

SIGNAL \Rightarrow anything that carries some information

eg: ECG, EEG, speech signal, seismic signal

signal is defined as any physical quantity that varies with time, space or any other independent variable.

1D signal \rightarrow depends on only one variable eg: speech, AC power supply

2D signal \rightarrow depends on two independent variables eg: image

multi-D signal \rightarrow more independent variables

eg: speed of wind (lat, long, elevation & time)

CLASSIFICATION OF SIGNALS

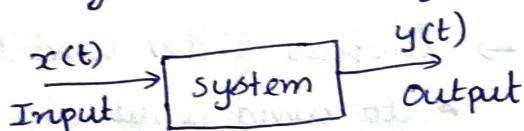
continuous-time signals: defined for every instant of time $x(t)$

discrete-time signals: defined at discrete instants of time $x(n)$

They are continuous in amplitude and discrete in time

Digital signals: discrete in time and quantized in amplitude

SYSTEM \rightarrow physical device that generates a response or output signal for a given input signal



Mathematically $y(t) = T[x(t)]$

System \rightarrow continuous time system \rightarrow operates on cont. time signal

* produces cont. time output i.e $y(t) = T[x(t)]$

System \rightarrow discrete-time system \rightarrow operates on discrete time signal

* produces discrete time output i.e $y(n) = T[x(n)]$

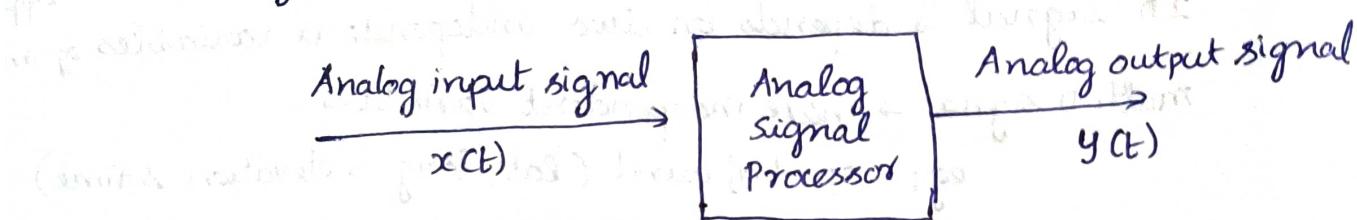
SIGNAL PROCESSING \rightarrow any operation that changes the characteristics (amplitude, shape, phase & frequency) of the signal

Analog signal processing

Real world is analog signal \rightarrow function of continuous variable
(time or space)

for processing analog signal \rightarrow amplifier, filter, freq. analyzer
respond to continuous variation of input signal inst. amplifies

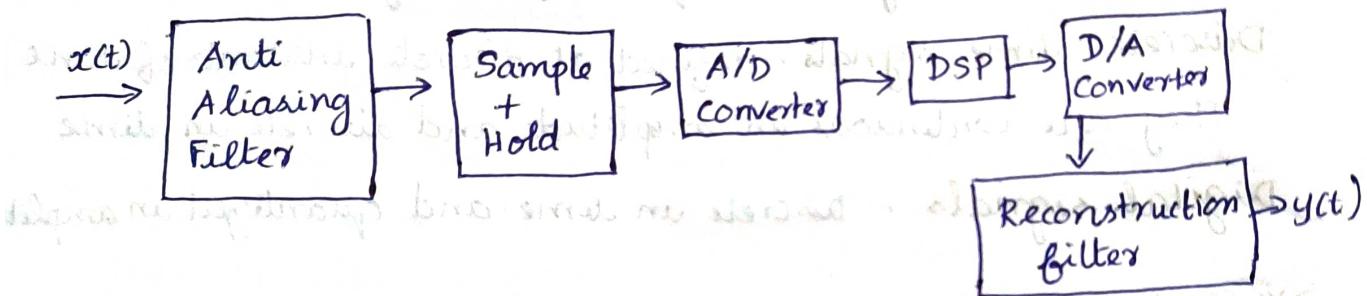
changes the characteristics of the signal or extract desired info



Digital signal processing

Convert to digital form

to digital form for processing benefits more accurate



Input signal $x(t)$ \rightarrow from transducer, communication signal
OR ECG, EEG

Anti-aliasing filter \rightarrow lowpass filter used to remove high freq.
 \star to band limit the signal

notch filter & amplifier may be used

\downarrow to remove power component \downarrow to bring signal upto voltage range
required for A/D conversion

Sample and hold \rightarrow provides input to ADC

input must remain relatively constant during A/D conversion

ADC \rightarrow output of sample & hold is input to ADC

output from ADC \rightarrow N-bit binary number depending on
the value of analog input signal

range of input to ADC 0 to $10V$ (unipolar)

-5 to $5V$ (bipolar)

Provided by preceding amplifier

Digital signal processor \rightarrow DSP may be large programmable digital computer or microprocessor programmed to perform desired operation on input signal.

Eg: ADSP2100, Motorola DSP 56000, TMS320C50

Configured to perform specific set of operations

D/A converter: Input \rightarrow digital signal from processor

Output \rightarrow continuous but not smooth signal

containing unwanted high frequency components

Reconstruction filter: To eliminate high frequency components

Output is a smooth continuous signal

Advantages of DSP

1. Greater accuracy: Tolerance of analog circuit components affects accuracy. DSP provides superior control of accuracy.

2. Cheaper: Digital realization is comparatively cheaper than analog.

3. Ease of data storage: Digital signal easily stored on magnetic media without loss of fidelity & can be processed remotely in offline.

4. Implementation of sophisticated algorithms

DSP allows to implement sophisticated algorithms than analog.

5. Flexibility in configuration:

DSP system \rightarrow easily reconfigured by changing program.

Analog system \rightarrow involve redesign of system hardware.

6. Applicability of VLF signals \rightarrow low frequency signals eg: seismic signal can be easily processed using DSP. ASP require very large LC components.

7. Time sharing: DSP allow sharing of processor among a no. of signals by time sharing thus reducing cost of processing.

LIMITATIONS

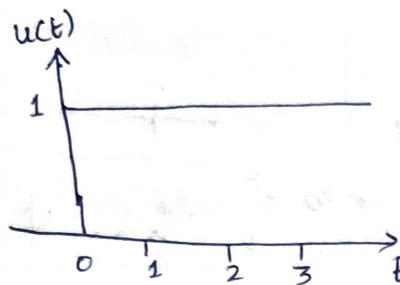
1. System complexity : System includes devices such as ADC, DAC & associated filters that increases complexity of DSP systems.
2. Bandwidth limited by sampling rate:
Band limited signal \rightarrow sampled without information if sampling rate $> 2B$
Wide bandwidth \rightarrow require fast sampling rate ADC & DSP.
Practical limitation in the speed of ADC & DSP
3. Power consumption:
ASP \rightarrow implemented using passive circuit elements (L, C, R) need less power
DSP \rightarrow contains over 4 lakh transistors dissipate more power (1 watt)

APPLICATIONS OF DSP

1. Telecommunication - Modem, video conferencing, cellular phone, FAX, line repeaters, channel Mux
2. Consumer electronics - Digital TV/Audio, electronic music, FM stereo
3. Instrumentation & control - spectrum analysis, PLL, Digital filter, servo control, Robot control, process control
4. Image processing : Image compression, enhancement, Image analysis & recognition
5. Medicine - CT, X-ray scanning, MRI, spectrum analysis of ECG, EEG
6. Speech processing - automatic speech recognition, text to speech
7. Seismology \rightarrow oil & gas exploration, earthquake monitoring
8. Military : Radar SP, sonar SP, navigation, secure communications

CONTINUOUS - TIME SIGNALS

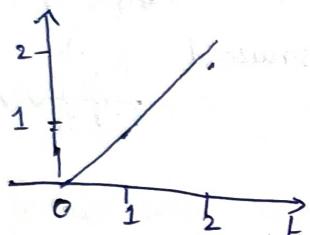
1. Unit step function



$$u(t) = 1 \text{ for } t \geq 0$$

$$0 \text{ for } t < 0$$

2. Unit ramp



$$r(t) = t \text{ for } t \geq 0$$

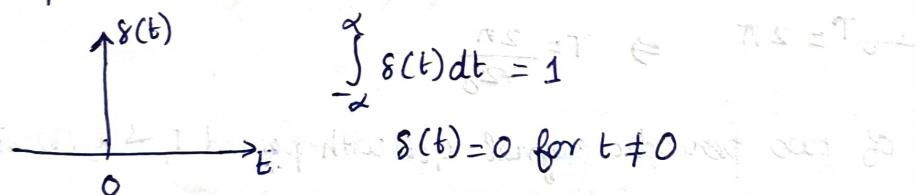
$$0 \text{ for } t < 0$$

$$(or) r(t) = t u(t)$$

$$u(t) \rightarrow \boxed{\int dt} \rightarrow \boxed{r(t)} \Rightarrow \int u(t) dt = \int d\tau = t = r(t)$$

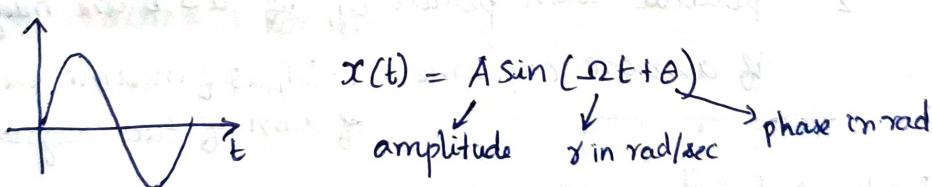
$$\frac{d}{dt} \rightarrow u(t) \Rightarrow \frac{d(r(t))}{dt} = \frac{dt}{dt} = u(t)$$

3. Impulse function



$$\int_{-\infty}^{\infty} \delta(t) dt = 1$$

4. Sinusoidal signal



$$x(t) = A \sin(\Omega t + \theta)$$

amplitude A in rad/sec phase in rad

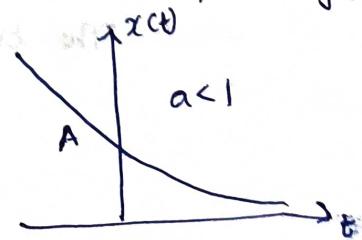
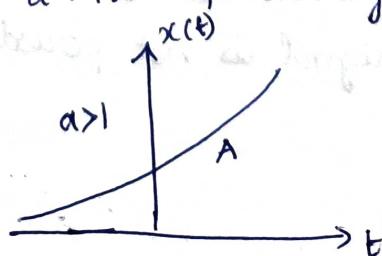
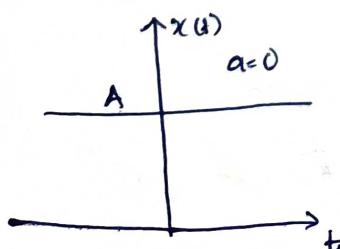
Periodic $\rightarrow x(t+T) = x(t)$ $T \rightarrow$ fundamental period

5. Real exponential signals

$$x(t) = A e^{at}$$

$a \rightarrow +ve$ exponentially \uparrow

$a \rightarrow -ve$ exponentially \downarrow



6. complex exponential signal
- $$x(t) = e^{st} \quad s = \sigma + j\omega \quad x(t) = e^{(\sigma+j\omega)t} = e^{\sigma t} e^{j\omega t}$$
- $$x(t) = e^{\sigma t} [\cos \omega t + j \sin \omega t]$$
- i) $\sigma = 0 \wedge \omega = 0 \quad x(t) = 1 \text{ DC}$
 - ii) $\sigma = 0 \wedge \omega \neq 0 \quad x(t) = e^{\sigma t} \quad \begin{matrix} \sigma > 0 \\ \sigma < 0 \end{matrix}$
 - iii) $\sigma \neq 0 \wedge \omega = 0 \quad x(t) \rightarrow \text{exponential}$
 - iv) $s = \sigma + j\omega \quad \begin{matrix} \sigma > 0 \text{ exponentially } \uparrow \text{ sinusoid} \\ \sigma < 0 \text{ exponentially } \downarrow \text{ sinusoid} \end{matrix}$
-

continuous time periodic signals

$$x(t+T) = x(t) \quad -\infty < t < \infty$$

eg: complex exponential \Rightarrow sinusoidal signals

$$x(t) = A \sin(\omega_0 t + \theta)$$

$$x(t) = A \sin[\omega_0(t+T) + \theta] = A \sin[\omega_0 t + \omega_0 T + \theta]$$

$$\omega_0 T = 2\pi \Rightarrow T = \frac{2\pi}{\omega_0}$$

sum of two periodic signal $x_1(t)$ with period T_1 , $x_2(t)$ with T_2

$$\frac{T_1}{T_2} = \frac{a}{b} \quad a \& b \text{ are integers}$$

$bT_1 = aT_2$ periodic with period bT_1 , if a & b are integer

if a & b are coprime ($\Rightarrow T = bT_1$ is fundamental period of sum of two signals)

$$\text{eg1: } T_1 = 5 \wedge T_2 = 7$$

$$\frac{T_1}{T_2} = \frac{5}{7} \text{ (integer)} \Rightarrow T_1 = 5T_2$$

$$\Rightarrow T = 7T_1 = 35 \text{ periodic signal}$$

$$\text{eg 2: } T_1 = 7 \wedge T_2 = \pi$$

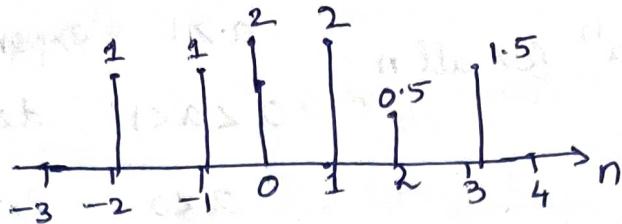
$$\frac{T_1}{T_2} = \frac{7}{\pi} \text{ (not integer)}$$

