

Quiz 1 (January 30, 2019 (Wednesday)) -

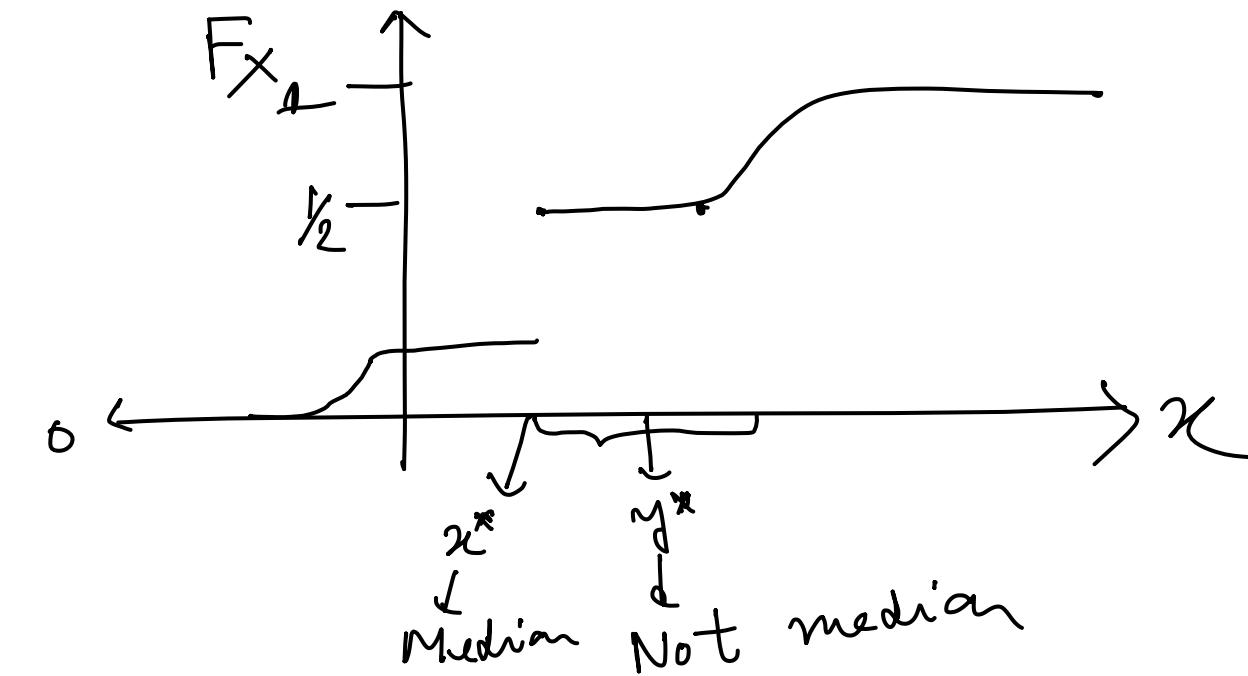
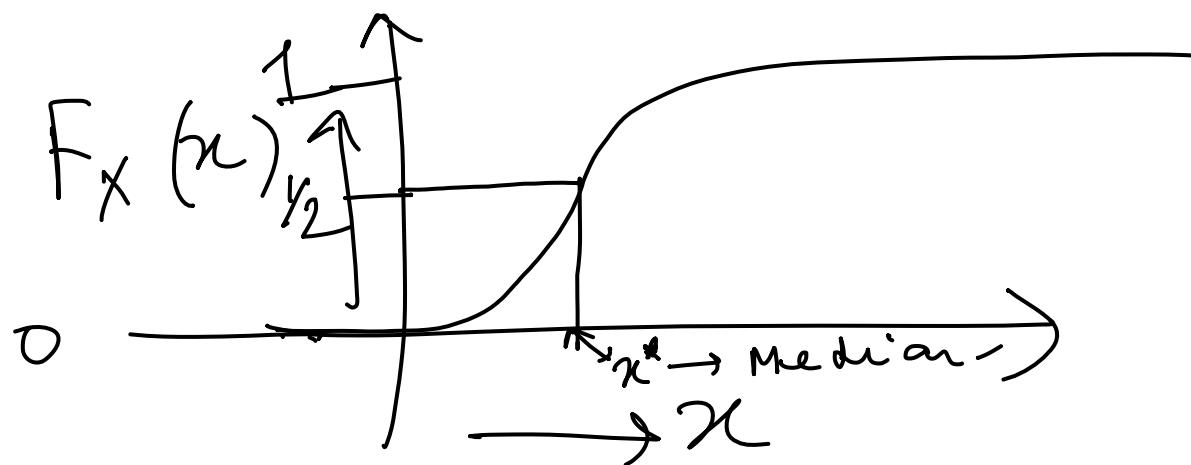
Time :- 6.35 p.m. to 7.15 p.m. Venue : L18, L19, L20.

