theta})| e^{-j\theta \sum}$$

$$\exists k \rightarrow G(z), \exists k \in \mathbb{Z} \rightarrow z^k G(z) \\ \therefore \text{for } z = e^{j\theta}, \text{ we have } G(z) e^{-j\sum}$$

i.e. the phase response is a linear function of freq. ( $\phi(\theta) = \frac{\sum}{2}\theta$ ) .

- For a linear phase filter, all the freq. component of the i/p signal are shifted in time by  $\frac{\sum}{2}$   
convolve  $\rightarrow$  no phase distortion in the o/p signal (i/p and o/p signals have to same shape)
- Linear phase is achieved if the filter coeff  $g_k$  is symmetric about  $S_N/2$ , i.e.

$$g_k = g_{N-k}$$

$$\begin{aligned} G(e^{j\theta}) &= \sum_{k=0}^N g_k e^{jk\theta} = g_0 e^{j0\theta} + g_1 e^{j1\theta} + g_2 e^{j2\theta} + g_3 e^{j3\theta} + \dots \\ &= e^{j\theta \frac{N}{2}} (g_0 e^{j0\frac{N}{2}} + g_1 e^{j1\frac{N}{2}} + g_2 e^{j2\frac{N}{2}-1} + g_3 e^{j3\frac{N}{2}-1} + \dots) \\ &= e^{j\theta \frac{N}{2}} (2g_0 \cos(0\frac{N}{2}) + 2g_1 \cos(\theta(\frac{N}{2}-1)) + \dots) \end{aligned}$$

→ Note that the window method gives linear phase filters ( $w_k = w_{N-k}$  by defn, so if the desired  $h_k$  is symmetric, i.e.  $h_k = h_{N-k}$ , then  $g_k = h_k w_k$  is also symmetric → linear phase).

Design by optimisation.

- Given the ideal filter  $H$  and the weighting function  $W$ , we can find the optimal filter  $G$  (and its coefficients  $g_k$ ) of length  $N$  by optimisation .

- We define a cost/error function  $E(\theta)$

$$E(\theta) = W(\theta) [H(\theta) - G(\theta)]$$

- $G$  can be found by minimising

(i) least squares error  $\int_{-\pi}^{\pi} E^2(\theta) d\theta$  i.e.  $G(\theta) = \arg \min_{G(\theta)} \int_{-\pi}^{\pi} E^2(\theta) d\theta$

(ii) max. error  $\max_{-\pi \leq \theta \leq \pi} |E(\theta)|$  i.e.  $G(\theta) = \arg \min_{G(\theta)} \max_{-\pi \leq \theta \leq \pi} |E(\theta)|$

equiripple filter .

# For Personal Use Only -bkwk2

## Design of IIR filters

Discretization by response matching.

- We can convert a continuous time filter  $G_c(s)$  into an equivalent discrete time filter  $G(z)$  by matching their response to a certain input (impulse, step, ramp).
- For a certain input, the response is matched if the dp y satisfies

$$\boxed{L^{-1}[Y(s)]|_{t=kT} = y(kT) = y_k = z^k [Y(z)]}$$

### ① Impulse invariance

$$G_c(s) \xrightarrow{\text{L}^{-1}} g(t) \xrightarrow{\text{sample}} g(kT) = y_k \xrightarrow{z^k} G(z)$$

$$\therefore G(z) = z [L^{-1}[G_c(s)]|_{t=kT}]$$

→ For bandlimited filters, the digital filter freq. response closely approx. the continuous freq. response

↳ Impulse invariance preserves stability, i.e.  $\text{Re}[p_a] < 0 \rightarrow |e^{j\omega T}| < 1$

### ② Step invariance

$$G_c(s) \xrightarrow{\text{step}} \frac{G_c(s)}{s} \xrightarrow{\text{L}^{-1}} y(t) \xrightarrow{\text{sample}} y(kT) = y_k \xrightarrow{z^k} Y(z) \xrightarrow{\text{step}} \frac{z^k}{z-1} Y(z)$$

$$\therefore G(z) = \frac{z^k}{z-1} z [L^{-1}\left(\frac{G_c(s)}{s}\right)|_{t=kT}]$$

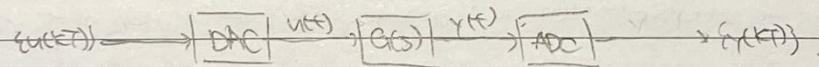
### ③ Ramp invariance

$$G_c(s) \xrightarrow{\text{ramp}} \frac{G_c(s)}{s^2} \xrightarrow{\text{L}^{-1}} y(t) \xrightarrow{\text{sample}} y(kT) = y_k \xrightarrow{z^k} Y(z) \xrightarrow{\text{ramp}} \frac{(z^k)^2}{z^2} Y(z)$$

$$\therefore G(z) = \frac{(z^k)^2}{z^2} z [L^{-1}\left(\frac{G_c(s)}{s^2}\right)|_{t=kT}]$$

choice of waveform for response matching.

- The choice of waveform to achieve response invariance depends on the hybrid analog/digital system.



- Within the hybrid system, it is the properties of the DAC and ADC that influence the choice of waveform for response invariance.

↳ Sample + hold DAC → step invariance.

↳ Moving average DAC → ramp invariance

# For Personal Use Only -bkwk2

Discretization by algebraic transformation

- Algebraic transformations are used to map systems that are not necessarily bandlimited or the cost of some distortion b/w the freq. response of the continuous/discrete time system.

$$H(z) = H(s) \Big|_{s=\gamma(z)}$$

where  $\gamma(z)$  is given by

Forward difference (Euler)

$$s = \frac{z-1}{T} \quad \tilde{x} = \frac{x(t+T) - x(t)}{T} \rightarrow sX(s) = \frac{z-1}{T} X(z)$$

Backward difference

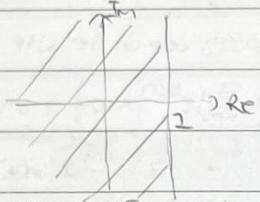
$$s = \frac{z-1}{zT} \quad \tilde{x} = \frac{x(t) - x(t-T)}{T} \rightarrow sX(s) = \frac{1-z^{-1}}{T} X(z)$$

Bilinear transform (Tustin)

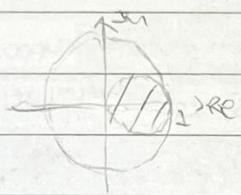
$$s = \frac{z-1}{T(z-1)} \quad z = e^{sT/2} = \frac{e^{sT/2} - 1}{e^{sT/2} + 1} = \frac{1 + sT/2}{1 - sT/2} = \frac{z - 1}{z + 1} \rightarrow s = \frac{z-1}{T(z+1)}$$

- Each of these transformations corresponds to a certain mapping b/w s-plane and z-plane.

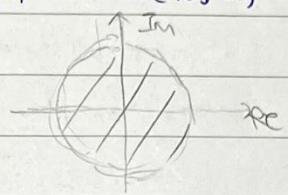
The shaded region below in the z-plane corresponds to the LHP of the s-plane (stable region)



forward difference  
(LHP shifted)



backward difference  
(LHP mapped to outer disk)



bilinear transform  
(LHP mapped to unit disk)

- Applying algebraic transformation on a stable continuous system results in a stable discrete time system for backward diff/bilinear transform, but is not necessarily true for forward diff.

