

Variables  $\rightarrow$  like Ages = {22, 25, 27, 30}

Variable  $\rightarrow$  like single Age = 22

## Variables

A variable is a property that can take on any value.  
Eg: Height, Weight

Types of variables:-  
Numerical.

Quantitative:- Measured numerically & perform arithmetic operations (Add, sub, multi)  
Eg:- Weight, height, Age.

Qualitative / Categorical var.:- Based on some characteristics <sup>not out grouped together</sup> we can derive categorical variable.  
Eg: Gender  $\begin{cases} M \\ F \end{cases}$ , IQ, Blood group, T-shirt (S, L, XL, XXL)

## Quantitative

Discrete Var.  $\rightarrow$  (whole number)

- 1) no. of children
- 2) no. of Bank A/c.
- 3) no. of vehicles

or countable entities

Continuous Var. <sup>(any value it can have)</sup>

Eg:- Height, 172.5,  
Weight, 95.5 kg

Example: 1.35 km

# → Central Tendency / Measure of Central Tendency

Refers to the measure, used to determine the center of the distribution of data.

Eg: mean, median, mode, Positional values  
 (divides into equal parts like median, quartile, decile, percentile)  
 Often called Mathematical average.  
Positional avg.

Ex: median is the positional no. in the short end data.  
 Positional value ex: data is dividing into 2 parts by median.



■ JUNE ■ 2018

with mean in feature engineering.  
 Trace NaN value.

156-199 / 100-24

15

Friday

## → Arithmetic Mean for Population & Sample

① Mean (Average) :-  $\mu$  (N)

a) 
$$\mu = \sum_{i=1}^N \frac{x_i}{N} \rightarrow \text{[in case of Population (N)]}$$

b) 
$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n} \rightarrow \text{[in case of sample (n)]}$$

a) {1, 1, 2, 2, 3, 3, 4, 5, 5, 6}

$$\mu = \sum_{i=1}^n \frac{x_i}{N} \Rightarrow \frac{1+1+2+2+3+3+4+5+5+6}{10} \Rightarrow \frac{32}{10} = 3.2$$

b)  $\bar{x} = \frac{32}{10} = 3.2$

•  $\mu$  can be  $\geq \bar{x}$

•  $\mu$  can never  $< \bar{x}$

c) <sup>with</sup> Outliers (outliers really have adverse impact on entire distribution)

Outlier is the no. which is "completely different from entire distribution"

{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100} → outlier

$$\mu = \frac{32+100}{11} \Rightarrow 12$$

$\mu = 3.2$  } huge diff. with resp. to mean due to outlier.  
 new  $\mu = 12$



09:00

10:00

11:00

12:00

13:00

14:00

15:00

16:00

17:00

18:00

19:00

20:00

17:16  
SCHEDULE

## ② Outliers in case of Median:-

a)  $\{1, 1, 2, 2, 3, \boxed{3}, 4, 5, 5, 6, \underline{100}\}$  with 1 outlier  
first sort the no.'s.

$$\text{Median} = \underline{3}$$

b)  $\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, \underline{100}, \underline{112}\}$  2 outliers

$$\text{Median} = \frac{\cancel{3} + 4}{2} \Rightarrow \underline{3.5}$$

• In case of mean, while adding outlier  
 $M = 3.2$   
 $M_1 = 12$  (with outlier)

• In case of median, while adding outlier  
 $\text{Med} = 3$   
 $\text{Med}_1 = 3.5$  (with outlier)

\* It implies, Median works well with outlier

17 Sunday

## ③ Mode with Outliers.

(we can use mode both in integer or categorical way but it works well with categorical variables.)

$\{1, 2, 2, 3, 4, 5, 6, 6, 6, 7, 8, \underline{100}\}$

In case of

$\{1, 2, 2, \dots, 7, 8, 100, 100, 100, 100\}$

we take 100 outlier as mode.



# → Measure of Dispersion

(Dispersion)  
↓  
(Spread)

- ① variance
- ② standard deviation.
- ③ Relative dispersion/coefficient of variation.
- ④ JQR

① variance :- concept of measure of dispersion.

Q. 1) Population variance (N)

$$\sigma^2 = \sum_{i=1}^N \frac{(x_i - \mu)^2}{N}$$

$$\Rightarrow \frac{10.84}{6} = 1.81$$

eg	x	x - μ	(x - μ) <sup>2</sup>
	1	-1.83	3.34
	2	-0.83	0.688
	2	0.83	0.688
	3	0.17	0.03
	4	1.17	1.37
	5	2.17	4.71
	μ = 2.83		Σ = 10.84

for grouped data  
(C-I)

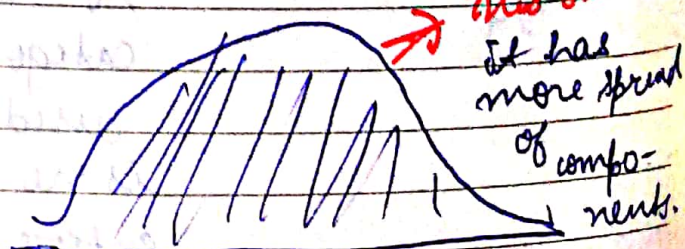
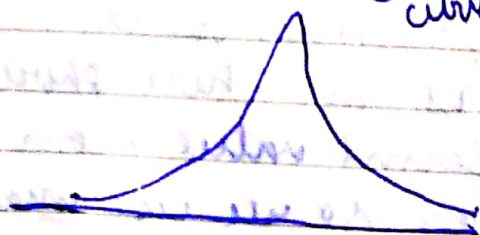
$$\sigma^2 = \sum_{i=1}^N \frac{f(x_i - \mu)^2}{N}$$

f = frequency.