The two signal is not periodic

DISCRETE TIME SIGNALS

Representation:

1. Graphical



2. Functional

$$x(n) = \begin{cases} 0.5 & \text{for } n=2 \\ 1.5 & \text{for } n=3 \\ 0 & \text{otherwise} \end{cases}$$

3. Tabular

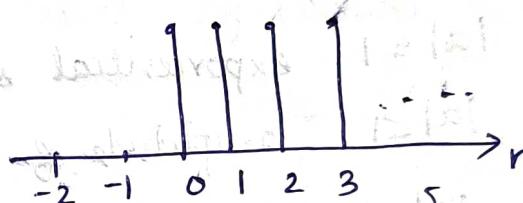
n	-1	0	1	2	3
$x(n)$	1	2	2	0.5	1.5

4. Sequence $x(n) = \{2, 4, 6, 8, -3\}$

Elementary Signals

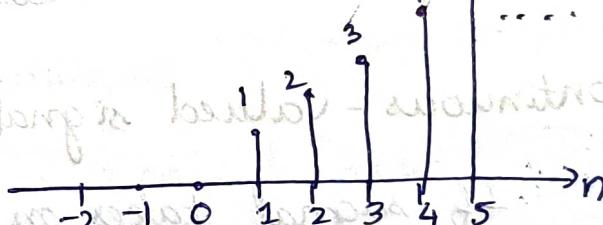
1. Unit step sequence

$$u(n) = \begin{cases} 1 & \text{for } n \geq 0 \\ 0 & \text{for } n < 0 \end{cases}$$



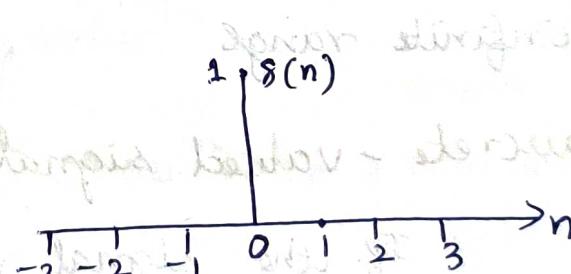
2. Unit ramp sequence

$$r(n) = \begin{cases} n & \text{for } n \geq 0 \\ 0 & \text{for } n < 0 \end{cases}$$



3. Unit impulse sequence

$$\delta(n) = \begin{cases} 1 & \text{for } n=0 \\ 0 & \text{for } n \neq 0 \end{cases}$$



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

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

$$\sum_{n=-\infty}^{\infty} x(n) \delta(n-n_0) = x(n_0)$$

4. Exponential sequence

$$x(n) = a^n \text{ for all } n$$

$a > 1$ exponentially grows

$0 < a < 1$ decays exponentially

$a < 0$ exponential signal take alternating signs

5. Sinusoidal signal

$$x(n) = A \cos(\omega_0 n + \phi)$$

$$\text{Euler's identity } x(n) = \frac{A}{2} e^{j\omega_0 n} e^{j\phi} + \frac{A}{2} e^{-j\omega_0 n} e^{-j\phi}$$

$$|e^{j\omega_0 n}|^2 = 1 \quad \text{energy} = \alpha \quad \text{power} = 1$$

6. Complex exponential signal

$$x(n) = a^n e^{j(\omega_0 n + \phi)} = a^n \cos(\omega_0 n + \phi) + j a^n \sin(\omega_0 n + \phi)$$

$|a| = 1$ exponential seq. is sinusoidal

$|a| < 1$ amplitude of sinusoidal sequence decays exponentially

$|a| > 1$ " increases exponentially"

Continuous-valued signal

If signal takes on all possible values on a finite or infinite range

Discrete-valued signal

If the signal takes on values from a finite set of possible values. Usually values are equidistant

Digital - Discrete time signal having a set of discrete values



quantization \rightarrow rounding or truncation

CLASSIFICATION OF DISCRETE TIME SIGNALS

Energy & power signals

$$E = \sum_{n=-\infty}^{\infty} |x(n)|^2 \quad \text{average power}$$

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

if total E is finite $\Rightarrow P=0$
average power is finite $\Rightarrow E=\infty$

Periodic & aperiodic signal

$x(n)$ is said to be periodic with period N if

$$x(N+n) = x(n) \text{ for all } n$$

smallest value of $N \rightarrow$ fundamental period

symmetric (even) and antisymmetric (odd) signal

Even signal $x(-n) = x(n)$ for all n eg: $A \cos \omega n$

odd signal $x(-n) = -x(n)$ eg: $A \sin \omega n$

$$x(n) = x_e(n) + x_o(n) \quad x(n) \rightarrow \text{represented by sum of even & odd components}$$

$$x(-n) = x_e(-n) + x_o(-n) = x_e(n) - x_o(n)$$

$$x(n) + x(-n) = 2x_e(n) \quad x_e(n) = \frac{1}{2} [x(n) + x(-n)]$$

$$x_o(n) = \frac{1}{2} [x(n) - x(-n)]$$

Causal and non-causal signal

$x(n)$ is said to be causal if its value is zero for $n < 0$

$$\text{eg: } x_1(n) = a^n u(n)$$

$$x_2(n) = \{1, 2, -3, -1, 2\}$$

$$\text{non causal} \quad \text{eg: } x_1(n) = a^n u(-n+1)$$

$$x_2(n) = \{1, -2, 1, 4, 3\}$$

Operation on Signals

Signal processing is a group of basic operations applied to an input signal resulting in another signal as the output

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

Basic operations are shifting, time reversal, time scaling, scalar multiplication, signal multiplier, signal addition

Shifting: takes input sequence & shifts the value by an integer increment of independent variable

$$y(n) = x(n-k)$$

k - tve delays the sequence k - +ve advances the seq

Time reversal \rightarrow by folding the sequence about $n=0$

Time scaling \rightarrow by replacing n by na

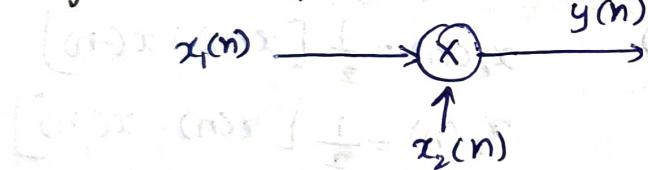
- $\xrightarrow{n \rightarrow \text{integer} \rightarrow \text{compress}}$
- $\xrightarrow{\text{scale multiple}}$
- $\xrightarrow{\frac{n}{a} \rightarrow \text{expands}}$

$x(2n) \rightarrow$ reducing sampling rate by a factor of 2
 ↳ down sampling or decimation

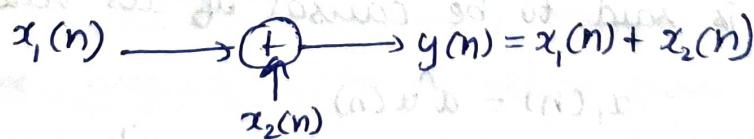
$x(n/2) \rightarrow$ increases sampling rate by 2 \rightarrow up sampling

Scalar multiplication $x(n)$ is multiplied by a scale factor a

Signal multiplier \rightarrow multiply two signal sequences

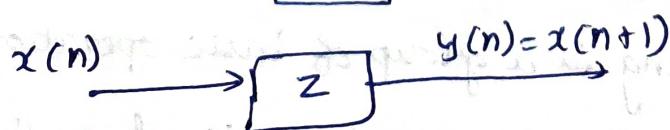


Addition operation: Two signal sequences can be added using adder



Unit delay $y(n) = x(n-1)$ (delayed by one sample)

↳ memory element

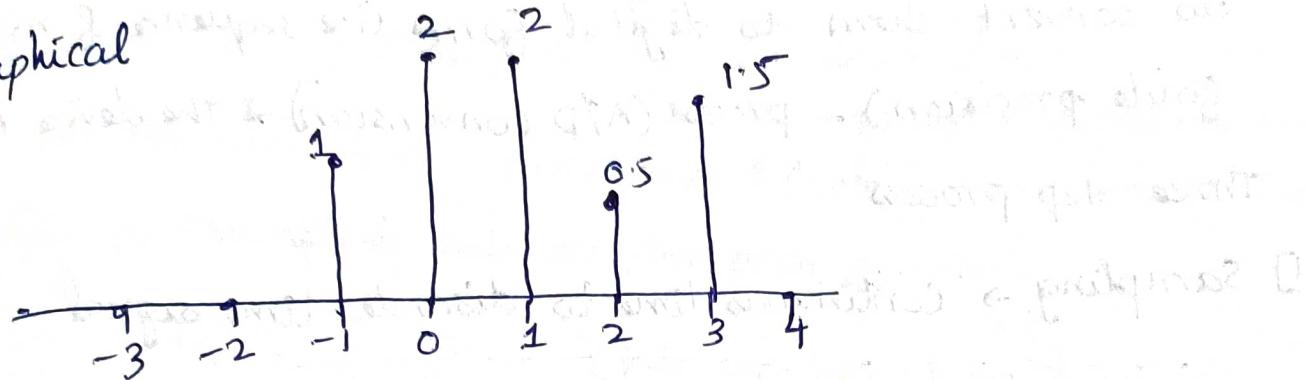


PROBLEMS

π Different representation of discrete-time signals.

$$x(-1) = 1, x(0) = 2, x(1) = 2, x(2) = 0.5, x(3) = 1.5$$

graphical



Functional

$$x(n) = \begin{cases} 1 & \text{for } n = -1 \\ 2 & \text{for } n = 0, 1 \\ 0.5 & \text{for } n = 2 \\ 1.5 & \text{for } n = 3 \\ 0 & \text{otherwise} \end{cases}$$

Tabular rep:

n	-1	0	1	2	3
x(n)	1	2	2	0.5	1.5

Sequence

$$x(n) = \{1, 2, 2, 0.5, 1.5\}$$

$$2. \sum_{n=-\infty}^{\infty} s(n-2) \sin 2n = \sin 2n \Big|_{n=2} = \sin 4 \quad \text{with } s(n-2) = 1 \text{ for } n=2 \quad 0 \text{ for } n \neq 2$$

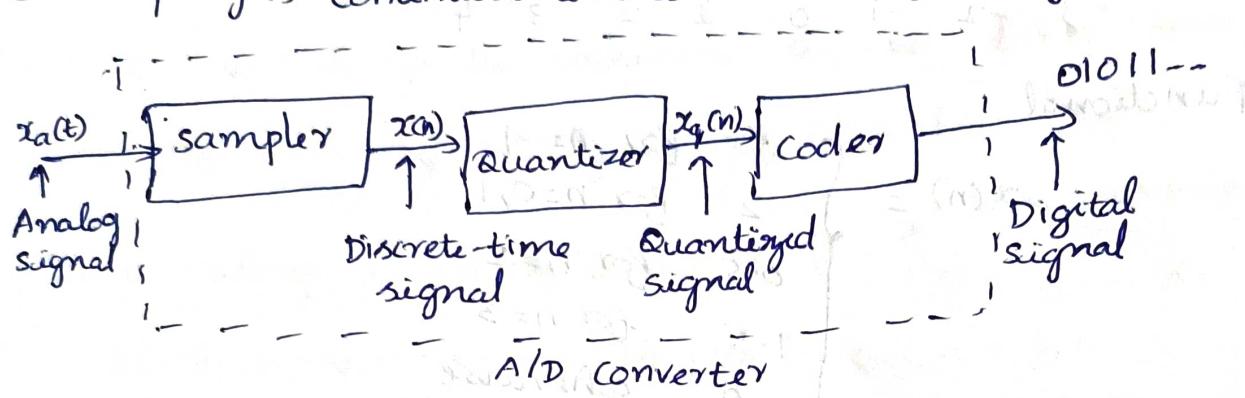
$$\sum_{n=-\infty}^{\infty} s(n+1) x(n) = x(n) \Big|_{n=-1} = x(-1)$$

ADC - Analog to Digital conversion

To process analog signal - by digital means \rightarrow necessary to convert them to digital form (i.e. sequence of nos. having finite precision) - process (A/D conversion) & the device A/D converter.