## Bilinear transform

- consider the bilinear transform  $s = \gamma(z) = \frac{z-1}{z+1}$

- Analog filter  $G_c(s)$  w/ freq. resp.  $G_c(j\omega)$  is Digital filter  $G(z)$  w/ freq. resp.  $G(e^{j\theta})$ .

Using the algebraic transformation  $\theta(z) = G_c(\gamma(z))$ , if freq. resp. of the digital filter for  $|z| < 1$  is

$$G(e^{j\theta}) = G_c(j\operatorname{tan}(\theta/2)) = G_c(j\omega) \quad \text{since } \gamma(e^{j\theta}) = \frac{e^{j\theta}-1}{e^{j\theta}+1} = \frac{e^{j\theta/2}(e^{j\theta/2}-e^{-j\theta/2})}{e^{j\theta/2}(e^{j\theta/2}+e^{-j\theta/2})} = j\operatorname{tan}(\theta/2)$$

- This means the freq. resp. of the analog filter at  $\omega$  is mapped to the freq. resp. of the digital filter at  $\theta$ , where  $\omega$  and  $\theta$  are related by

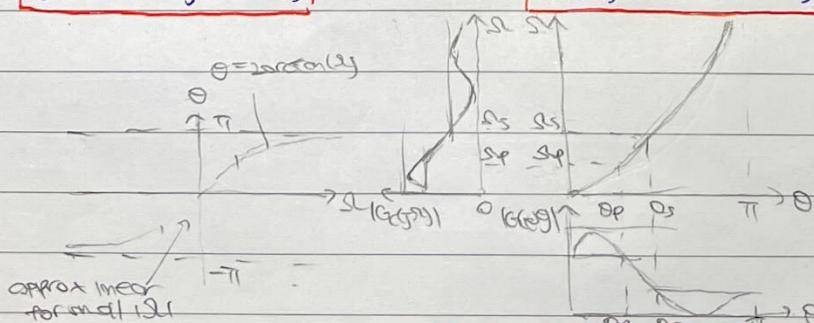
$$\omega_c = \operatorname{tan}(\theta/2)$$

$$\theta = 2\operatorname{atan}(\omega_c)$$

- Therefore, we can map b/w the discrete/continuous time filters using

$$G(e^{j\theta}) = G_c(j\operatorname{tan}(\theta/2))$$

$$G_c(j\omega) = G(e^{j2\operatorname{atan}(\omega_c)})$$



- \* Using the bilinear transform  $s = \gamma(z) = \frac{z-1}{z+1}$ ,  $\omega_c = \frac{\pi}{T} \operatorname{tan}(\theta/2)$ ,  $\theta = 2\operatorname{atan}(\omega_c)$ , and the mapping is

$$G(e^{j\theta}) = G_c(j\frac{z-1}{z+1} \operatorname{tan}(\theta/2))$$

$$G_c(j\omega) = G(e^{j2\operatorname{atan}(\omega_c)})$$

# For Personal Use Only -bkwk2

- consider the bilinear transform  $s = \gamma(z) = \frac{z-1}{z+1} \rightarrow z = \gamma(s) = \frac{1+s}{1-s}$   
 For  $s = \lambda + j\omega$ ,  $|z|^2 = z\bar{z} = \gamma(s)\gamma(s)^* = \frac{1+\lambda+j\omega}{1-\lambda+j\omega} \cdot \frac{1+\lambda-j\omega}{1-\lambda-j\omega} = \frac{(\lambda+1)^2 + \omega^2}{(\lambda-1)^2 + \omega^2}$

If  $\lambda = 0$ , then  $|z|^2 = \frac{1+\omega^2}{1+\omega^2} = 1$ , i.e. on unit circle

If  $\lambda < 0$ , then  $|z|^2 = \frac{(\lambda+1)^2 + \omega^2}{(\lambda-1)^2 + \omega^2} < 1$ , i.e. inside the unit circle.

→ The LHP in the s-plane is mapped exactly to the unit circle in the z-plane → stability preserved.

## Low pass filter design (IIR)

- consider the first order analog LPF  $G_a(s) = \frac{1}{1+s/\omega_c}$ , which has gain 1 at  $s=0$  and gain = 0.5 at  $s=j\omega_c$ . Applying the bilinear transform  $s = \gamma(z) = \frac{z-1}{z+1}$ ,

$$G(z) = G_a(\gamma(z)) = \frac{1}{1 + \gamma(z)/\omega_c} = \frac{1}{1 + \frac{z-1}{z+1}/\omega_c} = \frac{(z+1)\omega_c}{(z+1)\omega_c + z-1} = \frac{(z+1)\omega_c}{(z+1)z + (\omega_c - 1)} = \frac{(z+1)\omega_c}{z + \frac{\omega_c - 1}{z+1}}$$

- Given a desired cutoff freq. for the digital filter  $\omega_c$ , we can find the equivalent pre-winded analog filter cut-off freq.  $\omega_a$  using  $\omega_a = \tan(\theta_c/2)$ .

- e.g. Design a first order digital LPF w/ -3dB freq. of  $1000\pi$  rad/s and a sampling freq. of  $8000\pi$  rad/s.

We are given  $f_c = 1000\pi$ ,  $f_s = 8000\pi \rightarrow \omega_c = 2\pi f_c$ ,  $T = \frac{1}{f_s} \rightarrow \theta_c = \omega_c T = \frac{2\pi f_c}{f_s} = \frac{\pi}{4}$  rad/sample

$$\therefore \omega_a = \tan\left(\frac{\pi/4}{2}\right) = \tan(\pi/8) = 0.444 \text{ rad/s} \rightarrow G(z) = \frac{(z+1)\omega_c}{z + \frac{\omega_c - 1}{z+1}} = \frac{0.293(z+1)}{z - 0.444}$$

since  $G(z) = \frac{Y(z)}{U(z)}$ , we have  $Y(z)(1 - 0.444z^{-1}) = U(z)(0.293(1+z^{-1}))$ , so the difference eqn. is

$$Y_k = 0.444Y_{k-1} + 0.293(U_k - U_{k-1})$$

\* THIS IS AN IIR filter as the transfer function  $G(z)$  has a pole at  $z \neq 0$ .

## Bond transformation

- Analog prototypes are typically LPF. Standard transformations can be used to achieve HPF/BPF/BSF.

- Assuming we have a LPF w/ cutoff at  $\omega_c = 1 \text{ rad/s}$ .

↳ LPF → LPF : 
$$\boxed{s \rightarrow \frac{\omega}{\omega_c}}$$
 ]  $\omega_c$  is the cutoff freq.

↳ LPF → HPF : 
$$\boxed{s \rightarrow \frac{j\omega}{s}}$$
 ]  $\omega_c$  is the lower cutoff freq.

↳ LPF → BPF : 
$$\boxed{s \rightarrow \frac{s + j\omega_c}{s(j\omega_c - \omega_c)}}$$
 ]  $\omega_u$  is the upper cutoff freq.

↳ LPF → BSF : 
$$\boxed{s \rightarrow \frac{s(j\omega_c - \omega_u)}{s + j\omega_u}}$$
 ]  $\omega_u$  is the upper cutoff freq.

