

Quantization Noise:

For most of the engineering applications, input signal is continuous time signal (or) analog signal.

The signal is to be converted into digital by using ADC. (Sampling and quantization). The signal $x(t)$ is first sampled at regular intervals $t=nT$ where $n=0, 1, 2, \dots$ to create a sequence $x(n)$. This is done by sampler.

Then numeric equivalent of each sample $x(n)$ is expressed by a finite no. of bits giving $x_q(n)$.

$$\left. \begin{array}{l} \text{Quantisation noise} \\ (\text{or}) \text{ error} \\ (\text{or}) \\ \text{ADC noise} \end{array} \right\} e(n) = x_q(n) - x(n)$$

If $(b+1)$ bits are used to represent each sequence (including sign bit), then no. of levels available for quantizing $x(n)$ is $\underline{2^{b+1}}$.

Thus the interval between successive level is said to be quantization step size (q).

$$q = \frac{2}{2^{b+1}} = \frac{2}{2^b \cdot 2} = 2^{-b}$$

e.g.: when $b=3$, $q = 2^{-3} = 0.125$

The common methods of quantization are :

1) Truncation

2) Rounding.

Truncation:

→ process of discarding all bits less significant than LSB that is retained. Suppose if we truncate binary numbers from 8 bits to 4 bits,

$$0.00110011 \Rightarrow 0.0011$$

$$1.01001001 \Rightarrow 1.0100$$

When it is truncated, its value is approximated by the highest quantization level that is not greater than the signal.

Rounding:

→ process of rounding the number closest to the original number unrounded.

eg: 0.11010101 $\overset{8\text{-bits}}{\Rightarrow} \overset{4\text{-bits}}{0.1101}$ (or) $\overset{4\text{-bits}}{0.1110} \leftarrow \cancel{\overset{8\text{-bits}}{0.1101}}$

⇒ have negligible effect on accuracy of computation.

Error due to truncation of fixed point numbers:

If a binary number is truncated, its value is approximated by the nearest level that does not exceed it.

i.e. e.g.: decimal number $\Rightarrow (0.12890625)_{10}$

binary equivalent $\Rightarrow (0.00100001)_2 = x$

b=4: if x is truncated to 4-bits,

$$\xrightarrow{\text{truncated value}} x_T = (0.0010)_2 = (0.125)_{10}$$

$$q = x_T - x = (-0.00390625)_{10} \Leftrightarrow -2^{-4} = (-0.0625)_{10}$$

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

here, $|x| < 0$ (assumption)

$$\therefore \boxed{0 > (x_T - x) > -2^{-b}} \quad \begin{matrix} \textcircled{1} \\ \Leftarrow \text{inequality condn.} \end{matrix}$$

This condition is satisfied for sign-magnitude, 1's complement and 2's complement if $x > 0$ (+ve, fraction).

We have to find whether the above inequality condition is satisfied or not.

Let us consider first the two's complement representation,

22.

the magnitude of the negative number is given by,

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

If we truncate x to N' bits, then

$$x_T = 1 - \sum_{i=1}^{N'} d_i 2^{-i} \rightarrow 3$$

The change in magnitude is,

$$x_T - x = \sum_{i=1}^b d_i 2^{-i} - \sum_{i=1}^{N'} d_i 2^{-i}$$

$$= \sum_{i=N+1}^b d_i 2^{-i} \geq 0 \rightarrow 4$$

$x_T - x > 0 \Rightarrow$ ~~change in magnitude~~ is +ve due to truncation,

which implies that error is negative and satisfy the inequality,

$$0 \geq (x_T - x) > -2^{-b} \quad \begin{matrix} \rightarrow 2 \\ 2^1's. \text{ complement} \\ \text{error} \end{matrix}$$

For 1's complement representation, the magnitude of negative number with ' b ' bits is given by,

$$x = 1 - \sum_{i=1}^b d_i 2^{-i} - 2^{-b} \rightarrow 5$$

when the number is truncated to ' N' ' bits, then

$$x_T = 1 - \sum_{i=1}^{N'} d_i 2^{-i} - 2^{-N} \rightarrow 6$$

23.