Three step process

- 1 Sampling \rightarrow continuous-time to discrete-time signal



by taking samples of continuous-time signal at discrete-time instants
 $x_a(t)$ is input to sampler $x_a(nT) = x(n)$ \Rightarrow sampling interval

- 2 Quantization \rightarrow conversion of discrete time continuous valued signal into a discrete time, discrete valued (digital) signal each signal sample is represented by a value selected from a finite set of possible values. The difference between unquantized $x(n)$ & quantized output $x_q(n)$ \rightarrow quantization error

- 3 Coding - each discrete value $x_q(n)$ represented by b-bit binary sequence \rightarrow converting to a form acceptable to computer

In many practical cases, the processed digital signal needs to be converted back to analog form - process D/A conversion

Sampling \rightarrow no loss of information or doesn't introduce distortion
 \rightarrow analog signal can be reconstructed by carefully choosing sampling rate $f_s \geq 2B$

Quantization \rightarrow irreversible process that results in signal distortion (depends on accuracy of A/D converters)

with \uparrow in accuracy or $f_s \rightarrow$ cost of A/D \uparrow

Sampling of Analog signals

Periodic or uniform sampling described by the relation

$$x(n) = x_a(nT) \quad -\infty < n < \infty$$

discrete time signal ↓ sampling period (interval btw samples)
 analog signal (taking sample every T seconds)

$F_s = 1/T$ sampling rate or sampling frequency

$$t = nT = \frac{n}{F_s} \quad (\text{relation btw. } t \& n)$$

Relationship btwn. frequency variable F for analog & f for discrete

$$x_a(t) = A \cos(2\pi F t + \theta)$$

when sampled periodically at $F_s = 1/T$ samples per second

$$x_a(nT) = x(n) = A \cos(2\pi F nT + \theta) = A \cos\left(\frac{2\pi n F}{F_s} + \theta\right)$$

$$f = \frac{F}{F_s} \quad (\text{normalized frequency})$$

For continuous time sinusoid $-\infty < F < \infty$

$$\cos(\omega_0 + 2\pi n) n + \theta = \cos(\omega_0 n + \theta)$$

For discrete time sinusoid $-1/2 < f < 1/2$ (freq separated by integer multiple of 2π → identical)

$$-\frac{1}{2} < \frac{F}{F_s} < \frac{1}{2} \Rightarrow -\frac{1}{2} = -\frac{F_s}{2} \leq F \leq \frac{F_s}{2} = \frac{1}{2T}$$

Moving infinite & orange for the variable F into a finite & range for

$$f_{\max} = \frac{F_s}{2} = \frac{1}{2T} \quad \text{highest } f \text{ in a continuous time}$$

signal that can be uniquely distinguished when sampled at

$$F_s = 1/T \text{ is } f_{\max} = \frac{F_s}{2}$$

Eg: $x_1(t) = \cos 2\pi (10)t$ which are sampled at $F_s = 40 \text{ Hz}$
 $x_2(t) = \cos 2\pi (50)t$

$$x_1(n) = \cos 2\pi \left(\frac{10}{40}\right)n = \cos \frac{\pi}{2}n$$

$$x_2(n) = \cos 2\pi \left(\frac{50}{40}\right)n = \cos \frac{5\pi}{2}n$$

$$\cos \frac{5\pi n}{2} = \cos \left(2\pi n + \frac{\pi n}{2}\right) = \cos \frac{\pi n}{2} \Rightarrow x_1(n) = x_2(n)$$

They are identical & indistinguishable: sampled values generated by $\cos(\pi f_2)n \xrightarrow{\text{ambiguity}}$ sample belong to $x_1(t)$ or $x_2(t)$

$F_2 = 50\text{Hz}$ is alias of $F_1 = 10\text{Hz}$ at $F_s = 40 \text{ samples/sec}$

$$\cos 2\pi (F_1 + 40k)t, k=1, 2, 3, 4$$

$k=1 \rightarrow 50\text{Hz}$ $k=2 \rightarrow 90\text{Hz}$ $k=3 \rightarrow 130\text{Hz}$ all alias of $F_1 = 10\text{Hz}$

$$x_a(t) = A \cos(2\pi F_k t + \theta) \quad F_k = F_0 + kF_s \quad k=\pm 1, \pm 2, \dots$$

$$x(n) = A \cos\left(2\pi \frac{F_0 + kF_s}{F_s} n + \theta\right)$$

$$= A \cos(2\pi n F_0/F_s + \theta + 2\pi kn)$$

$$= A \cos(2\pi f_0 n + \theta)$$

Highest f is $F_s/2$ then any f above F_s is mapped to equivalent f below $F_s/2$

Problem:

Consider analog signal $x_a(t) = 3 \cos 100\pi t$

a) Determine the minimum sampling rate required to avoid aliasing

$$x_a(t) = 3 \cos(2 \times 50\pi t) \quad f = 50\text{Hz} \quad \text{minimum } F_s = 100\text{Hz}$$

b) If the signal sampled at $F_s = 200\text{Hz}$ what is the discrete time signal obtained after sampling?

$$x(n) = 3 \cos \frac{100\pi}{200} n = 3 \cos \frac{\pi}{2} n$$

c) suppose that the signal is sampled at $F_s = 75\text{Hz}$. What is the discrete time signal obtained after sampling?

$$x(n) = 3 \cos \frac{100\pi}{75} n = 3 \cos \frac{4\pi}{3} n$$

$$= 3 \cos\left(2\pi - \frac{2\pi}{3}\right) n = 3 \cos \frac{2\pi}{3} n$$

d) What is the frequency $0 < F < F_s/2$ of a sinusoid that yields samples identical to those obtained in part (c)?

$$\text{For } F_s = 75 \text{ Hz} \quad F = f F_s = 75 f$$

$$\text{From part c } f = \frac{1}{3} \Rightarrow F = 25 \text{ Hz}$$

$$y_a(t) = 3 \cos 2\pi F t = 3 \cos 50\pi t$$

$$\text{when sampled at } F_s = 75 \text{ Hz} = 3 \cos \frac{2\pi}{3} n \text{ (identical to part c)}$$

Hence $F = 50 \text{ Hz}$ is an alias of $F = 25 \text{ Hz}$ for $F_s = 75 \text{ Hz}$

SAMPLING THEOREM

$$-\frac{F_s}{2} \leq F \leq \frac{F_s}{2} \quad \frac{F_s}{2} > F_{\max}$$

$$\Rightarrow F_s > 2F_{\max} \text{ (to avoid aliasing)}$$

If highest γ in analog signal $x_a(t)$ is $F_{\max} = B$ & signal is sampled at $F_s > 2F_{\max} = 2B$, then $x_a(t)$ can be recovered using the interpolation function

$$\sin 2\pi B t = g(t) = \frac{\sin 2\pi B t}{2\pi B t}$$

$$x_a(t) = \sum_{n=-\infty}^{\alpha} x_a\left(\frac{n}{F_s}\right) g\left(t - \frac{n}{F_s}\right) \quad x_a\left(\frac{n}{F_s}\right) = x_a(nT) = x(n)$$

$$x_a(nT) g(t - nT)$$

sampling of $x_a(t)$ is performed at minimum $F_s = 2B$

$$x_a(t) = \sum_{n=-\infty}^{\alpha} x_a\left(\frac{n}{2B}\right) \frac{\sin 2\pi B \left(t - \frac{n}{2B}\right)}{2\pi B \left(t - \frac{n}{2B}\right)}$$

$$F_N = 2B = 2F_{\max} \text{ is the Nyquist rate}$$

The reconstruction of $x_a(t)$ from $x_a(n)$ is complicated - involves weighted sum of interpolation function $g(t)$ & its time shifted version $g(t - nT)$ for $-\infty < n < \alpha$, weighting factors are samples of $x(n)$

- Because of complexity + infinite no. of samples required these reconstruction formulae are of theoretical interest

Problem

consider analog signal

$x_a(t) = 3 \cos 50\pi t + 10 \sin 300\pi t - \cos 100\pi t$. What is the Nyquist rate for this signal

$$F_1 = 25 \text{ Hz} \quad F_2 = 150 \text{ Hz} \quad F_3 = 50 \text{ Hz} \Rightarrow F_{\max} = 150 \text{ Hz}$$

$$F_s > 2F_{\max} = 300 \text{ Hz} \quad \text{Nyquist rate } F_N = 300 \text{ Hz}$$

Discussion:

while sampling $10 \sin 300\pi t$ at $F_N = 300 \text{ Hz}$ we will get $10 \sin \pi t = 0$
i.e. we are sampling at zero-crossing point \Rightarrow miss the signal component
solution \rightarrow sample the analog signal at rate higher than F_N

Problem

consider analog signal

$$x_a(t) = 3 \cos 2000\pi t + 5 \sin 6000\pi t + 10 \cos 12,000\pi t$$

a) what is the Nyquist rate for this signal?

$$F_1 = 1 \text{ kHz} \quad F_2 = 3 \text{ kHz} \quad F_3 = 6 \text{ kHz} \Rightarrow F_{\max} = 6 \text{ kHz}$$

$$F_s > 2F_{\max} = 12 \text{ kHz} \quad F_N = 12 \text{ kHz}$$

b) Assume that we sample this signal using a sampling rate

$F_s = 5000 \text{ samples/sec}$. What is the discrete-time signal obtained after sampling?

$$F_s = 5 \text{ kHz}, \text{ folding} \Rightarrow \frac{f_s}{2} = 2.5 \text{ kHz}$$

$$x(n) = x_a(nT) = x_a\left(\frac{n}{F_s}\right) = 3 \cos 2\pi \left(\frac{1}{5}\right)n + 5 \sin 2\pi \left(\frac{3}{5}\right)n + 10 \cos 2\pi \left(\frac{6}{5}\right)n$$

$$= 3 \cos 2\pi \left(\frac{1}{5}\right)n + 5 \sin 2\pi \left(1 - \frac{2}{5}\right)n + 10 \cos 2\pi \left(1 + \frac{1}{5}\right)n$$

$$= 3 \cos 2\pi \left(\frac{1}{5}\right)n + 5 \sin^2 \left(\frac{-2}{5}\right)n + 10 \cos 2\pi \left(\frac{1}{5}\right)n$$

$$x(n) = 13 \cos 2\pi \left(\frac{1}{5}\right)n + 5 \sin 2\pi \left(\frac{2}{5}\right)n$$

Q) What is the analog signal $y_a(t)$ we can reconstruct from the samples if we use ideal interpolation?

since $f = 1\text{ kHz} \rightarrow 2\text{ kHz}$ are present in sampled signal

the analog signal that we can recover is

$$y_a(t) = 13 \cos 2\pi(1000)t - 5 \sin 2\pi(2000)t$$

$= 13 \cos 2000\pi t - 5 \sin 4000\pi t$ which is different from due to low sampling rate \rightarrow resulting in aliasing effect.

Quantization of continuous-amplitude signals

discrete time continuous amplitude signal into digital signal

- sample value as a finite no. of digits \rightarrow quantization

error while representing discrete value level - quantization error

$$x_q(n) = Q[x(n)]$$

quantized sample → quantizer operation

$$e_q(n) = x_q(n) - x(n)$$

actual sample value

quantization error

$$\text{eg } x(n) = 0.43046721 \quad n=8 \quad x(n) = 0.4 \quad x_q(n) = 0.4$$

rounding

excess digits discarded \rightarrow truncation

rounding the resulting number

values allowed - quantization levels, $\Delta \rightarrow$ distance b/w two successive quantization levels
 rounding quantizer \rightarrow assigns to nearest quantization level
 truncation \rightarrow assigned to quantization level below it

$$-\frac{\Delta}{2} \leq e_q(n) \leq \frac{\Delta}{2}$$

error can't exceed half of the quantization step

$$\Delta = \frac{x_{\max} - x_{\min}}{L-1}$$

dynamic range
 no. of quantization levels

$$\text{eg: } x_{\max} = 1 \quad x_{\min} = 0 \quad L = 11 \quad \Delta = 0.1$$

DR - fixed then $\uparrow L \rightarrow \downarrow \Delta \rightarrow e_q \downarrow \rightarrow$ accuracy \uparrow

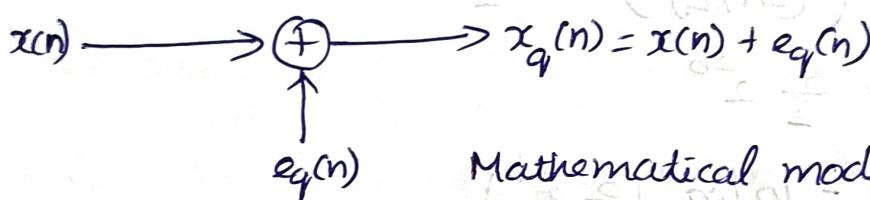
quantization \rightarrow loss of information, irreversible process
 since all sample at a distance $\Delta/2$ about certain quantization level assigned same value \rightarrow results in ambiguity

ANALYSIS OF QUANTIZATION ERROR

Deterministic analysis of quantization error is not possible since the error depends on characteristic of input signal and nonlinear nature of the quantizer.