## Discrete Fourier Transform

## Discrete Fourier Transform (DFT)

- The DFT is given by

$$X_k = \sum_{n=0}^{N-1} x_n e^{-j \frac{2\pi}{N} k n}, \quad \text{for } 0 \leq k \leq N-1$$

- since the transform maps a finite length sequence  $\{x_n\}_{n=0}^N$  to a finite length sequence  $\{X_k\}_{k=0}^{N-1}$ , we can rewrite the DFT as a matrix-vector operation.

$$\begin{bmatrix} X_0 \\ X_1 \\ X_2 \\ \vdots \\ X_{N-1} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & w_N & w_N^2 & \dots & w_N^{N-1} \\ 1 & w_N^2 & w_N^4 & \dots & w_N^{2(N-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & w_N^{N-1} & w_N^{2(N-1)} & \dots & w_N^{(N-1)(N-1)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

where  $w_N = e^{-j \frac{2\pi}{N}}$  is a primitive  $N$ th root of unity (i.e.  $w_N^N = 1$  and  $w_N^M \neq 1$  for  $M=1, 2, \dots, N-1$ )

- The  $N$  freq. readings obtained from the DFT are "sampled" freq. b/wn 0 and  $2\pi$  in normalized freq. ( $2\pi$  corresponds to sampling freq.  $f_s$ ).

↳  $X_0 \leftrightarrow$  freq. 0 ;  $X_{N/2} \leftrightarrow$  freq.  $\pi$  (i.e.  $f_s/2$ )

↳  $X_1, \dots, X_{N/2-1} \leftrightarrow$  equally spaced freq. pts b/wn 0 and  $\pi$  (i.e.  $f_s/2$ )

↳  $X_{N/2+1}, \dots, X_{N-1} \leftrightarrow$  equally spaced -ve freq. pts b/wn 0 and  $-\pi$  (i.e.  $-f_s/2$ )

Cyclic properties of the powers of  $w_N$ 

- Any integer  $k$  can be written as

$$k = qN + r, \quad \text{w/ } 0 \leq r < N$$

where  $q$  is the quotient and  $r = R_N(k)$  is the remainder

- We can therefore write  $w_N^k$  as

$$w_N^k = w_N^{qN+r} = w_N^q \cdot w_N^r = 1 \cdot w_N^r = w_N^r \rightarrow w_N^k = w_N^{R_N(k)}$$

so the DFT Vandermonde matrix  $E$  can be rewritten s.t. powers of  $w_N$  higher than  $N-1$  can be expressed as powers b/wn 0 and  $N-1$ .

Constructing the DFT Vandermonde matrix  $E$ .

- The DFT Vandermonde matrix  $E$  can be constructed as follows:

$$E = \begin{bmatrix} 1 & 1 & \dots & (\text{row 0: } N \text{ ones}) \\ 1 & w_N & \dots & (\text{row 1: powers 0 to } N-1 \text{ of } w_N) \\ 1 & w_N^2 & \dots & (\text{row 2: every 2nd element of row 1}) \\ 1 & w_N^3 & \dots & (\text{row 3: every 3rd element of row 1}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & w_N^{N-1} & \dots & (\text{row } N-1: \text{row 1 cyclically shifted}) \end{bmatrix}$$

where "every  $k$ th element" is understood cyclically.

\* Note the last row is "every  $(N+1)th$  element", which is the sole os going backwords.

# For Personal Use Only -bkwk2

Properties of the DFT

① Periodicity:  $X_n = X_{n+N}$

$$X_{n+N} = \sum_{k=0}^{N-1} X_k e^{-j\frac{2\pi}{N} k(n+N)} = \sum_{k=0}^{N-1} X_k e^{-j\frac{2\pi}{N} kn - j\frac{2\pi}{N} k^2 N} = X_n$$

② Linearity:  $\text{DFT}[x_k + y_k] = \text{DFT}[x_k] + \text{DFT}[y_k]$

$$\text{DFT}[x_k + y_k] = \sum_{k=0}^{N-1} (x_k + y_k) e^{-j\frac{2\pi}{N} kn} = \sum_{k=0}^{N-1} x_k e^{-j\frac{2\pi}{N} kn} + \sum_{k=0}^{N-1} y_k e^{-j\frac{2\pi}{N} kn} = \text{DFT}[x_k] + \text{DFT}[y_k]$$

③ Conjugate symmetry:  $X_{-n} = X_n^*$  for real  $x_k$ 's  $X_k = X_k^*$  for real  $x_n$

$$X_n = \sum_{k=0}^{N-1} X_k e^{-j\frac{2\pi}{N} k(n)} = \left( \sum_{k=0}^{N-1} X_k e^{-j\frac{2\pi}{N} kn} \right)^* = X_n^*$$

$$X_k = \frac{1}{N} \sum_{n=0}^{N-1} X_n e^{-j\frac{2\pi}{N} kn} = \left( \frac{1}{N} \sum_{n=0}^{N-1} X_n e^{-j\frac{2\pi}{N} kn} \right)^* = X_k^*$$

i.e. real for time/freq.  $\leftrightarrow$  freq/time has CS; complex for time/freq.  $\leftrightarrow$  freq/time does not have CS

④ Time/frequency shift:  $\text{DFT}[x_{R(k+d)}] = W_N^{-d} X_n$ ;  $\text{IDFT}[x_{R(k+d)}] = W_N^{kd} X_k$

$$\text{DFT}[x_{R(k+d)}] = \sum_{k=0}^{N-1} x_{R(k+d)} W_N^{kn} = \sum_{k=0}^{N-1} x_k W_N^{nR(k+d)} = \sum_{k=0}^{N-1} x_k W_N^{n(k+d)} = W_N^{-d} \sum_{k=0}^{N-1} x_k W_N^{nk} = W_N^{-d} X_n$$

$$\text{IDFT}[x_{R(k+d)}] = \frac{1}{N} \sum_{n=0}^{N-1} x_{R(n+d)} W_N^{-kn} = \frac{1}{N} \sum_{n=0}^{N-1} x_n W_N^{-k(n+d)} = \frac{1}{N} \sum_{n=0}^{N-1} x_n W_N^{-k(n+d)} = W_N^{kd} \sum_{n=0}^{N-1} x_n W_N^{-kn} = W_N^{kd} X_k$$

⑤ Cyclic convolution:  $\text{DFT}[\{x_k\} \otimes \{y_k\}] = X_n Y_n$  |  $\{z_k\} = \{x_k\} \otimes \{y_k\}$ ,  $P_k = \sum_{m=0}^{N-1} X_m Y_m R(k-m)$

$$\text{DFT}[\{x_k\} \otimes \{y_k\}] = \sum_{k=0}^{N-1} (x_k y_k) W_N^{kn} = \sum_{k=0}^{N-1} \left( \sum_{m=0}^{N-1} x_m y_{k+m} \right) W_N^{kn} = \sum_{m=0}^{N-1} x_m \sum_{k=0}^{N-1} y_k W_N^{k(m+n)} W_N^{kn}$$

$$= \sum_{m=0}^{N-1} x_m \sum_{k=0}^{N-1} y_k W_N^{k(m+n)} = \sum_{m=0}^{N-1} x_m \sum_{k=0}^{N-1} y_k W_N^{k(m+n)} = \sum_{m=0}^{N-1} x_m W_N^m \sum_{k=0}^{N-1} y_k W_N^{k(m+n)} = X_n Y_n$$

i.e. circular convolution in time domain is equivalent to point-wise multiplication in freq. domain,

The DFT Vandermonde matrix  $E$ .