The change in magnitude is,

$$x_T - x = \sum_{i=N}^b d_i 2^{-i} - (2^{-N} - 2^{-b}) < 0 \quad \rightarrow ⑧$$

\therefore the magnitude decreases with truncation which implies that error is positive and satisfy the inequality,

$$0 \leq (x_T - x) < 2^{-b} \quad \text{⑨}$$

\leftarrow 1's complement & sign-magnitude error

The above condition holds for sign-magnitude representation also.

Range of errors in truncation of fixed pt nos.:

number & its representation	error range
-----------------------------	-------------

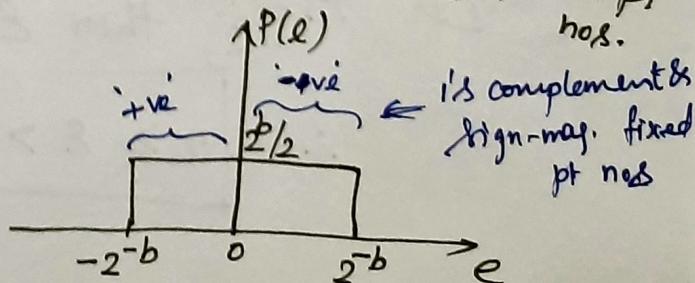
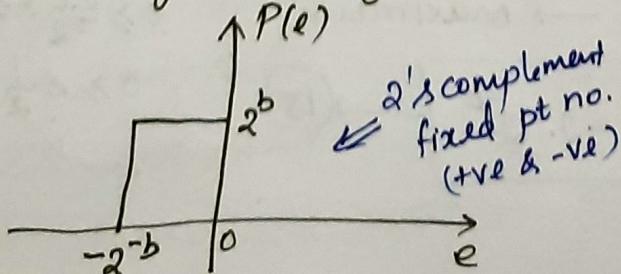
positive number (all 3 rep.)	$0 \geq e > -2^{-b}$
------------------------------	----------------------

sign-magnitude negative no.	$0 \leq e < 2^{-b}$
-----------------------------	---------------------

1's complement " "	$0 \leq e \leq 2^{-b}$
--------------------	------------------------

2's complement " "	$0 \geq e > -2^{-b}$
--------------------	----------------------

probability density function $P(e)$ for truncation of fixed pt nos.



Errors due to truncation of floating point numbers:

In floating point systems, the truncation effect is visible only in the mantissa.

Let the mantissa is truncated to 'n'-bits,

If $x = M \cdot 2^{\text{exp}}$, then

$$x_T = 2^{\text{exp}} M_T$$

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

(10)

From equation number ⑤, with 2's complement representation of mantissa, we have

$$0 \geq M_T - M > -2^{-b} \quad (11)$$

$$0 \geq e > -2^{-b} \cdot 2^{\text{exp}} \quad (12)$$

$$\therefore e = 2^{\text{exp}} (M_T - M)$$

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

$$\Rightarrow \boxed{e = \epsilon x}$$

$$\therefore (12) \Rightarrow 0 \geq \epsilon x > -2^{-b} \cdot 2^{\text{exp}} \rightarrow (13)$$

$$0 \geq \epsilon \cdot 2^{\text{exp}} M > -2^{-b} \cdot 2^{\text{exp}} \quad \because x = 2^{\text{exp}} \cdot M$$

$$0 \geq \epsilon \cdot M > -2^{-b} \rightarrow (14)$$

If $M = \frac{1}{2}$, then $\epsilon \rightarrow \text{maximum}$

$$0 \geq \epsilon \cdot \frac{1}{2} > -2^{-b}$$

$$\therefore \boxed{0 \geq \epsilon > -2 \cdot 2^{-b}} \rightarrow (15)$$

$$0 \geq \epsilon \cdot \frac{1}{2} > -2 \cdot 2^{-b}$$

←

25.