Statistical approach is adopted to determine the effects of quantization on performance of A/D converter.

It is assumed that quantization error is random in nature and it is modelled as noise that is added to the original signal.



Effect of noise ($e_q(n)$) on desired signal is quantified by Signal to Quantization noise ratio (SQNR).

$$\text{SQNR} = \frac{P_x}{P_n}$$

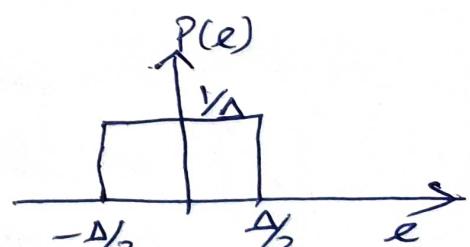
$$\text{Signal power } P_x = \sigma^2 = E[x^2(n)]$$

$$\text{Quantization noise power } P_n = \sigma^2 = E[e_q^2(n)]$$

$$P_n = \sigma_e^2 = \int_{-\Delta/2}^{\Delta/2} e^2 p(e) de = \frac{1}{\Delta} \int_{-\Delta/2}^{\Delta/2} e^2 de = \frac{\Delta^2}{12}$$

$$\text{SQNR} = 12 \frac{P_x}{\Delta^2}$$

CODING



For each quantization level, unique binary number is assigned. Each discrete value $x_q(n)$ is represented by b bits binary sequence.

$$L = 2^b \quad L \rightarrow \text{no. of quantization levels}$$

b → no. of bits

$$b = \log_2 L$$

For a given sinusoidal signal $x(t) = A \cos(\omega t)$

$$x_{\max} - x_{\min} = 2A$$

$$\Delta = \frac{x_{\max} - x_{\min}}{L} = \frac{2A}{L}$$

$$P_X = \frac{A^2}{A/2}$$

We know that $SQNR = 12 \frac{P_x}{N^2}$

$$SQNR = 12 \left(\frac{A^2 / 2}{(2A/L)^2} \right) = \frac{3}{2} L^2$$

$$SQRN = \frac{3}{2} \cdot \frac{2^b}{2}$$

$$SQNR(dB) = 10 \log_{10} \left(\frac{3}{2} 2^{2b} \right)$$

$$y = 1.76 + 6.02b$$

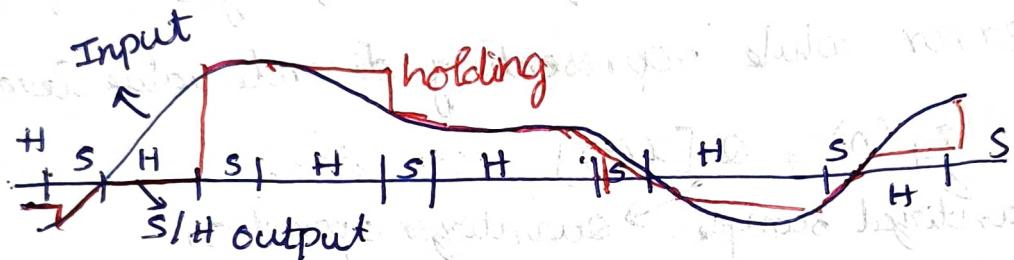
For every bit added to the word length, SNR increases approximately by 6 dB.

SAMPLE AND HOLD CIRCUIT

Sampling is performed by sample and hold (S/H) circuit which is digitally controlled analog circuit. There are two modes of operation.

Sample mode - tracks the analog input signal

Hold mode - holds the instantaneous value of the signal when the system is switched from sample to hold mode.



S/H allow A/D converter to operate more slowly. In the absence of S/H, input signal will not change more than Δ_2 during conversion.

Sampling does not introduce distortion in the conversion process but time-related degradations occur in practical devices.

OVERSAMPLING A/D CONVERTERS

The basic idea of oversampling A/D converter is to increase the sampling rate of a signal to a point where a low-resolution quantizer suffices.

By oversampling, the dynamic range of the signal values between successive samples is reduced. Thus resolution requirement of the quantizer is reduced.

For oversampling A/D converter, differential quantization is employed to reduce dynamic range between successive samples.

$$d(n) = x(n) - x(n-1)$$

The variance of $d(n)$ is

$$\sigma_d^2 = E[d^2(n)] = E[(x(n) - x(n-1))^2]$$

$$= E[x^2(n)] + E[x^2(n-1)] - 2E[x(n)x(n-1)]$$

$$= 2\sigma_x^2 - 2\sigma_x^2 R_{xx}(1) \quad \therefore E[x^2(n)] \approx E[x^2(n-1)]$$

$$\sigma_d^2 = 2\sigma_x^2 [1 - R_{xx}(1)]$$

$R_{xx}(1)$ is the value of autocorrelation sequence of $x(n)$.

For best performance $R_{xx}(1) > 0.5$ only then $\sigma_d^2 < \sigma_x^2$.

Under this condition, it is better to quantize the difference $d(n)$ and to recover $x(n)$ from the quantized value $\{d_q(n)\}$.

To obtain high correlation between successive samples of signal, we require sampling rate to be higher than Nyquist rate.

To make variance range broader and to overcome the constraint $R_{xx}(1) > 0.5$, we adopt an even better approach.

$$d(n) = x(n) - \alpha x(n-1)$$

α is the parameter selected to minimize variance in $d(n)$

$$\sigma_d^2 = E[d^2(n)] = E[(x(n) - \alpha x(n-1))^2]$$

$$= E[x^2(n)] + E[\alpha^2 x^2(n-1)] - 2E[\alpha x(n) x(n-1)]$$

$$= \sigma_x^2 + \alpha^2 \sigma_x^2 - 2\alpha \sigma_x^2 R_{xx}(1)$$

For best minimization we choose $R_{xx}(1) = \alpha$

$$\sigma_d^2 = \sigma_x^2 + \alpha \sigma_x^2 - 2\alpha \sigma_x^2$$

$$\sigma_d^2 = \sigma_x^2(1-\alpha^2)$$

$\sigma_d^2 < \sigma_x^2$ if the value of α is in the range $0 < \alpha < 1$.

For predictor of order P

$$\hat{x}(n) = \sum_{k=1}^P a_k x(n-k)$$

The goal of this predictor is to provide an estimate

$\hat{x}(n)$ of $x(n)$ from a linear combination of past values of $x(n)$.

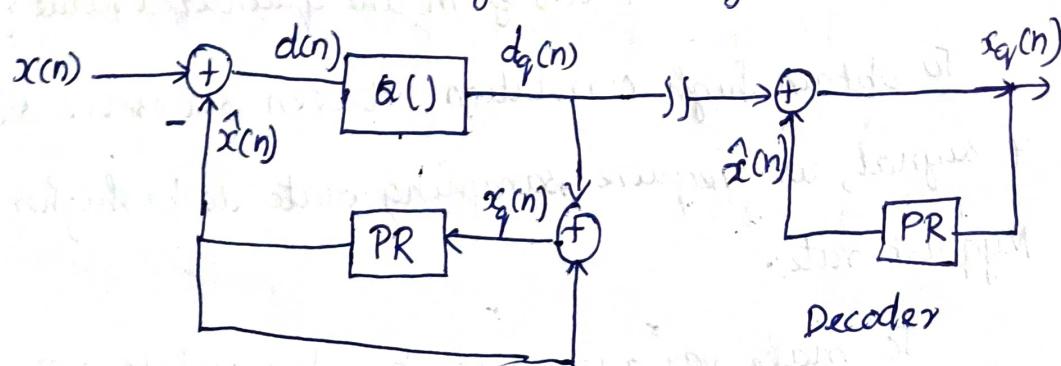
Quantization error is given by

$$e(n) = d(n) - d_q(n)$$

$$= x(n) - \hat{x}(n) - d_q(n)$$

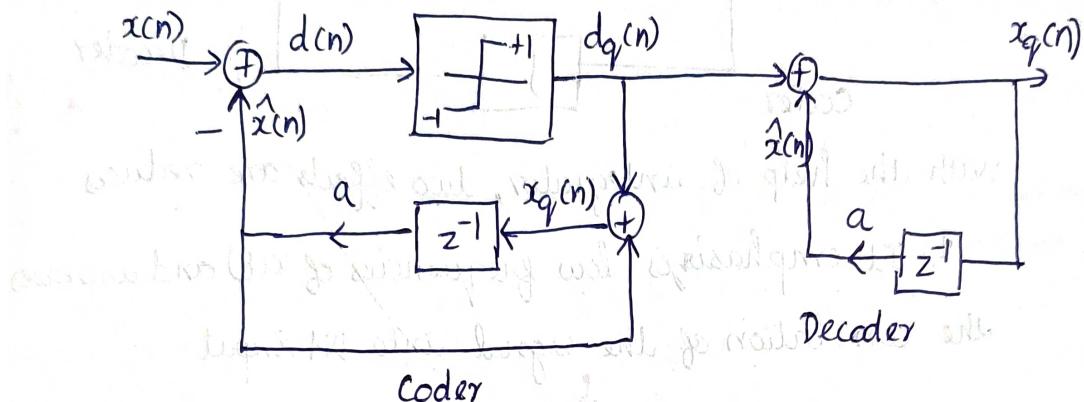
$$e(n) = x(n) - x_q(n)$$

Differential Predictive signal Quantizer



The decoder reconstructs the signal from the quantized value.

The simplest form of differential predictive quantization is called delta modulation (DM). This is also called as 1-bit differential quantizer.



In Delta modulation, the quantizer is a simple 1-bit quantizer and the predictor is a first-order predictor.

DM provides staircase approximation of input signal. At every sampling instant, the difference between input sample $x(n)$ and its most recent staircase approximation $\hat{x}(n) = a x_q(n-1)$ is determined and then the signal is updated by a step Δ in the direction of difference.

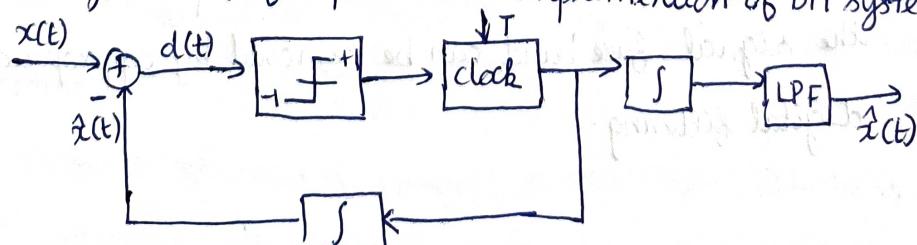
$$x_q(n) = d_q(n) + a x_q(n-1)$$

which is discrete-time equivalent of analog integrator

If $a=1 \rightarrow$ ideal accumulator (Integrator)

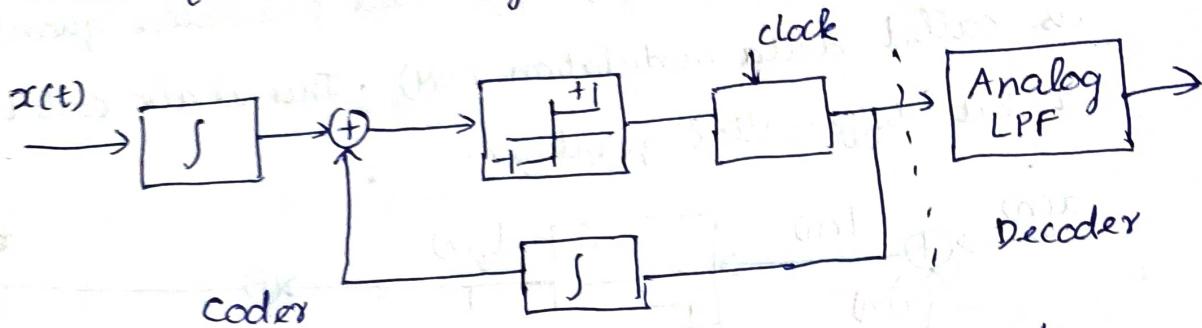
If $a < 1 \rightarrow$ leaky integrator.

Analog model for practical implementation of DM system



Analog LPF is necessary for rejection of out-of-band components.