- The sum of column  $m$  of the DFT Vandermonde matrix  $E$  is given by  $\sum_{k=0}^{N-1} W_N^{km}$

$$m=0: \quad \sum_{k=0}^{N-1} W_N^{km} = \sum_{k=0}^{N-1} 1 = N$$

$$m \neq 0: \quad \sum_{k=0}^{N-1} W_N^{km} = \frac{(W_N^m)^N - 1}{1 - W_N^m} = \frac{1 - W_N^{m(N)}}{1 - W_N^m} = \frac{1 - W_N^N}{1 - W_N^m} = 0$$

i.e. the sum of any column of  $E$  except column 0 is zero.

- Let  $f_l$  be the  $l$ th column of the DFT Vandermonde matrix  $E$ . For  $l > m$ ,

$$\text{complex dot product} \rightarrow f_l \cdot f_m = \sum_{k=0}^{N-1} W_N^{lk} (W_N^{mk})^* = \sum_{k=0}^{N-1} (W_N^{lmk})$$

i.e. the sum of the  $(l-m)$ th column of  $E$ , and hence 0 when  $l > m$  /  $N$  when  $l = m$

- This means the columns of  $E$  are orthogonal, and their magnitude is  $\sqrt{N}$ ,

$$f_i \cdot f_j = \begin{bmatrix} N & i=j \\ 0 & i \neq j \end{bmatrix}$$

$\therefore$  The DFT is an orthogonal basis transform (w/ scaling factor  $N$ ).

- If  $X = E x$  and  $Y = E y$ , the scaling factor of  $\sqrt{N}$  on the basis vectors give

$$x \cdot y = (\frac{1}{\sqrt{N}} X) \cdot (\frac{1}{\sqrt{N}} Y) = \frac{1}{N} X \cdot Y$$

$$\|x\|_2^2 = (\frac{1}{\sqrt{N}} X) \cdot (\frac{1}{\sqrt{N}} X) = \frac{1}{N} \|X\|_2^2$$

# For Personal Use Only -bkwk2

## INVERSE DISCRETE FOURIER TRANSFORM (IDFT)

- The IDFT is given by

$$x_k = \frac{1}{N} \sum_{n=0}^{N-1} x_n e^{-j\frac{2\pi}{N} kn} \quad \text{for } 0 \leq k \leq N-1$$

- Since the transform maps a finite length sequence  $\{x_n\}_{n=0}^{N-1}$  to a finite length sequence  $\{X_k\}_{k=0}^{N-1}$ , we can rewrite the IDFT as a matrix-vector operation.

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{N-1} \end{bmatrix} = \frac{1}{N} \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & W_N^0 & W_N^1 & \dots & W_N^{N-1} \\ 1 & W_N^{-1} & W_N^{-2} & \dots & W_N^{-N+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & W_N^{N-1} & W_N^{2(N-1)} & \dots & W_N^{(N-1)^2} \end{bmatrix} \begin{bmatrix} X_0 \\ X_1 \\ X_2 \\ \vdots \\ X_{N-1} \end{bmatrix}$$

where  $W_N = e^{-j\frac{2\pi}{N}}$  is a primitive  $N$ th root of unity (i.e.  $W_N^0 = 1$  and  $W_N^m \neq 1$  for  $m=1, 2, \dots, N-1$ )

- By inspection, the IDFT is simply DFT but  $W_N \rightarrow \frac{1}{N} W_N^{-1}$

Alternatively, we can obtain the IDFT by i/p conjugation  $\rightarrow$  DFT  $\rightarrow$  o/p conjugation  $\rightarrow$  divide by  $N$ .

$$\frac{1}{N} [DFT[x_n^*]]^* = \frac{1}{N} \left[ \sum_{n=0}^{N-1} x_n^* e^{-j\frac{2\pi}{N} kn} \right]^* = \frac{1}{N} \sum_{n=0}^{N-1} x_n e^{+j\frac{2\pi}{N} kn} = IDFT[x_n]$$

- Considering the IDFT matrix, it is simply the inverse of the DFT Vandermonde matrix  $E$ ,  $E^{-1}$ . Since  $E$  is orthogonal, it is easy to determine its inverse.

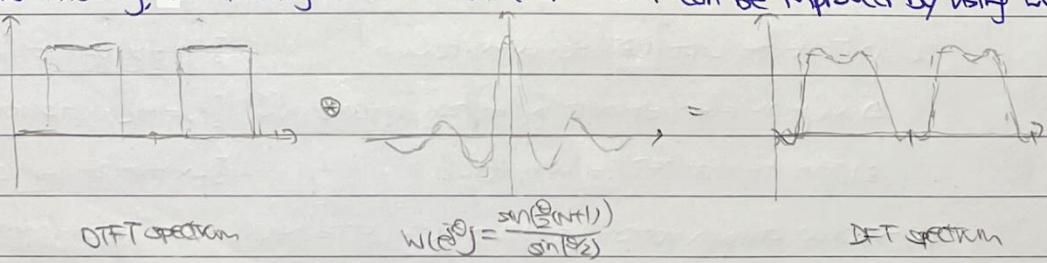
$$EE^H = N I_N \quad \rightarrow \quad E^{-1} = \frac{1}{N} E^H.$$

Note that  $E$  is symmetric, so we just need to conjugate the entries of  $E$ , i.e.  $E^{-1} = \frac{1}{N} E^*$

- As  $W_N^k = (W_N^*)^*$ , the inverse matrix  $E = E^*$  has the same row 0 as  $E$  and its rows 1 to  $N-1$  are the corresponding rows of  $E$  in reverse order (counting down from  $N-1$  to 1).

## DFT and DTFT frequency comparison

- The DFT can be interpreted as an approx. of the continuous DTFT spectrum by truncation + sampling.
- Truncation of the sequence  $\{x_n\}_{n \geq 0}$  ( $DTFT \rightarrow DFT$ ) is equivalent to a multiplication of the sequence by a rectangular window of length  $N$ . / convolution of the spectrum w/  $W(e^{j\theta}) = \frac{\sin(\frac{\theta}{2}(N+1))}{\sin(\frac{\theta}{2})}$ .
- Therefore the value of  $X_N$  depends on the whole DTFT spectrum by convolution (instead of only over the ideal interval  $[\frac{-\pi}{N}(N-1/2), \frac{\pi}{N}(N+1/2)]$ )  $\rightarrow$  leakage from nearby bin freq. bins.
- As for filtering, the matching of DFT to DTFT spectrum can be improved by using windowing.



# For Personal Use Only -bkwk2

## Fast Fourier Transform

### Fast Fourier Transform (FFT)

- FFT algorithms exploit properties of the DFT Vandermonde matrix  $E$  to perform the DFT in  $O(N \log N)$  operations instead of  $O(N^2)$  normally needed for matrix-vector multiplication.
- ↳ FFT relies on redundancy in the calculation of the basic DFT  $\rightarrow$  reuse past computations.
- ↳ FFT is recursive  $\rightarrow$  logarithmic complexity.
- There are many FFT algorithms that work for length  $N$  w/ specific properties.
  - $\hookrightarrow$  (General) Cooley-Tukey :  $N = LM$ , where  $L, M \in \mathbb{Z}^+$ .
  - $\hookrightarrow$  Radix-2 :  $N = 2^m$ , where  $m \in \mathbb{Z}^+$ .

