

show
fig 6.39
6.40
6.41
0 & 5.

Binary Representation of #s.

Many Formats for representation of binary #s.

- one's complement

- sign & mag.

- two's

complement.

most commonly used.

- Real # in two's complement with ∞ precision.

$$x = X_m \left(-b_0 + \sum_{i=1}^{\infty} b_i 2^{-i} \right)$$

$X_m =$ arbitrary scale factor $|x| < X_m$

$b_i =$ either zero or 1

$$b_0 = \text{sign bit} \rightarrow \begin{cases} b_0 = 0 & 0 \leq x \leq X_m \\ b_0 = 1 & -X_m \leq x < 0 \end{cases}$$

with finite # of bits $(B+1)$ we get

representation:

$$\hat{x}_B = Q_B[x] = X_m \left(-b_0 + \sum_{i=1}^B b_i 2^{-i} \right)$$

\hat{x} is quantized version of x .

Smallest difference between any 2 \hat{x} 's in quantized domain $\Delta = X_m 2^{-B}$

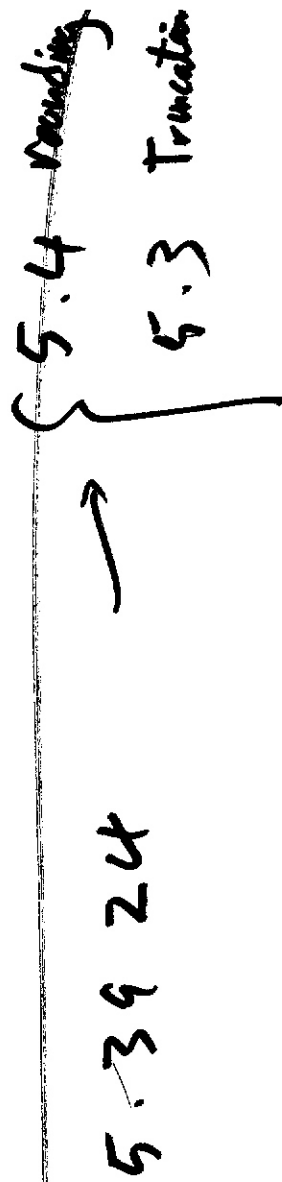
5.87923, only use 2 bits 5.87924

- Quantized #s are in the range.

$$-X_M \leq \hat{x} \leq X_M$$

$$\hat{x}_B = b_0 \square \uparrow b_1 b_2 b_3 \dots b_B$$

Binary
Point.



Start a real number x
to get \hat{x}_B , one can either
rounding or Truncation:

Show Fig 6.37 (a), 6.37(b) in OXS.

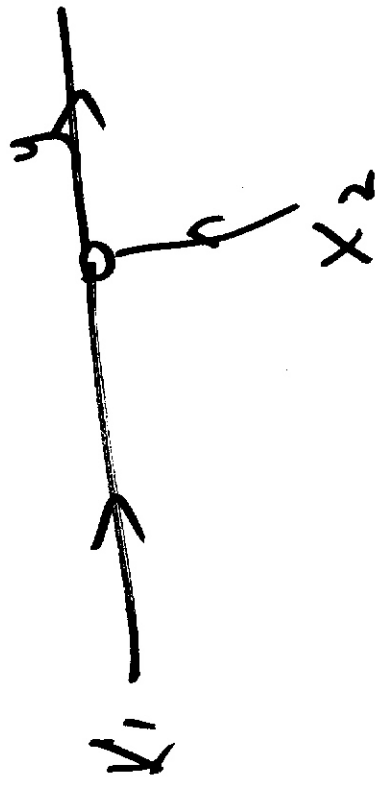
Quantization error

$$e = Q_B[\hat{x}] - x \\ = \hat{x}_B - x$$

2's complement : Rounding error $-\frac{\Delta}{2} < e < +\frac{\Delta}{2}$
Truncation error $0 \leq e \leq \Delta$

$$x_1 + x_2 < 10000$$

$$y = x_1 + x_2$$



$$x_1 < 5000$$

$$x_2 < 5000$$

Overflow.

Natural overflow

b.38
085

Saturated

- Interesting property of two's complement.
+ natural overflow:

Add few #s, if the final sum doesn't overflow, then result is correct even though the intermediate results overflow. The intermediate overflow & rounding error.

- Tradeoff between overflow
 $X_M \uparrow \rightarrow$ overflow is less likely, $e \uparrow$
but $\Delta \uparrow$,

$X_n \downarrow$ \rightarrow overflow is more likely.
but $A \downarrow$ $e \downarrow$

BUT B can cause overflow.

- keep k_m large to minimize chances of overflow

- But keep B large to keep

D, e small.

Multiplication also introduces $\left\{ \begin{array}{l} \text{- overflow} \\ \text{- rounding error} \end{array} \right.$