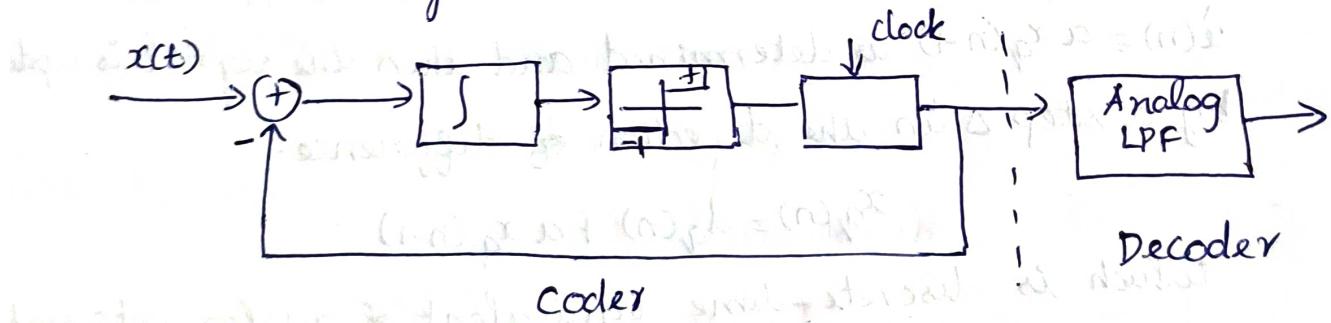
There are two types of quantization error [Granular noise and slope overload distortion]. These distortion can be avoided if we use integrator in front of DM.



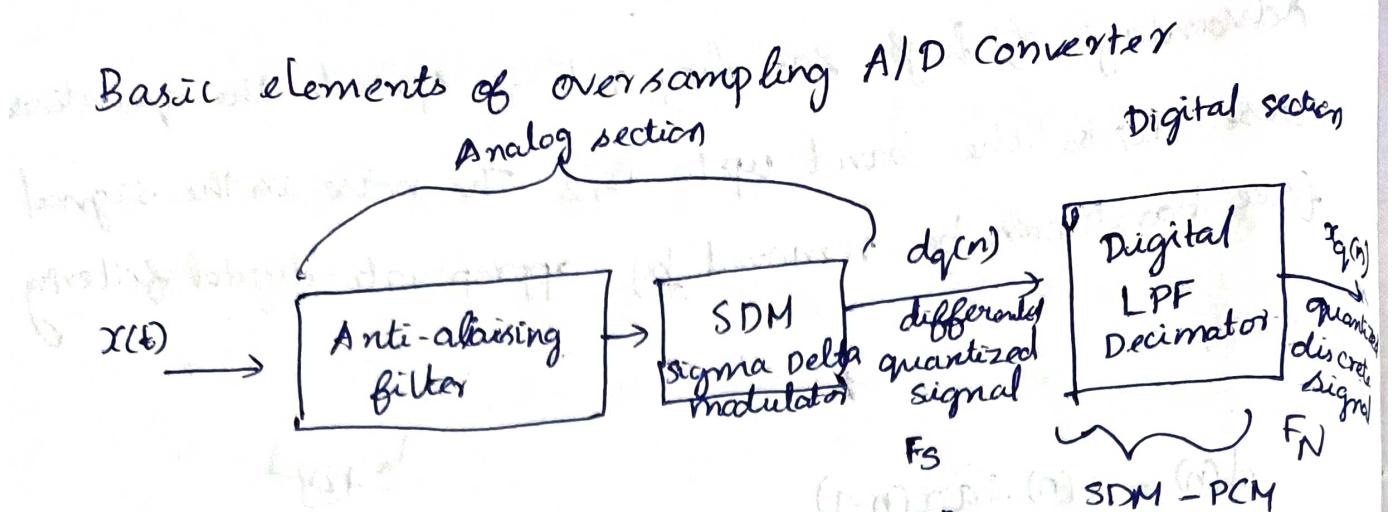
With the help of integrator, two effects are reduced

1. It emphasizes low frequencies of $x(t)$ and increases the correlation of the signal into DM input.
2. Simplifies DM decoder - the differentiator required at decoder is cancelled by DM integrator. Hence decoder is simply LPF.

The two integrator at the encoder can be replaced by a single integrator placed before the comparator. This system is known as Sigma-Delta modulation (SDM).

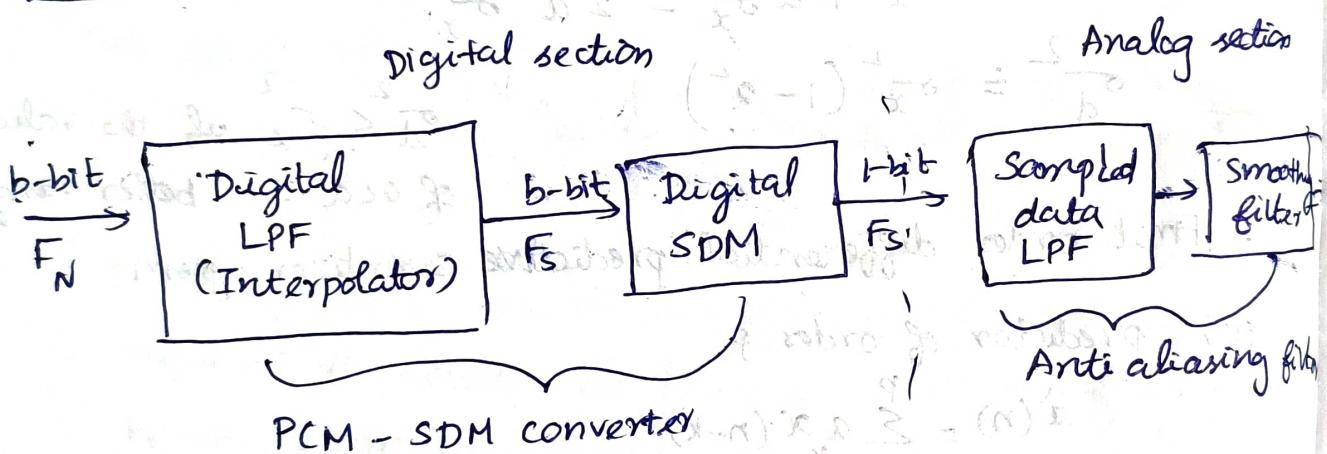


SDM is an ideal candidate for A/D converter. It takes the advantage of high sampling rate and spreads the quantization noise across the band upto $F_s/2$. The noise in the signal free band can be removed by appropriate digital filtering.

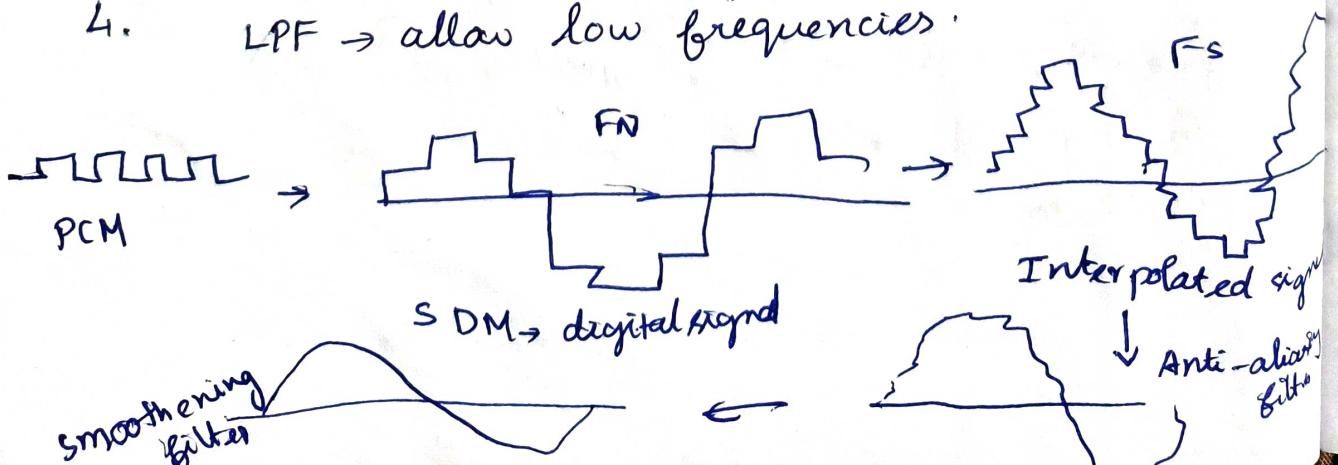


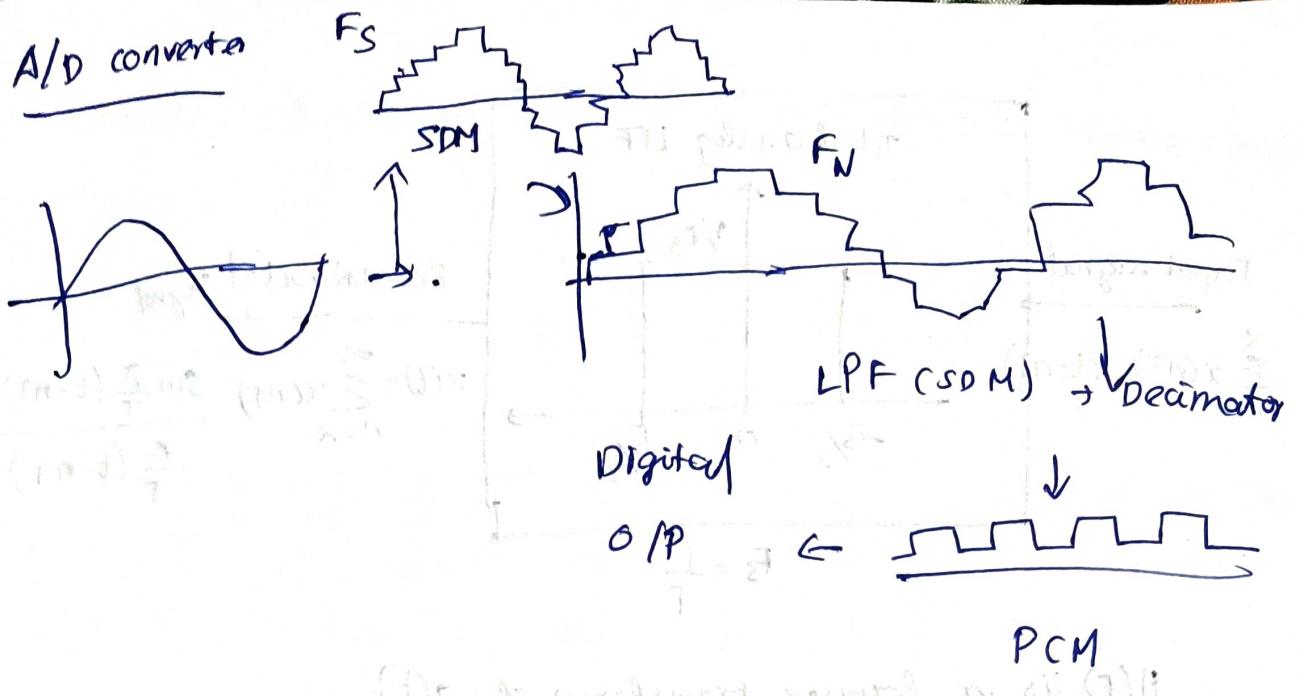
1. $F_s \rightarrow$ sampling frequency $F_s >> 2B$
2. $F_N \rightarrow$ Nyquist frequency $F_N = 2B$
3. Decimator or down sampling reduces the frequency from f_s to f_N for better bit rate efficiency.

D/A converter



1. $F_N \rightarrow$ Nyquist $= 2B = 2 \times$ signal γ
2. Interpolator upsample the signal from f_N to $f_s >> 2^B$
3. PCM \rightarrow acts as encoder / decoder
4. LPF \rightarrow allow low frequencies.





DIGITAL TO ANALOG CONVERSION

When band limited lowpass analog signal sampled at Nyquist rate, it can be reconstructed without distortion.