### Cooley-Tukey algorithm.

- The Cooley-Tukey algorithm works for composite numbers  $N = LM$ , where  $L, M \in \mathbb{Z}^+$ .
- If further factorisation of  $L/M$  is possible, the method can be applied recursively.
- Re-consider the  $N$ -length i/p and o/p vectors  $\underline{x}, \underline{X}$  as a  $L \times M$  array w/  $L$  rows and  $M$  columns, i.e.

$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{N-1} \end{bmatrix} \rightarrow \begin{array}{ccccccccc} x_0 & x_1 & x_2 & \dots & x_M \\ x_M & x_{M+1} & x_{M+2} & \dots & x_{2M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{(L-1)M} & \dots & \dots & \dots & x_{LM} \end{array} \quad L$$

This is row major form — consecutive entries in a row are consecutive in  $\underline{x}$ ; consecutive entries in a column are  $M$  apart in  $\underline{x}$ .

- Reconsidering the time/freq. domain index  $k/n$  using a 2D index, we have

$$k = lM + r$$

$$n = ML + r$$

where  $l, r$  indicate the row no.;  $M, l$  indicate the column no.

- Using the 2D index, we can rewrite the DFT as.

$$X_n = X_{ML+r} = \sum_{k=0}^{N-1} x_k w_N^k = \sum_{s=0}^{M-1} \sum_{l=0}^{L-1} x_{lM+s} w_N^{(lM+s)(ml+r)} = \sum_{s=0}^{M-1} \sum_{l=0}^{L-1} w_N^{lM+s} w_N^{ml+r} = w_N^{ML+r} \sum_{s=0}^{M-1} w_N^{lM+s} w_N^{ml+r}$$

Noting that  $w_N^{lM+r} = e^{-j\frac{2\pi}{N}(lM+r)}$   $= 1$ .  $w_N^L = e^{-j\frac{2\pi}{N}L} = e^{-j\frac{2\pi}{N}M} = W_M$ ,  $w_N^M = e^{-j\frac{2\pi}{N}M} = e^{-j\frac{\pi}{N}} = w_L$ ,

$$X_{ML+r} = \sum_{s=0}^{M-1} W_M^{lM+s} w_N^{ml+r} \quad \boxed{\text{Twiddle factor } L\text{-pt DFT}}$$

$\underbrace{\qquad\qquad\qquad}_{M\text{-pt DFT}}$

i.e. 1) Take the  $L$ -pt DFT for every column

$[M \times L$ -pt DFT  $\rightarrow ML^2$  computations]

2) Multiply every element by its twiddle factor  $w_N^{ml+r}$

$[N = ML$  multiplications  $\rightarrow N$  computations]

3) Take the  $M$ -pt DFT of every row.

$[L \times M$ -pt DFT  $\rightarrow LM^2$  computations]

- The total no. of computations is  $ML + ML^2 + LM^2 = ML(1 + L + M) = N(L + L + M)$

The total no. of computations for the  $N$ -pt DFT is  $N^2$ .

# For Personal Use Only -bkwk2

## Radix-2 algorithm

- The Radix-2 algorithm works for powers of 2,  $N=2^M$ , where  $M \in \mathbb{Z}^+$ .

- We can split the summation of the  $N$ -pt DFT into two parts.

$$X_n = \sum_{k=0}^{N-1} X_k e^{-j\frac{2\pi}{N} nk} = \sum_{k=0}^{\frac{N}{2}-1} X_{2k} e^{-j\frac{2\pi}{N}(2k)n} + \sum_{k=0}^{\frac{N}{2}-1} X_{2k+1} e^{-j\frac{2\pi}{N}(2k+1)n} = \sum_{k=0}^{\frac{N}{2}-1} X_{2k} e^{-j\frac{2\pi}{N/2} nk} + e^{j\frac{\pi}{2}} \sum_{k=0}^{\frac{N}{2}-1} X_{2k+1} e^{-j\frac{2\pi}{N/2} nk}$$

DFT of even index      DFT of odd index

$$\therefore X_n = A_n + W_N^n B_n, \text{ where } A_n, B_n \text{ are } \frac{N}{2}\text{-pt DFTs.}$$

- consider  $X_{n+\frac{N}{2}}$

$$X_{n+\frac{N}{2}} = \sum_{k=0}^{\frac{N}{2}-1} X_{2k} e^{-j\frac{2\pi}{N/2}(n+\frac{N}{2})k} + e^{j\frac{\pi}{2}} \sum_{k=0}^{\frac{N}{2}-1} X_{2k+1} e^{-j\frac{2\pi}{N/2}(n+\frac{N}{2})k} = \sum_{k=0}^{\frac{N}{2}-1} X_{2k} e^{-j\frac{2\pi}{N} nk} - e^{j\frac{\pi}{2}} \sum_{k=0}^{\frac{N}{2}-1} X_{2k+1} e^{-j\frac{2\pi}{N} nk}$$

DFT of even index      DFT of odd index

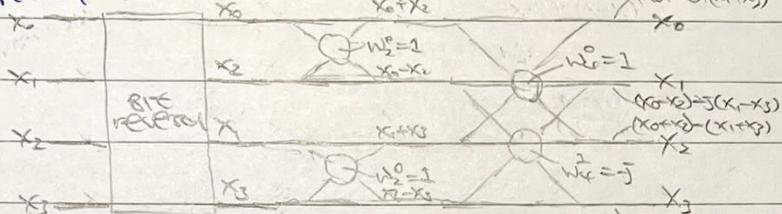
$$\therefore X_{n+\frac{N}{2}} = A_n - W_N^n B_n, \text{ where } A_n, B_n \text{ are } \frac{N}{2}\text{-pt DFTs.}$$

- We have recursion + redundancy by decomposing a  $N$ -pt DFT into  $2 \times \frac{N}{2}$ -pt DFTs + reusing past computations.

$$A_n \quad W_N^n \quad B_n \quad X_n = A_n + W_N^n B_n$$

$$B_n \quad X_{n+\frac{N}{2}} = A_n - W_N^n B_n$$

- e.g.: Consider the 4-pt DFT



**Bit reversal**: convert index from decimal to binary,  $\rightarrow$  reverse bits  $\rightarrow$  convert from binary to decimal.

## Applications of the DFT

### Cyclic convolution

- Consider independent RVs  $X_1, X_2, \dots, X_n$  defined over an alphabet  $\Sigma = \{0, 1, \dots, p-1\}$ , where  $p \in \mathbb{Z}^+$ .

consider the RV  $Z$

$$Z = X_1 \oplus X_2 \oplus \dots \oplus X_n, \text{ where } \oplus \text{ is modular addition } (x \oplus y = R_p(x+y))$$

For all  $n$  and  $p$ , the pmf of  $Z$ ,  $P_Z$ , is given by the cyclic convolution of the pmfs of  $X_i$ ,  $P_{X_i}$ .

$$P_Z = P_{X_1} \oplus P_{X_2} \oplus \dots \oplus P_{X_n}$$

- For  $p=2$ ,  $X_1, X_2, \dots, X_n$  are binary RVs w/ pmf  $P_{X_i} = [1-p_i, p_i]$

$$1) 2\text{-pt DFT of } P_{X_i}: Q_{X_i} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1-p_i \\ p_i \end{bmatrix} = \begin{bmatrix} 1 \\ 1-2p_i \end{bmatrix}$$

$$2) \text{Component wise multiplication: } Q_Z = \begin{bmatrix} \frac{1}{2} & 1 \\ \frac{1}{2} & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 1-2p_1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & 1 \\ \frac{1}{2} & 1-2p_1 \end{bmatrix}$$

$$3) \text{Inverse 2-pt DFT: } P_Z = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 1 \\ \frac{1}{2} & 1-2p_1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} + \frac{1}{2} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} - \frac{1}{2} \frac{1}{2} & \frac{1}{2} - 1 + \frac{1}{2} \frac{1}{2} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} - 1 + \frac{1}{2} \frac{1}{2} \end{bmatrix}$$

- In general, we evaluate the circular convolution using FFT  $\rightarrow$  componentwise multiplication  $\rightarrow$  iFFT.