If $M = \frac{-1}{2}$, ϵ will be

$$0 \geq \epsilon > -2 \cdot 2^{-b} \quad \rightarrow 16$$

$$0 \geq \epsilon \times \frac{-1}{2} > -2^{-b}$$

$$0 \geq \epsilon \times \frac{1}{2} > -2 \cdot 2^{-b}$$

In 1's complement ^{& sign-magnitude} representation, the error for truncation of positive values of mantissa is

$$0 \geq M_T - M > -2^{-b} \text{ (or)}$$

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

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

when $M = \frac{1}{2}$, $0 \geq \epsilon \cdot M > -2^{-b}$

$$\text{with } \epsilon = \frac{e}{x}$$

$$\Rightarrow e = \epsilon x = \epsilon \cdot M \cdot 2^{\exp}$$

we get maximum range of relative error for positive M

as

$$0 \geq \epsilon > -2 \cdot 2^{-b} \quad \rightarrow 17$$

For negative mantissa, error value is

$$0 \leq M_T - M < 2^{-b} \text{ (or)}$$

$$0 \leq e < 2^{\exp} \cdot 2^{-b} \quad \rightarrow 18$$

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

$$0 \geq \epsilon > -2 \cdot 2^{-b} \quad \rightarrow 19$$

Range of errors in truncation of floating point nos:

Type of representation for
mantissa

• 2's complement positive
mantissa numbers

Range of error

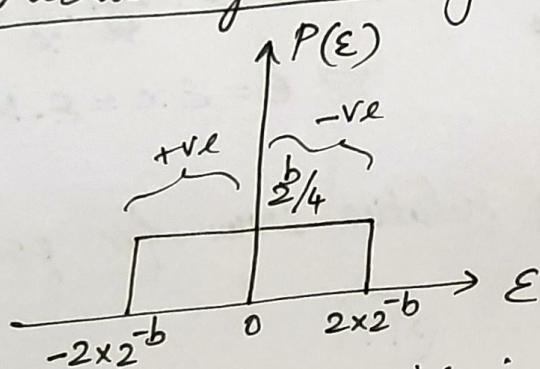
$$0 \geq \epsilon > -2 \cdot 2^{-b}$$

$\cdot 2^b$'s complement -ve mantissa } $0 \leq \varepsilon < 2 \times 2^{-b}$

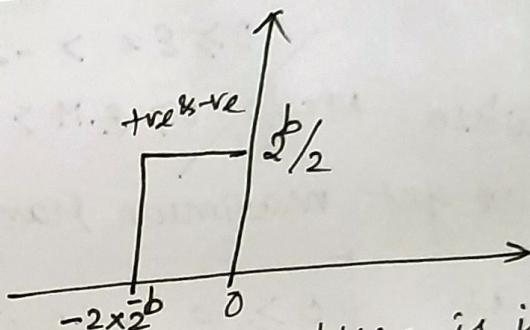
1^b 's complement +ve & -ve mantissa } $0 \leq \varepsilon < -2 \times 2^{-b}$

sign-magnitude +ve & -ve mantissa } $0 \leq \varepsilon < -2 \times 2^{-b}$

Probability density functions $P(\varepsilon)$:



when mantissa is in
2's complement form



when mantissa is in
1's complement or sign-magnitude form.

Error due to rounding of fixed point numbers:

In fixed point arithmetic, error due to rounding a number to ' b '-bits produces an error $\dot{x} = x_R - x$ which satisfies inequality,

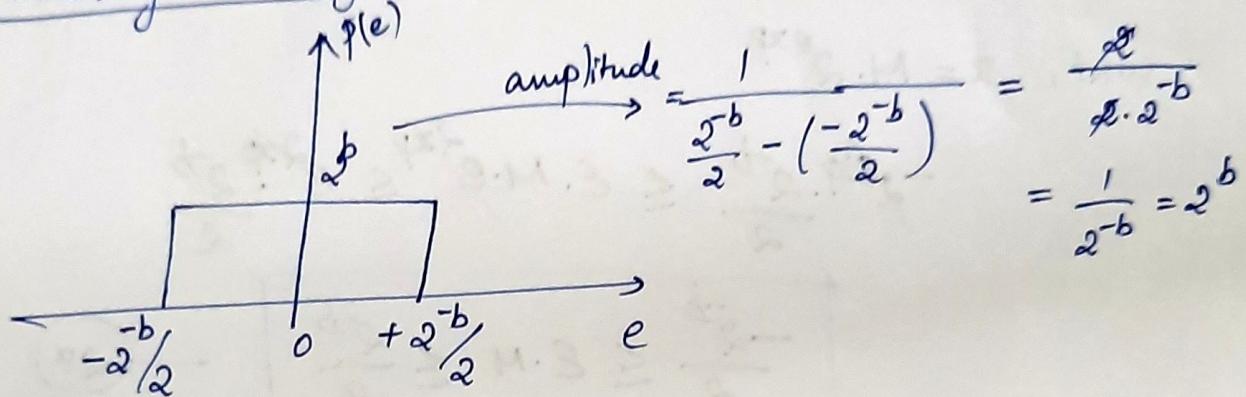
$$\left[-\frac{2^{-b}}{2} \leq x_R - x \leq \frac{2^{-b}}{2} \right] \rightarrow (20)$$