The ideal reconstruction formula is

$$x(t) = \sum_{n=-\infty}^{\infty} x(nT) \frac{\sin(\pi/T)(t-nT)}{(\pi/T)(t-nT)}$$

where $T = 1/F_s = 1/2B$

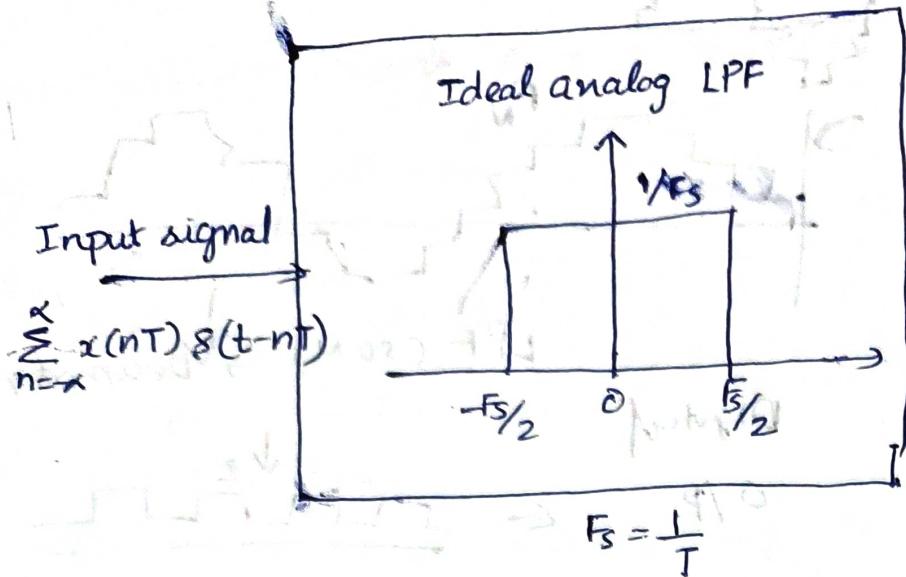
The ideal interpolation function is

$$g(t) = \frac{\sin(\pi t/T)}{\pi t/T}$$

The interpolation formula for $x(t)$ is linear superposition of time-shifted version of $g(t)$, with each $g(t-nT)$ weighted by corresponding signal sample $x(nT)$.

Alternatively, reconstruction of signal can also be considered as a linear filtering process where sampled signal excites an analog filter. The frequency response of ideal analog filter is

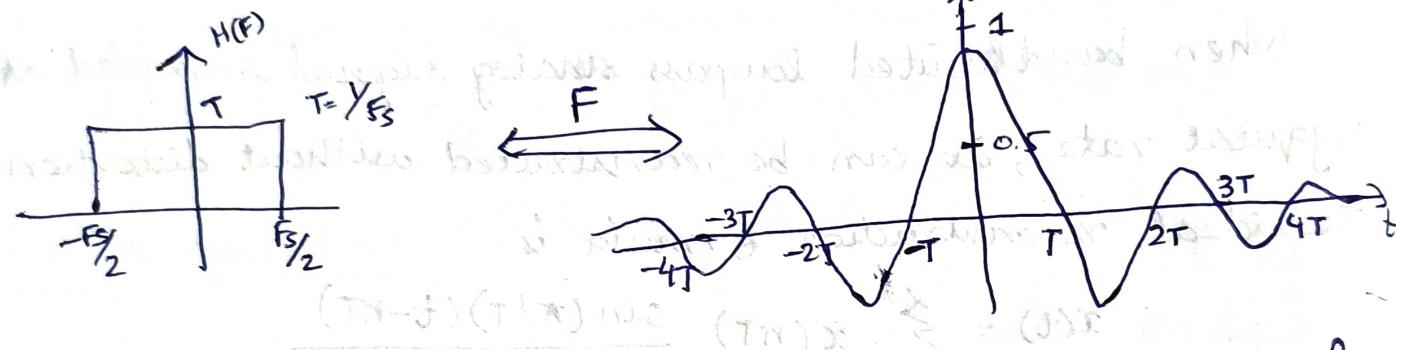
$$H(F) = \begin{cases} T, & |F| \leq \frac{1}{2T} = \frac{F_s}{2} \\ 0, & |F| > \frac{1}{2T} \end{cases}$$



Reconstructed signal

$$x(t) = \sum_{n=-\infty}^{\infty} x(nT) \frac{\sin \frac{\pi}{T}(t-nT)}{\frac{\pi}{T}(t-nT)}$$

$H(F)$ is a Fourier transform of $g(t)$

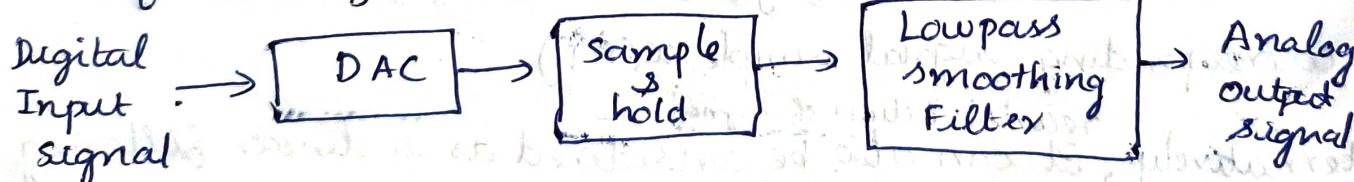


The ideal reconstruction filter is an ideal LPF & its impulse response extends for all time. Hence it is non-causal & physically non-realizable.

SAMPLE & HOLD (S/H)

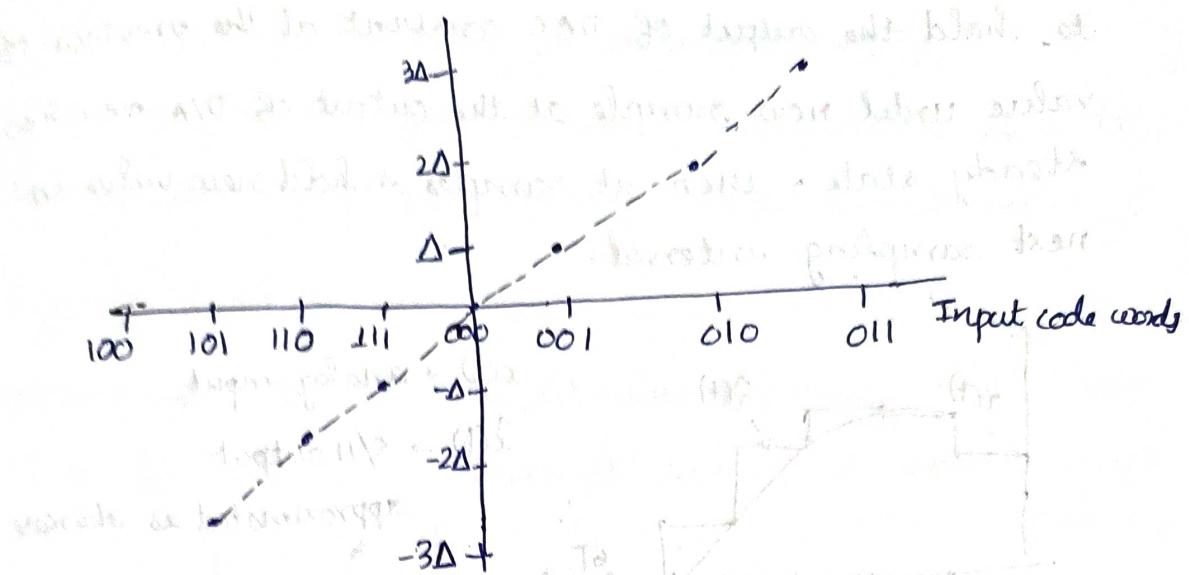
D/A conversion is performed by D/A converter with S/H

& followed by a LPF



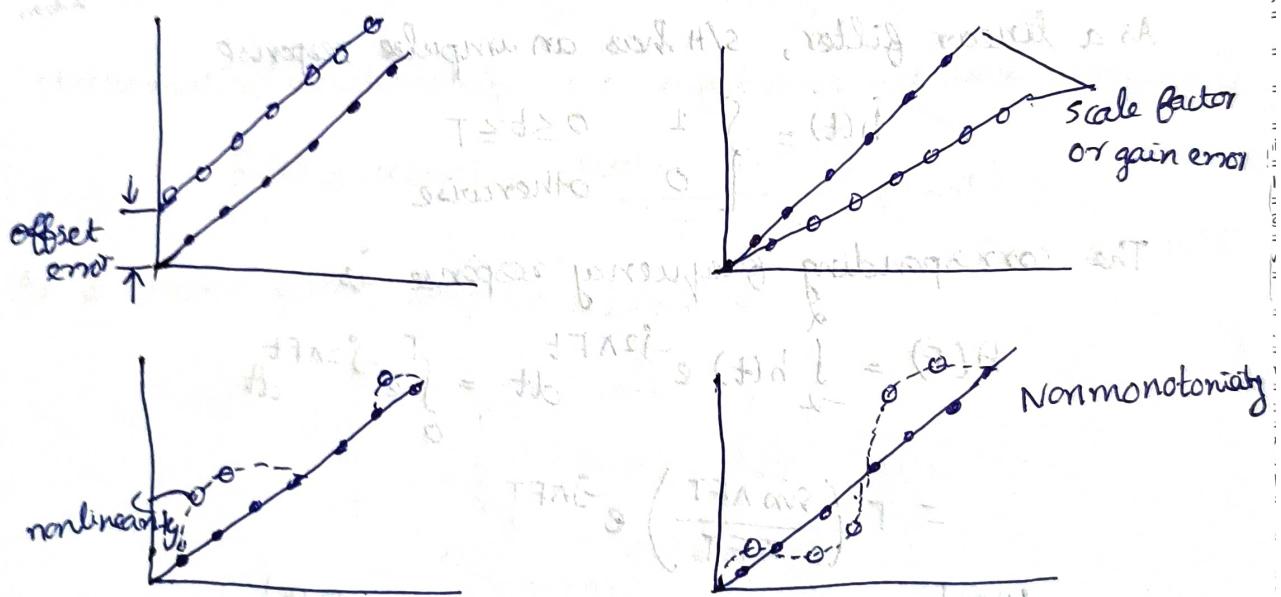
The input to D/A converter is electrical signal that corresponds to binary word & produces output voltage or current that is proportional to the value of binary word.

Input \Rightarrow output characteristics of DAC for 3-bit bipolar signals.



The line connecting dots is straight line through origin whereas in practical DAC, the line may deviate from ideal case.

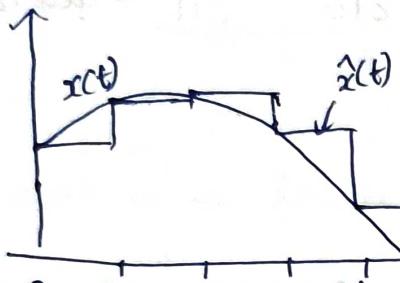
The different types of errors are



Settling time \Rightarrow time required for the output of D/A converter to reach & remain within a given fraction of final value ($\pm \frac{1}{2} \text{ LSB}$), after application of input code word.

Glitch \Rightarrow a high amplitude transient that occurs after application of input code word especially when two consecutive code words to A/D differ by several bits.

S/H circuit serve as a deglitcher. The basic task is to hold the output of DAC constant at the previous value until new sample at the output of D/A reaches steady state, then it samples & hold new value in next sampling interval.



$x(t)$ → analog input

$z(t)$ → S/H output

approximated as staircase

Takes signal sample from DAC & holds it for T seconds

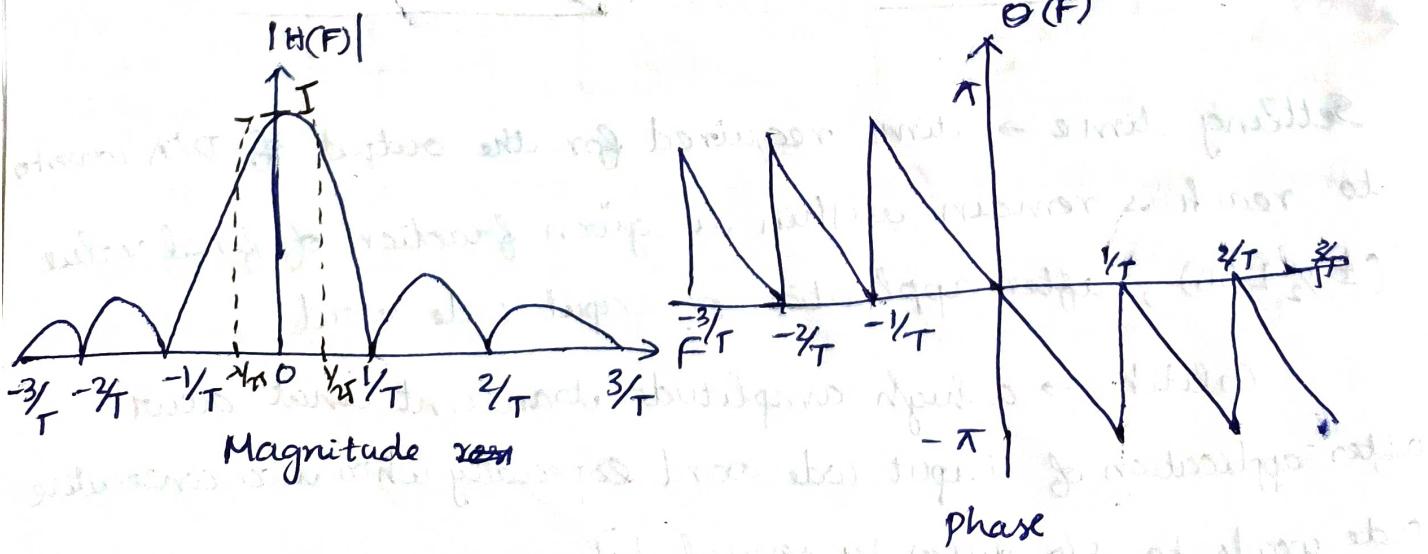
when next sample arrives, it jumps to next value & holds for T sec

As a linear filter, S/H has an impulse response

$$h(t) = \begin{cases} 1 & 0 \leq t \leq T \\ 0 & \text{otherwise} \end{cases}$$

The corresponding frequency response is

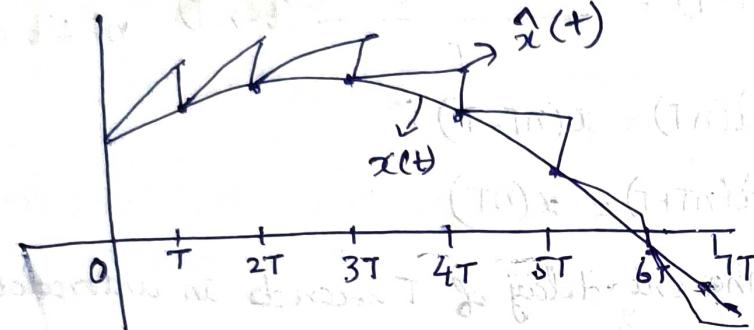
$$\begin{aligned} H(F) &= \int_{-\infty}^{\infty} h(t) e^{-j2\pi F t} dt = \int_0^T e^{-j2\pi F t} dt \\ &= T \left(\frac{\sin \pi F T}{\pi F T} \right) e^{-j\pi F T} \end{aligned}$$



S/H does not possess sharp cut-off frequency response characteristics
 so passes aliased components to its output ($\gamma > f_{s/2}$)
 Hence LPF is used \rightarrow highly attenuates & components above $f_{s/2}$
 & smooths the signal by removing sharp discontinuities.