OFDM (Fitting a convolution channel in doing cyclic convolution).

(Refer to FFT notes)

# For Personal Use Only -bkwk2

## Linear convolution

- Linear convolutions of arbitrary sequences cannot in general be computed via an FFT.

However, consider the case of an FIR filter of length  $L+1$ .

infinite finite semi-infinite

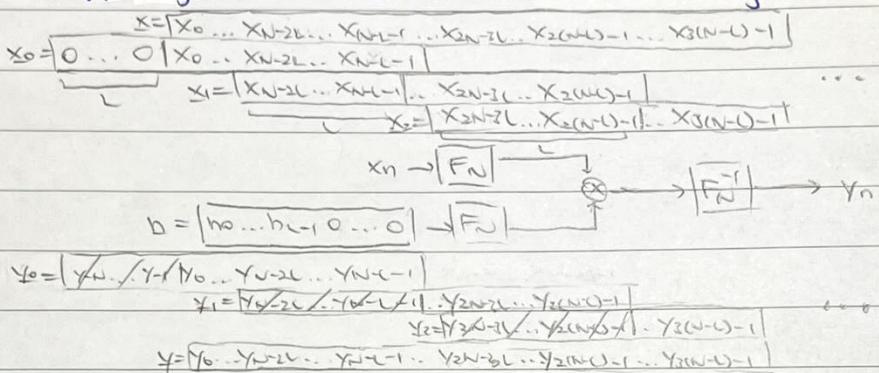
$$\{Y_k\}_{k \geq 0} = \{h_k\}_{k \geq 0} * \{x_k\}_{k \geq 0}, \quad \text{i.e. } Y_k = \sum_{m=0}^L h_m x_{k-m}$$

- If we take a block  $x = [x_0 \dots x_{L-1}]$  and evaluate the circular convolution  $y = h \otimes x$ ,

the first  $L$  ops differ from linear convolution because they depend on the last  $L$  entries of the impulse, instead of  $x_L \dots x_{-1} = 0$  for linear convolution.

- For  $N \gg L$ , the ops beyond the first  $L$  will be identical to linear convolution.

- We can use length  $N$  FFT to do linear convolutions w/ FIR of length  $L+1 \ll N$  w/ overoptimal methods.



## Interpolation and decimation.

- Signal processing systems often req. a change of sampling freq.  $f_s$  so processing can be done at a higher  $f_s$  (better accuracy) or lower  $f_s$  (reduce computations)

- Interpolation (upping to a higher freq.) is done by first zero stuffing: insert  $k$  zeros b/w every two samples, then filtering: low-pass and scale the resulting signal.

- Decimation (reducing to a lower freq.) is done by dropping samples: keep only one in every  $k$  samples. (To avoid aliasing, this can only be done when  $R_s < \frac{f_s}{2}$ )

- Interpolating / decimating via the FFT gives a graceful transition b/w freq., and no longer req.  $f_s^{new}$  to be a multiple or divisor of  $f_s^{old}$ .

- To interpolate/decimate, we pick  $N$  and  $N'$  s.t.  $\frac{N}{N'} = \frac{f_s^{new}}{f_s^{old}}$ .

1) Take the  $N$ -pt FFT of the i/p  $x$ .

2) Interpolation: Insert zeros b/w  $X_{N/2+1}$  and  $X_{N/2+1}$  (odd  $N$ ) / either side of  $X_{N/2}$  (even  $N$ )

Decimation: Delete components symmetrically around  $X_{N/2}$ .

3) Take the  $N'$ -pt IFFT of the modified spectrum

Interpolation

$$\begin{array}{c} [x_0 \dots x_{N/2-1} | x_{N/2+1} \dots x_{N-1}] \\ \downarrow \\ [x_0 \dots x_{N/2-1} | 0 \ 0 \ 0 \ 0 \ 0 \ 0 | x_{N/2+1} \dots x_{N-1}] \end{array}$$

Decimation

$$\begin{array}{c} [x_0 \dots x_{N/2-1} | x_{N/2} \dots x_{N/2+k-1} | x_{N/2+k+1} \dots x_{N-1}] \\ \downarrow \\ [x_0 \dots x_{N/2-1} | x_{N/2+k+1} \dots x_{N-1}] \end{array}$$

\* The FFT method has the usual cost of a cut in the spectrum being equivalent to multiplying by a rectangular window  $\rightarrow$  can be improved using window-based approaches.

## Continuous time random signals

### Random process

- A random signal is a time-dependent RV,  $X(t)$ , which is an instance drawn randomly from a set of possible signal waveforms. (ensemble)
- The ensemble of random signals + their associated probability distributions is known as a random process.
- We denote the random process as  $\{X(t|\alpha)\}$  and an instance of the random process, the random signal as  $X(t, \alpha)$ , where  $\alpha$  is a draw from some set  $A$ .

### Correlation and covariance functions

- The autocorrelation function of a random process  $\{X(t|\alpha)\}$  is

$$r_{XX}[t_1, t_2] = E_A[X(t_1, \alpha) X(t_2, \alpha)] = \int_A X(t_1, \alpha) X(t_2, \alpha) f_A(\alpha) d\alpha$$

The autocovariance function of a random process  $\{X(t|\alpha)\}$  is

$$c_{XX}[t_1, t_2] = E_A[(X(t_1, \alpha) - M_{t_1})(X(t_2, \alpha) - M_{t_2})] = \int_A (X(t_1, \alpha) - M_{t_1})(X(t_2, \alpha) - M_{t_2}) f_A(\alpha) d\alpha$$

- The crosscorrelation function between random processes  $\{X(t|\alpha)\}$  and  $\{Y(t|\alpha)\}$  is

$$r_{XY}[t_1, t_2] = E_A[X(t_1, \alpha) Y(t_2, \alpha)] = \int_A X(t_1, \alpha) Y(t_2, \alpha) f_A(\alpha) d\alpha$$

The crosscovariance function between random processes  $\{X(t|\alpha)\}$  and  $\{Y(t|\alpha)\}$  is

$$c_{XY}[t_1, t_2] = E_A[(X(t_1, \alpha) - M_{X,t_1})(Y(t_2, \alpha) - M_{Y,t_2})] = \int_A (X(t_1, \alpha) - M_{X,t_1})(Y(t_2, \alpha) - M_{Y,t_2}) f_A(\alpha) d\alpha$$

\* Note the expectations  $E_A[\cdot]$  are over the whole ensemble.