$X$  : Random variable .

$F_x$  : CDF of  $X$  .

$\text{Med } F = \inf_x \left\{ x : F(x) \geq \frac{1}{2} \right\}$  → General def<sup>n</sup>.

If  $F$  is cont<sup>n</sup>, we then have  $\text{Med } F = F^{-1}\left(\frac{1}{2}\right)$  .



Advantage of using median! -

$$X = \{x_1, \dots, x_n\} \rightarrow \text{data}$$

Want to compute median ( $x$ ) .

Middle most observation .

Suppose  $X^* = \{x_{(1)}, \dots, x_{(n)}\}$

The smallest  $\downarrow$   $\rightarrow$  Ordered data

$x_{(i)} \rightarrow i\text{-th smallest}, i=1, \dots, n$ .

$$\text{Med}(x) = x_{\left(\frac{n+1}{2}\right)} \quad \text{if } n \text{ is odd} .$$

$$= \frac{x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2} + 1\right)}}{2} \quad \text{if } n \text{ is even} .$$

$$A = \{3, 1, 2, 4\}, A^* = \{1, 2, 3, 4\} \rightarrow \text{ordered}$$

$$\text{Med}(A) = \frac{2+3}{2} = 2.5.$$

Suppose, there is one contamination

$$A_{\text{con}} = \{3, 2, 400, 1\}, A^*_{\text{con}} = \{1, 2, 3, 400\}.$$

$$\text{Med}(A^*_{\text{con}}) = \frac{2+3}{2} = 2.5.$$

In other words, this numerical example indicates that median does not get contaminated by single obs.

However,

$$\text{Mean}(A) = \frac{1+2+3+4}{4} = 2.5$$

$$\text{but } \text{Mean}(A_{\text{con}}) = \frac{1+2+3+400}{4} = 101.5.$$

That mean is "highly" affected by a single contamination.

Comment:- Overall, median is more "Robust" against the influential observations.

Quantile:-

$X$ : Random variable

$F$ : CDF of  $X$

$$\text{Quantile}_h F(x) = \inf_{x} \left\{ x : F(x) \geq \alpha \right\}, \quad \alpha \in (0, 1).$$

If  $\alpha = \frac{1}{2}$ , Then Quantile = Median.

Some well-known distributions (discrete & cont<sup>n</sup>):-

Discrete Distributions:-

## 1. Discrete uniform :-

Suppose  $X = \{x_1, \dots, x_n\}$ .

$X$ : Random variable

P.M.F:-  $P[X = x_i] = \frac{1}{n} \quad \forall i = 1, \dots, n.$

$P[\cdot]$  is a proper p.m.f. (Try to convince yourself).

$$E(X) = \sum_{i=1}^n x_i \times \left(\frac{1}{n}\right) = \frac{1}{n} \sum_{i=1}^n x_i \text{ (average).}$$

$$E(X^K) = \sum_{i=1}^n x_i^K \times \left(\frac{1}{n}\right) = \frac{1}{n} \sum_{i=1}^n x_i^K$$

$$M_X(t) = \sum_{i=1}^n e^{tx_i} \times \left(\frac{1}{n}\right) = \frac{1}{n} \sum_{i=1}^n e^{tx_i}$$

$$\text{Var}(x) = E(x^2) - \{E(x)\}^2$$

(check  
 $\text{Var}(x) = E[x - E(x)]^2$ )

$$= \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2$$

Check that

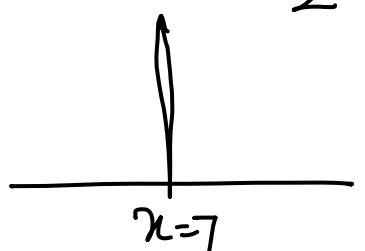
$$\frac{1}{n} \sum_{i=1}^n (x_i - \frac{1}{n} \sum_{i=1}^n x_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2$$