First-order hold

Approximates $x(t)$ by straight-line segments which have
 a slope determined by current sample $x(nT)$ & previous sample $x(nT-T)$

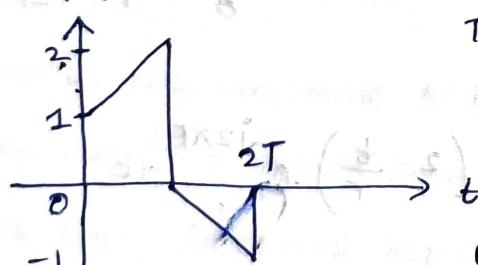


Mathematical relationship b/wn input & output waveform is

$$\hat{x}(t) = x(nT) + \frac{x(nT) - x(nT-T)}{T} (t-nT) \quad nT \leq t < (n+1)T$$

As a linear filter, impulse response is

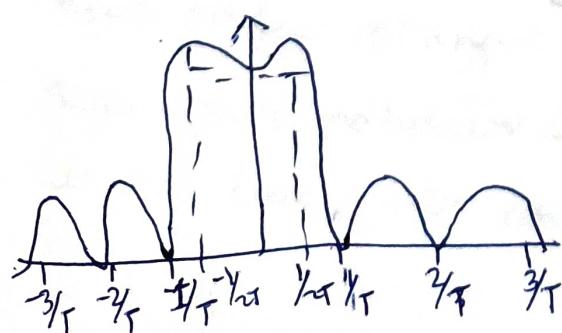
$$h(t) = \begin{cases} 1 + \frac{t}{T}, & 0 \leq t \leq T \\ 1 - \frac{t}{T}, & T \leq t < 2T \\ 0, & \text{otherwise} \end{cases}$$



Taking Fourier transform

$$H(F) = T \left(1 + 4\pi^2 F^2 T^2\right)^{1/2} \left(\frac{\sin \pi F T}{\pi F T}\right)^2 e^{j\theta(F)}$$

$$\text{where } \theta(F) = -\pi F T + \tan^{-1} 2\pi F T$$



First-order hold exhibits better frequency response characteristics

Linear Interpolation with delay

First order hold - attempts to linearly predict next sample based on $x(nT) \approx x(nT-T)$. The estimated $\hat{x}(t)$ contains jumps at the sample points:

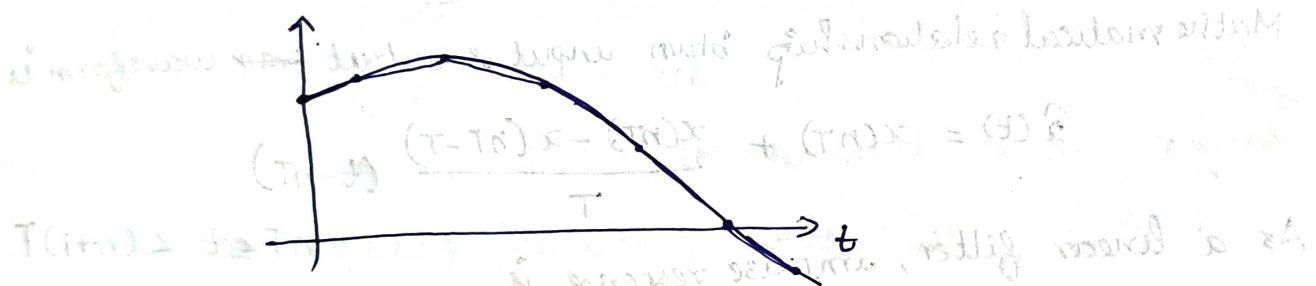
By providing one sample delay, the successive sample points can be connected by straight line segments. The interpolated signal $\hat{x}(t)$ is expressed as

$$\hat{x}(t) = x(nT-T) + \frac{x(nT) - x(nT-T)}{T} (t-nT) \quad nT \leq t < (n+1)T$$

at $t=nT$, $\hat{x}(nT) = x(nT-T)$

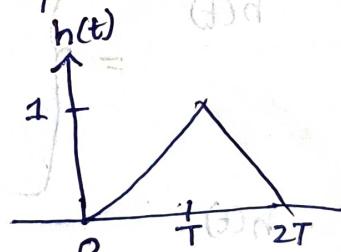
at $t=nT+T$, $\hat{x}(nT+T) = x(nT)$

It has an inherent delay of T seconds in interpolating actual signal



As a linear filter, the impulse response is

$$h(t) = \begin{cases} t/T & 0 \leq t < T \\ 2-t/T & T \leq t < 2T \\ 0 & \text{otherwise} \end{cases}$$



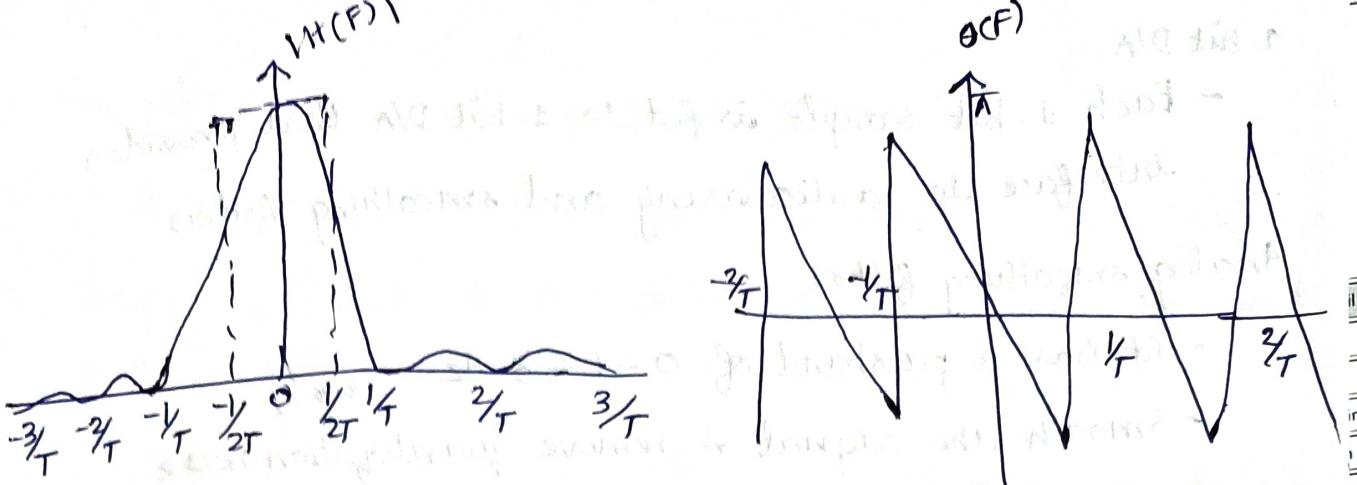
Frequency response is

$$H(F) = \int_0^T \frac{t}{T} e^{-j2\pi F t} dt + \int_T^{2T} \left(2 - \frac{t}{T}\right) e^{-j2\pi F t} dt$$

$$= T \left(\frac{\sin \pi F T}{\pi F T} \right)^2 e^{-j2\pi F T}$$

Had distortion with overshoot?

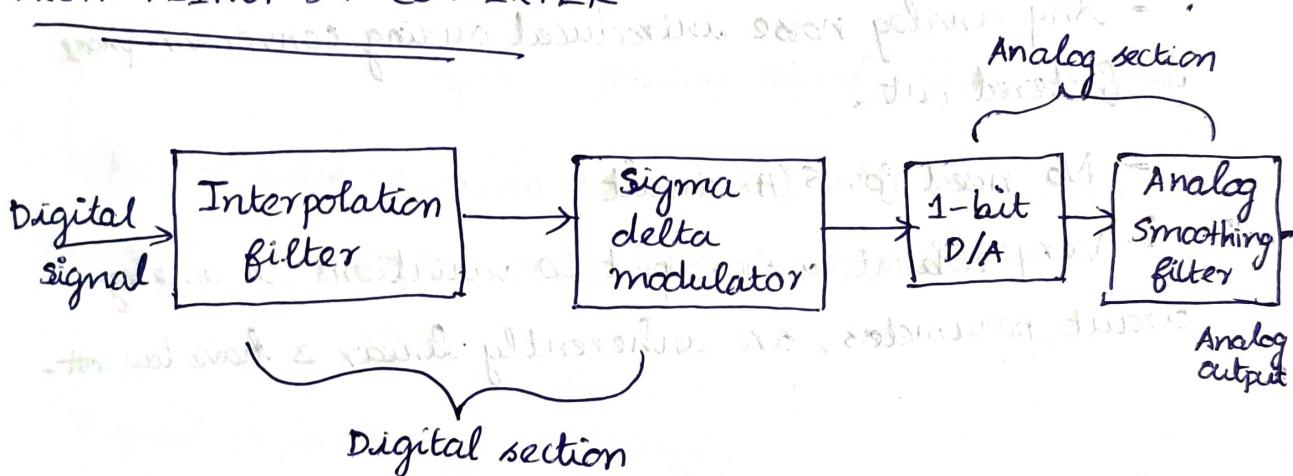
What happens if we prefilter?



$|H(F)|$ the curve falls off rapidly
 $\theta(F)$ is linear due to delay T
 \therefore contains small sidelobes beyond F

Following interpolator, LPF with sharp cutoff frequency is used
 to remove high frequency components in $\hat{x}(t)$.

OVERSAMPLING D/A CONVERTER



Oversampling D/A converter subdivided into

- digital front end
- analog section

Interpolator - increase sampling rate by a factor I by inserting $I-1$ zeros between successive low rate samples.

It is then processed by a digital filter ($f_c = B/F_s$) to reject replicas of input signal spectrum.

Sigma-delta modulator - the higher rate signal is then fed to SDM that creates noise shaped 1-bit sample.

1-bit D/A

- Each 1-bit sample is fed to 1-bit D/A that provides interface to antialiasing and smoothing filters.

Analog smoothing filter

- It has a passband of $0 \leq F \leq B$ Hz
- Smooth the signal & remove quantization noise in the frequency band $B \leq F \leq F_s/2$

Advantage of oversampling A/D or D/A converter:

- High sampling rate & subsequent filtering minimizes or removes the need for complex and expensive analog antialiasing filter

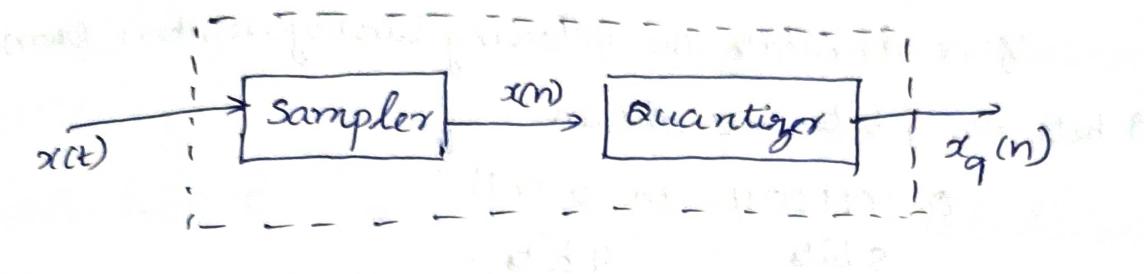
- Any analog noise introduced during conversion phase is filtered out.