### Stationarity

- stationarity means the statistical characteristics of a signal do not change over time.
  - A random process is strict sense stationary (SSS) if its probability distribution does not change over time.
- Formally, for all finite  $N$  and all sets of time pts  $\{t_1, \dots, t_N\}$ , the  $N$ th order joint pdf is invariant for any time shift  $T$ .

$$f_{X(t_1), \dots, X(t_N)}(x_1, \dots, x_N) = f_{X(t_1+T), \dots, X(t_N+T)}(x_1, \dots, x_N)$$

- A random process is wide sense stationary (WSS) iff

$$(i) E[X(t)] = M \quad \forall t \quad (\text{const. mean})$$

$$(ii) \text{Var}[X(t)] < \infty / E[X(t)^2] < \infty \quad \forall t \quad (\text{finite variance/power}) \quad [\text{ignored in 3F1}]$$

$$(iii) r_{XX}[t_1, t_2] \rightarrow r_{XX}[t_2 - t_1] \quad (r_{XX} = r_{XX}(t_2 - t_1))$$

- For WSS  $\{X(t|\alpha)\}$ ,  $r_{XX}[t] = r_{XX}[-t]$ , i.e.  $r_{XX}[t]$  is even

$$r_{XX}[t] = E[X(t) X(t + T)] = E[X(t + T) X(t)] = r_{XX}[-t]$$

Moreover, as  $T$  becomes large,  $X(t)$  and  $X(t+T)$  tends to become decorrelated. ( $r_{XX}[T] \rightarrow E[X(t)] E[X(t+T)] = M^2$ )

$$\text{Therefore, } |r_{XX}[T]| = r_{XX}[0] \quad (\text{formal proof in 3F3 notes})$$

The width of the autocorrelation function  $r_{XX}[t]$  tells us how slowly the signal fluctuates in time (width of  $\frac{1}{r_{XX}[t]}$ )

# For Personal Use Only -bkwk2

## Power spectral density (PSD)

- For WSS  $\{x(t)\}$ , the PSD function is the FT of the autocorrelation function  $R_{xx}[t]$

$$S_x(\omega) = \int_{-\infty}^{\infty} R_{xx}[t] e^{-j\omega t} dt$$

$$R_{xx}[t] = \frac{1}{2\pi} \int_{-\infty}^{\infty} S_x(\omega) e^{j\omega t} d\omega$$

- For jointly WSS  $\{x(t)\}, \{y(t)\}$ , the cross PSD function is the FT of the crosscorrelation function  $R_{xy}[t]$

$$S_{xy}(\omega) = \int_{-\infty}^{\infty} R_{xy}[t] e^{-j\omega t} dt$$

$$R_{xy}[t] = \frac{1}{2\pi} \int_{-\infty}^{\infty} S_{xy}(\omega) e^{j\omega t} d\omega$$

## Ergodicity

- If we can exchange ensemble averages  $E[\cdot]$  for time averages  $\langle \cdot \rangle_T$ , the random process is ergodic.

- This means the ensemble of the random process simply composed of all possible time shifts of a single random signal.

- Ensemble average  $E_A[\cdot]$  and time averages  $\langle \cdot \rangle_T$  are defined as

$$E_A[X(t+\alpha)] = \int_A x(t+\alpha) f_A(\alpha) d\alpha \quad \langle X(t+\alpha) \rangle_T = \lim_{T \rightarrow \infty} \frac{1}{2T} \int_{-T}^T X(t+\alpha) dt$$

- A WSS  $\{x(t)\}$  could be mean ergodic/correlation ergodic, defined as

↳ mean ergodic

$$M = E_A[X(t+\alpha)] = \langle X(t+\alpha) \rangle_T$$

↳ correlation ergodic

$$R_{xx}[t] = E_A[X(t+\alpha)X(t+\tau+\alpha)] = \langle X(t+\alpha)X(t+\tau+\alpha) \rangle_T$$

white noise

- In reality, no process is truly stationary (therefore cannot be ergodic), but many noise processes are approx stationary → use ergodicity to simplify measurement of WSS process.

- We may estimate quantities from a finite period  $[T_1, T_2]$  of a ergodic process using

$$M \approx \frac{1}{T_2 - T_1} \int_{T_1}^{T_2} X(t+\alpha) dt$$

$$R_{xx}[t] \approx \frac{1}{T_2 - T_1} \int_{T_1}^{T_2} X(t+\alpha) X(t+\tau+\alpha) dt$$

## System response to random signals

### Linear systems and random processes

- If the ip  $\{x(t)\}$  of a LTI system w/ impulse resp  $h(t)$  is WSS, then the op  $\{y(t)\}$  is also WSS.

(i) Consider the mean  $E[y(t)]$

$$E[Y(t)] = E\left[\int h(k) x(t+k) dk\right] = \int h(k) E[x(t+k)] dk = M x \int h(k) dk \rightarrow \text{const. over all } t.$$

(ii) consider the autocorrelation  $R_{yy}[t]$

$$R_{yy}[t] = E\left[\int h(k) x(t+k) dk, \int h(k_2) x(t+k_2) dk_2\right] = E\left[\int \int h(k) h(k_2) x(t+k) x(t+k_2) dk dk_2\right]$$

$$= \int \int h(k) h(k_2) E[x(t+k) x(t+k_2)] dk dk_2 = \int \int h(k) h(k_2) R_{xx}[t+k+k_2] dk dk_2 = \int h(k) [h(k) R_{xx}[t+k+k_2]] dk dk_2,$$

$$= \int h(k) [h(t+k) + R_{xx}[t+k]] dk, = h(-t) * h(t) + R_{xx}[t] \rightarrow \text{only a function of } t.$$

- The ip  $\{x(t)\}$  and op  $\{y(t)\}$  of the LTI system is jointly WSS. We already know  $\{x(t)\}, \{y(t)\}$  are WSS.

$$R_{xy}[t] = E[x(t)y(t+\tau)] = E[x(t) \int h(k) y(t+k+\tau) dk] = \int h(k) E[x(t)x(t+k+\tau)] dk = \int h(k) R_{xx}[t+k+\tau] dk = h(\tau) + R_{xx}[\tau]$$

- Notice that the correlation functions can be expressed as convolutions

$$R_{xy}[t] = h(t) * h(-t) * R_{xx}[t]$$

$$R_{yy}[t] = h(t) * R_{xx}[t]$$

# For Personal Use Only -bkwk2

- Taking the FT of the correlation function expressed as convolutions,

$$S_Y(\omega) = H(\omega) H^*(\omega) S_X(\omega) = |H(\omega)|^2 S_X(\omega) \quad S_{XY}(\omega) = H(\omega) S_X(\omega),$$

Rearranging for  $H(\omega)$ ,

$$\boxed{H(\omega) = \frac{S_Y(\omega)}{S_X(\omega)}} \quad \text{only magnitude}$$

$$\boxed{H(\omega) = \frac{S_{XY}(\omega)}{S_X(\omega)}} \quad \text{both magnitude and phase}$$

- If an important system is subject to measurable random fluctuations at the input, we can measure the system freq. response w/o taking the plant offline. (estimate  $r_{XX}(t)$ ,  $r_{YY}(t)$ ,  $r_{XY}(t)$  w/ time averages  $\rightarrow S_X(\omega)$ ,  $S_Y(\omega)$ ,  $S_{XY}(\omega) \rightarrow H(\omega)$ )

Physical interpretation of PSD.