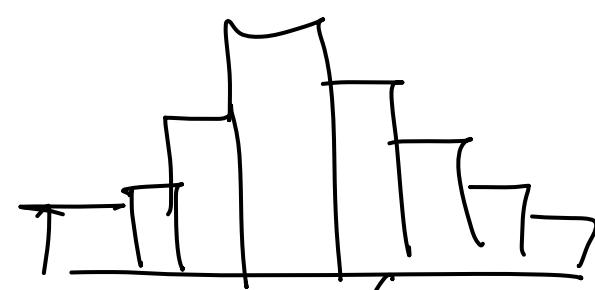
Remark:-

Under Discrete uniform distribution,  
 we are having the sample version of different measure  
 in descriptive statistics.

2. Degenerative dist<sup>n</sup>. :-



$x$ : Random variable.  
 $P[x = c] = 1$  and  $P[x = x] = 0 \quad \forall x \neq c$ .



Histogram of  
 the data  
 obtained  
 from  $N(a)$

$$E(x) = cx_1 = c$$

$$E(x^2) = c^2 x_1 = c^2$$

$$M_x(t) = e^{tc} x_1 = e^{tc},$$

$$\text{Var}(x) = E(x^2) - \{E(x)\}^2 = c^2 - c^2 = 0$$

$$F_x(x) = P[x \leq x] = \begin{cases} 0 & \text{if } x < c \\ 1 & \text{if } x \geq c \end{cases}$$

3. Two points distribution (e.g. Bernoulli dist<sup>n</sup>)  
mass points

$$x : \quad x_1 \quad x_2$$

$$P[x=x] : p \quad p_2 = (1-p)$$

$x_1$  and  $x_2$

$$E(X) = x_1 \times p + x_2 \times (1-p) \quad \text{---} \quad ①$$

$$E(X^2) = x_1^2 \times p + x_2^2 \times (1-p) \quad \text{---} \quad ②$$

$$\text{Var}(X) = E(X^2) - \{E(X)\}^2 = ② - ①^2 \text{ chuck} = p(1-p)(x_1 - x_2)^2$$

$$M_X(t) = E(e^{tx}) = e^{tx_1} \times p + e^{tx_2} \times (1-p).$$

Remark:-  $E(X) = \frac{d}{dt} M_X(t) \Big|_{t=0} = x_1 p + x_2 (1-p)$

↓  
This is another way of deriving  $E(X)$

or any other moments

Note:- If  $x_1 = 0$  &  $x_2 = 1$ , it is called Bernoulli dist<sup>n</sup>-  $P[X=x] = p^x (1-p)^{1-x}$ ,  $x=0 \text{ or } 1$ ,  $p \in (0,1)$ .

## Binomial dist<sup>n</sup>.

Formulation:-

$n$  indep.

Bernoulli trials with success prob  
 $= p$ .

$X$ : No. of success in  $n$  trials.

$\mathcal{X}$ : Range of  $X$

$$\mathcal{X} = \{0, 1, \dots, n\}.$$

$p(x) = P[X=x] = \begin{cases} 0 & \text{if } x \notin \mathcal{X} \\ \binom{n}{x} p^x (1-p)^{n-x} & \text{if } x \in \mathcal{X}. \end{cases}$

$$X \sim \text{Bin}(n, p)$$

$$\mathcal{X} = \{0, \dots, n\}.$$

$$P[X=x] = \begin{cases} 0 & \text{if } x \notin \mathcal{X} \\ \binom{n}{x} p^x (1-p)^{n-x} & \text{if } x \in \mathcal{X} \end{cases}$$

Note  $P[X=x] \geq 0 \quad \forall x = 0, \dots, n$

And,  $\sum_{x=0}^n \binom{n}{x} p^x (1-p)^{n-x} = (1+p)^n = 1$

Then,  $P[X=x]$  is a f.m.f.