- No need for S/H circuit

- Very robust with respect to variations in analog-circuit parameters, are inherently linear & have low cost.

Quantization Noise:

The process of converting an analog signal to digital is given below



Signal $x(t)$ is sampled at an interval $T = nT$, where $n=0, 1, 2$ to create sequence $x(n)$ → sampler

Each sample $x(n)$ is expressed by a finite number of bits by Quantizer to get quantized sequence $x_q(n)$

The difference between $x_q(n)$ & $x(n)$, is called quantization noise $e(n) = x_q(n) - x(n)$

Let us assume sinusoidal signal varying between $+1 \times -1$ having dynamic range 2. To convert sinusoidal signal ($b+1$) bits are employed including sign bit

No. of quantization levels $\Delta L = 2^{b+1}$

Dynamic range $x_{\max} - x_{\min} = 2$

The quantization step size $\Delta = \frac{x_{\max} - x_{\min}}{L}$

$$\text{Case no. of digits as } b=3 \text{ bits } \Delta = \frac{2}{2^{b+1}}$$

$$\Delta = \frac{2}{2^4} = \frac{1}{16}$$

If $b=3$ bits then $q = 2^{-3} = 0.125$

The most common methods of Quantization are

1. Truncation

2. Rounding

Truncation:

It is a process of discarding all bits less significant than least significant bit that is retained.

If we truncate the following binary numbers from 8 bits to 4 bits, we obtain

$$0.00110011 \text{ to } 0.0011$$

8 bits 4 bits

When we truncate a number, the signal value is approximated by highest quantization level that is not greater than signal

Rounding:

Rounding of a number of b bits is accomplished by choosing the rounded result as b bit number closest to original number unrounded.

e.g.: 0.11010 rounded to 3 bits is either 0.110 or 0.111

Rounding up or down will have negligible effect on accuracy of computation.

Error due to truncation and rounding

TRUNCATION:

If quantization method is truncation, the number is approximated by the nearest level that doesn't exceed it.

In this case the error $x_T - x$ is negative or zero

where x_T is truncation value of x and it is assumed $|x| \leq 0$

The error made by truncating a number to b bits following binary point satisfies the inequality,

$$0 \geq x_T - x > -2^{-b} \quad \rightarrow ①$$

Eg: consider decimal number 0.12890625

Its binary equivalent is 0.00100001

Truncating to 4 bits $x_T = (0.0010)_2 \approx (0.125)_{10}$

$$\text{error } x_T - x = 0.12890625 - 0.125 = -0.00390625$$

$$-2^{-b} = -2^{-4} = -0.0625 \text{ which satisfies (1)}$$

Eqn (1) holds for both sign-magnitude, one's complement and two's complement if $x > 0$

If $x < 0$, we have to verify if it holds for all types of representation.

considering two's complement representation, the magnitude of the negative number is same as the positive number

$$x = 1 - \sum_{i=1}^b c_i 2^{-i}$$

If we truncate the number to N bits, then

$$x_T = 1 - \sum_{i=1}^N c_i 2^{-i}$$

The change in magnitude

$$\begin{aligned} x_T - x &= \sum_{i=1}^b c_i 2^{-i} - \sum_{i=1}^N c_i 2^{-i} \\ &= \sum_{i=N+1}^b c_i 2^{-i} \geq 0 \end{aligned}$$

We find that due to truncation, the change in magnitude is positive, which implies that error is negative and satisfy the inequality $0 \geq x_T - x > -2^{-b}$ (2)

Now consider one's complement representation, the magnitude of negative number with b bits is given by

$$x = 1 - \sum_{i=1}^b c_i 2^{-i} - 2^{-b}$$

When the number is truncated to N bits, then

$$x_T = 1 - \sum_{i=1}^N c_i 2^{-i} - 2^{-N}$$

The change in magnitude due to truncation is

$$x_T - x = \sum_{i=1}^b c_i 2^{-i} - (2^{-N} - 2^{-b}) < 0$$

Therefore the magnitude decreases with truncation which implies that error is positive & satisfy the inequality

$$0 \leq x_T - x < 2^{-b}$$

This inequality holds for sign magnitude representation

In floating point systems the effect of truncation is visible only in mantissa. Let the mantissa is truncated to N bits.

If $x = 2^c \cdot M$, then

$$x_T = 2^c \cdot M_T$$

$$\text{Error } e = x_T - x = 2^c (M_T - M)$$

From ①, with 2's complement representation of mantissa, we have

$$0 \geq M_T - M > -2^{-b}$$

$$0 \geq \frac{e}{2^c} > -2^{-b}$$

$$0 \geq e > -2^{-b} 2^c \quad \text{--- ③}$$

We define relative error $\epsilon = \frac{x_T - x}{x} = \frac{e}{x}$

③ can be written as

$$0 \geq \epsilon x > -2^{-b} 2^c$$

$$\text{or } 0 \geq \epsilon 2^c M > -2^{-b} 2^c$$

$$\text{or } 0 \geq \epsilon M > -2^{-b}$$

If $M = \frac{1}{2}$, the relative error is maximum

Therefore rounding off products for binary numbers

$$0 \leq e > -2 \cdot 2^{-b}$$

for $M = \frac{1}{2}$

If $M = -\frac{1}{2}$, the relative error range is

$$0 \leq e < 2 \cdot 2^{-b}$$

In one's complement representation, the error for truncation of positive values of the mantissa is

$$0 \geq M_T - M > -2^{-b}$$

or

$$0 \geq e > -2^{-b} \cdot 2^C$$

with $e = ex = e \cdot 2^C \cdot M \propto M = \frac{1}{2}$

We get maximum range of the relative error for positive M .

$$0 \geq e > -2 \cdot 2^{-b}$$

For negative mantissa values, the error is

$$0 \leq M_T - M < 2^{-b}$$

or

$$0 \leq e < 2 \cdot 2^{-b}$$

With $M = -\frac{1}{2}$, the maximum range of the relative error for negative M is

$$0 \geq e > -2 \cdot 2^{-b}$$

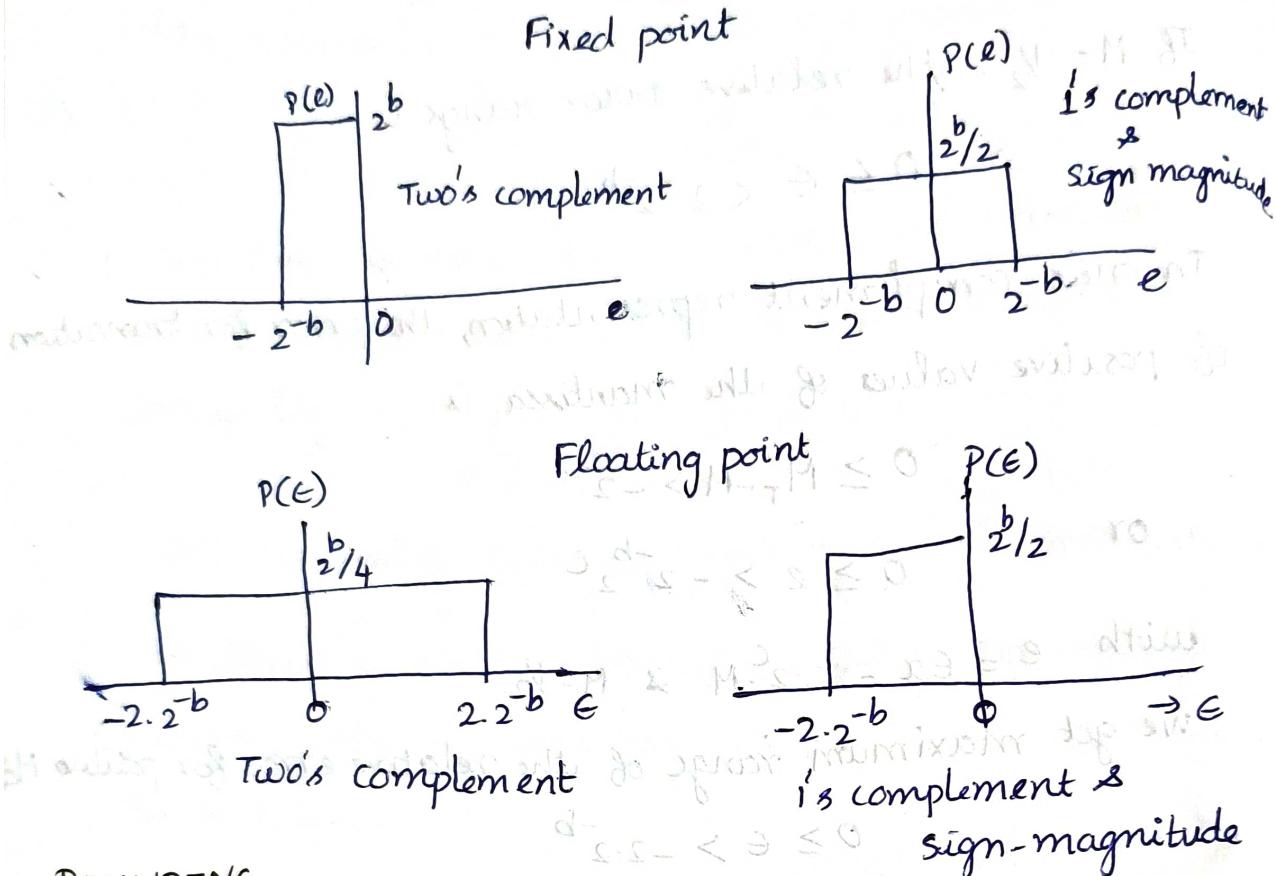
which is same as positive M .

So we can say that rounding off products for binary numbers with one's complement representation gives minimum error.

Maximum error is given by $-2 \cdot 2^{-b}$.

Similarly, with two's complement binary products for rounding off products for binary numbers give minimum error.

The probability density function $P(e)$ for truncation of fixed point and floating point numbers.



ROUNDING

In fixed point arithmetic the error due to rounding a number to b bits produces an error $e = x_T - x$ which satisfies the inequality

$$-\frac{b}{2} \leq x_T - x \leq \frac{b}{2}$$

This is because with rounding, if the value lies half way between two levels, it can be approximated to either nearest higher level or by the nearest lower level.

This inequality satisfies regardless of sign-magnitude or one's complement used for negative numbers.

In floating point arithmetic, only the mantissa is affected by quantization.

If $x = M \cdot 2^c$

and $x_T = M_T \cdot 2^c$

then $e = x_T - x = (M_T - M) 2^c \quad \epsilon = \frac{x_T - x}{x} = \frac{e}{x}$

But for rounding

$$-\frac{2^{-b}}{2} \leq M_T - M \leq \frac{2^{-b}}{2}$$

$$-2^c \frac{-2^{-b}}{2} \leq x_T - x \leq 2^c \frac{2^{-b}}{2}$$

or

$$-2^c \frac{-2^{-b}}{2} \leq \epsilon x \leq 2^c \frac{2^{-b}}{2}$$

we have

$$x = 2^c \cdot M$$

$$\text{then } -2^c \frac{-2^{-b}}{2} \leq \epsilon \cdot 2^c \cdot M \leq 2^c \frac{2^{-b}}{2}$$

which gives

$$-\frac{2^{-b}}{2} \leq \epsilon \cdot M \leq \frac{2^{-b}}{2}$$

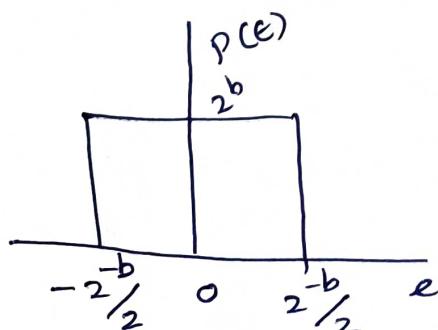
The Mantissa satisfies

$$\frac{1}{2} \leq M < 1$$

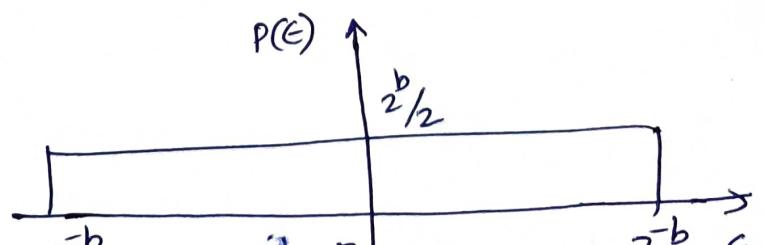
If $M = \frac{1}{2}$ we get maximum range of relative error

$$-2^{-b} \leq \epsilon \leq 2^{-b}$$

The probability density function for rounding is



Fixed point



Floating point