- For a deterministic signal  $x(t)$ , the instantaneous power is  $x(t)^2$  and the average power is  $\lim_{T \rightarrow \infty} \frac{1}{2T} \int_0^T x(t)^2 dt$
- For a random process  $\{x(t)\}$ , the expected instantaneous power is  $E[x(t)^2]$  and the expected average power is  $E\left[\lim_{T \rightarrow \infty} \frac{1}{2T} \int_0^T x(t)^2 dt\right] = \lim_{T \rightarrow \infty} \frac{1}{2T} \int_0^T E[x(t)^2] dt$
- If  $\{x(t)\}$  is WSS, then  $E[x(t)^2] = r_{xx}[0]$  for all  $t$ , so the expected average power is  $r_{xx}[0]$ .

$$r_{xx}[0] = \frac{1}{2\pi} \int_{-\infty}^{\infty} S_X(\omega) e^{j\omega t} d\omega = \frac{1}{2\pi} \int_{-\infty}^{\infty} S_X(\omega) d\omega$$

i.e. the area under  $S_X(\omega)$  gives the total power of  $\{x(t)\}$ .

- To find the power of a WSS process  $\{x(t)\}$ , in a particular freq. band, say  $[w_1, w_2]$  s.t.  $0 < w_1 < w_2$ , consider passing  $\{x(t)\}$  through an ideal bandpass filter w/ TF  $H(\omega) = \begin{cases} 1 & \text{if } w_1 \leq \omega \leq w_2 \\ 0 & \text{otherwise} \end{cases}$ , then find the avg power  $r_{yy}[0]$ .

$$r_{yy}[0] = \frac{1}{2\pi} \int_{-\infty}^{\infty} S_Y(\omega) d\omega = \frac{1}{2\pi} \int_{-\infty}^{\infty} (H(\omega))^2 S_X(\omega) d\omega = \frac{1}{2\pi} \left[ \int_{w_1}^{w_2} S_X(\omega) d\omega + \int_{w_2}^{\infty} S_X(\omega) d\omega \right] = \frac{1}{\pi} \int_{w_1}^{w_2} S_X(\omega) d\omega$$

$\rightarrow$  The power of  $\{x(t)\}$  in a particular freq. band  $w_1 < w < w_2$  is  $\boxed{\frac{1}{\pi} \int_{w_1}^{w_2} S_X(\omega) d\omega}$ .

White noise process.

- A WSS process  $\{x(t)\}$  is termed white noise if

$$\boxed{c_{xx}[\tau] = \sigma^2 \delta(\tau)}$$

where  $\delta(\tau)$  is the delta function.

- The PSD  $S_X(\omega)$  for white noise is given by

$$S_X(\omega) = \int_{-\infty}^{\infty} c_{xx}[\tau] e^{j\omega\tau} d\tau = \int_{-\infty}^{\infty} (\sigma^2 \delta(\tau) + \mu^2) e^{j\omega\tau} d\tau = \int_{-\infty}^{\infty} \sigma^2 \delta(\tau) e^{-j\omega\tau} d\tau + \mu^2 \int_{-\infty}^{\infty} e^{j\omega\tau} d\tau = \sigma^2 + \mu^2 2\pi \delta(\omega)$$

For zero-mean white noise,  $\mu = 0$ , so  $\boxed{S_X(\omega) = \sigma^2}$ , i.e. flat across all freq.

- Note that power for white noise is  $E[x(t)^2] = r_{xx}[0] = \sigma^2 + \mu^2 2\pi \delta(0) \rightarrow \infty$ , so pure white noise is unrealizable in practice. (However, we can have "nearly white" processes).

# For Personal Use Only -bkwk2

From discrete time to continuous time.

- consider a discrete time i/p  $x_k$  that is zero for  $k \notin [-\frac{N}{2}+1, \frac{N}{2}]$  and each  $x_k$  is sampled independently from a distribution w/  $M=0, \sigma^2=1$ .

$$|x(e^{j\theta})|^2 = x(e^{j\theta}) x^*(e^{j\theta})$$

$$|x(e^{j\theta})|^2 = \left( \sum_k x_k e^{-jk\theta} \right) \left( \sum_l x_l e^{jl\theta} \right) = \sum_k w_k^2 + \sum_{k \neq l} w_k w_l e^{j\theta(l-k)}$$

Taking expectations, noting that  $E[x_k^2] = 1$  and  $E[x_k x_l] = E[x_k]E[x_l] = 0$  for  $k \neq l$ ,

$$E[|x(e^{j\theta})|^2] = N$$

$$\therefore \frac{1}{N} E[|Y(e^{j\theta})|^2] = \frac{1}{N} E[|H(e^{j\theta})|^2 |x(e^{j\theta})|^2] = \frac{1}{N} |H(e^{j\theta})|^2 E[|x(e^{j\theta})|^2] = |H(e^{j\theta})|^2$$

- We can derive the continuous case by taking the limit as sampling time tends to zero.

- consider a continuous time i/p  $x_n(t)$  b/w  $-T/2$  and  $T/2$ , given by

$$x_n(t) = \sum_{k=-\frac{N}{2}}^{\frac{N}{2}-1} p(t-kh) w_k$$

where  $h = T/N$  is the sampling time, and each  $w_k$  is sampled independently from a distribution w/  $M=0, \sigma^2=\frac{1}{h}$ .

$p(t) = H(t) - H(t-h)$  is the pulse function.

- The FT of  $x_n(t)$ ,  $X_n(j\omega)$  is given by

$$X_n(j\omega) = FT \left[ \sum_{k=-\frac{N}{2}}^{\frac{N}{2}-1} p(t-kh) w_k \right] = \sum_{k=-\frac{N}{2}}^{\frac{N}{2}-1} FT[p(t-kh)] w_k = \sum_{k=-\frac{N}{2}}^{\frac{N}{2}-1} \frac{w_k}{h} \left( \frac{1}{j\omega - jkh} \right) e^{-j\omega kh}$$

$FT[H(t)]$   $FT[H(t-h)]$  time shift by  $kh$

Consider  $|X_n(j\omega)|^2 = X_n(j\omega) X_n^*(j\omega)$

$$|X_n(j\omega)|^2 = \left( \sum_k w_k \left( \frac{1-e^{j\omega kh}}{j\omega} \right) e^{-j\omega kh} \right) \left( \sum_l w_l \left( \frac{1-e^{j\omega lh}}{j\omega} \right) e^{j\omega lh} \right) = \sum_k w_k^2 \left( \frac{1-e^{j\omega kh}}{j\omega} \right) \left( \frac{1-e^{-j\omega kh}}{j\omega} \right) + \sum_{k \neq l} w_k w_l \left( \frac{1-e^{j\omega kh}}{j\omega} \right) \left( \frac{1-e^{-j\omega lh}}{j\omega} \right) j\omega h$$

Taking expectations, noting that  $E[w_k^2] = \frac{1}{h}$  and  $E[w_k w_l] = E[w_k]E[w_l] = 0$  for  $k \neq l$ ,  $h = T/N \rightarrow T = Nh$

$$E[|X_n(j\omega)|^2] = \sum_k E[w_k^2] \left( \frac{(1-e^{j\omega kh})(1-e^{-j\omega kh})}{h^2} \right) = \sum_k \left( \frac{1}{h} \right) \left( \frac{(1-(1-j\omega h))(1-(1+j\omega h))}{h^2} \right) = \sum_k h = T$$

$\therefore$  Define  $x(t) = \lim_{N \rightarrow \infty} x_n(t)$ , i.e. white noise, then we have

$$\frac{1}{T} E[|Y(j\omega)|^2] = \frac{1}{T} E[|H(j\omega)|^2 |x(j\omega)|^2] = \frac{1}{T} |H(j\omega)|^2 E[|x(j\omega)|^2] = |H(j\omega)|^2$$