This is because with rounding, the error lies ~~at~~ half way between 2 levels, i.e. it can be approximated

27.

to either nearest higher level or by the nearest lower level. For fixed point number, the above inequality condition is satisfied for all the 3 types of representation. (Sign-mag, 1C₂⁸)

Probability density function P(e):



Error due to rounding of floating point nos:

In floating point arithmetic, only mantissa is affected by rounding. (ie. quantization)

$$\text{If } x = M \cdot 2^{\text{exp}} \text{ & } x_R = M_R \cdot 2^{\text{exp}}$$

$$\text{then } e = x_R - x = (M_R - M) \cdot 2^{\text{exp}} \rightarrow ①$$

But for rounding,

$$\frac{-2^{-b}}{2} \leq (M_R - M) \leq \frac{2^{-b}}{2} \rightarrow ②$$

From ①, $(M_R - M) = \frac{e}{2^{\text{exp}}} = \frac{x_R - x}{2^{\cancel{\text{exp}}} \cancel{2^{\text{exp}}}} \rightarrow ③$

$$-\frac{2^{\exp} \cdot 2^{-b}}{2} \leq (x_R - x) \leq 2^{\exp} \cdot \frac{2^{-b}}{2} \rightarrow 24$$

let relative error $\epsilon = \frac{x_R - x}{x} = \frac{e}{x} \rightarrow 25$

$$e = \epsilon x$$

$$-\frac{2^{\exp} \cdot 2^{-b}}{2} \leq \epsilon x \leq \frac{2^{\exp} \cdot 2^{-b}}{2} \rightarrow 26$$

wkt, $x = M \cdot 2^{\exp}$

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

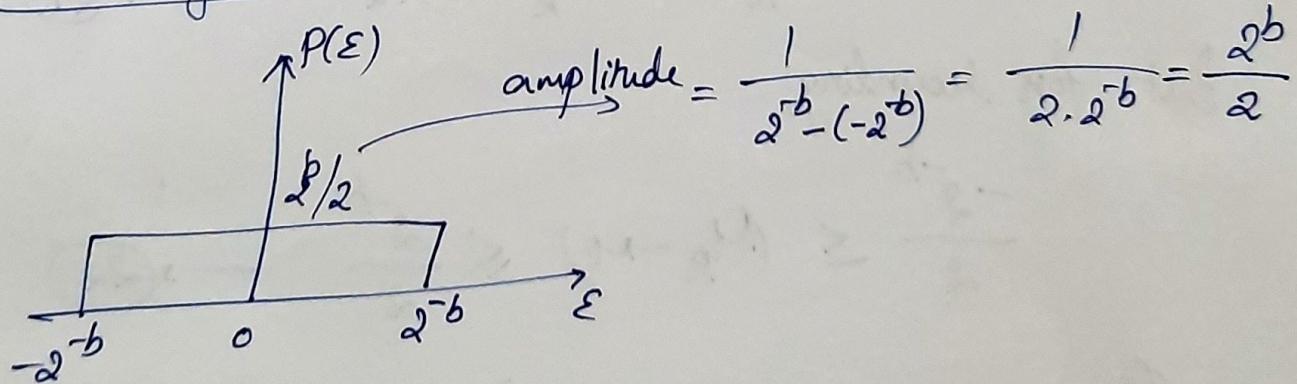
$$\boxed{-\frac{2^{-b}}{2} \leq \epsilon \cdot M \leq \frac{2^{-b}}{2}} \rightarrow 27$$

^{wkt} The mantissa range $M \Rightarrow \frac{1}{2} \leq M < 1$

If $M = \frac{1}{2}$, maximum range of relative error is obtained

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

Probability density function $P(\epsilon)$:



for all 3 types of representation