\* "The more variance means the data is more dispersed."

ex which graph has more variance?

(Bell curve)





SCHEDULE

## ② Standard deviation ( $\sigma$ )

$$\sigma = \sqrt{\text{variance}}$$

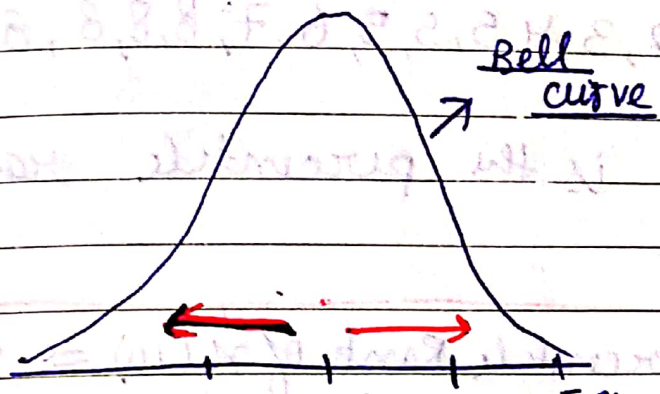
OR

$$\sigma = \sqrt{\sigma^2}$$

$$\sigma = \sqrt{1.81} = 1.345 \text{ (s.d.)}$$

$$\begin{array}{r} 2.87 \\ + 1.34 \\ \hline 4.17 \end{array}$$

$$\begin{array}{r} 2.83 \\ - 1.34 \\ \hline 1.49 \end{array}$$



(1 s.d. to the left)

1.49    2.83 ( $\mu$ )    4.17    5.51 (1 s.d. to the right)

[specifying measure of cent. tendency, center is here for specific distribution.]

\* Basically variance is rectifying how much spread present & S.D tells the particular element is how far from the mean.

## b) Sample Variance ( $n$ )

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{[n-1]}$$

for grouped data:-

$$s^2 = \frac{\sum f(x_i - \bar{x})^2}{n-1}$$

$f = \text{frequency}$

for Population

③

$$C.V = \frac{\sigma}{\mu} \times 100$$

for sample  $\Rightarrow$



→ Percentile & Quartile } how to find outliers  
(helps to remove outliers)

A Percentile is a value below which certain percentage of observation lie.

• Ex- 99 percentile: - A student got better marks than 99% of entire class

Reg <sup>sm</sup>  
Data: 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

→ Q1 what is the percentile ranking of 10?

Sol  
 $n=20$

$$\text{Percentile Rank of } x(10) = \frac{\text{no. of values below } x}{n}$$

$$= \frac{8}{20} \times 100$$

→ 80 percentile

\* means 80% of entire distribution is less than 10

→ Q2 what value exist at percentile ranking of 25%? (same data:  $n=20$ )

Sol

$$\text{Value} = \frac{\text{Percentile} \times n+1}{100}$$

$$\text{Value} = \frac{25}{100} \times 21 = 5.24 \text{ (index)} \\ \text{value at 5th index}$$



→ Five Number Summary is used to determine the outlier

- ① minimum
- ② First Quartile ( $Q_1$ )
- ③ Median
- ④ Third Quartile ( $Q_3$ )
- ⑤ Maximum

\* Removing the outlier (i.e.) using Five no. summary

{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27 }

→ when we need to remove outlier, we ~~have~~ need to define a lower fence & higher fence. Value ~~in~~ present should b/w lower fence & higher fence.

• Lower fence - After a ~~to~~ smaller no. all the no. below the smaller no. will be lower fence

• Higher fence - After a higher no. all the no. above the higher no. will be higher fence.

$$\text{Lower fence} = Q_1 - 1.5(IQR)$$

$IQR = \text{Inter Quartile Range.}$

$$\text{Upper fence} = Q_3 + 1.5(IQR)$$

$$IQR = Q_3 - Q_1$$

$$Q_1 = 25\%$$

$$Q_3 = 75\%$$



23

174-101/Wk-25

2018 ■ JUNE ■

Saturday

data: {1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27}

09.00

 $Q_1$ what is 25% percentile ( $Q_1$ ) & 75% ( $Q_3$ )

10.00

 $Q_1$ 

11.00

$$\bullet \text{ Value} = \frac{25}{100} \times (19+1) \Rightarrow \frac{205}{100}$$

12.00

13.00

$$\Rightarrow \boxed{3} = (Q_1)$$

14.00

15.00

16.00

$$\bullet \text{ Value} = \frac{375}{100} \times 205 \Rightarrow 15^{\text{th}}$$

17.00

18.00

$$\Rightarrow \boxed{7} = (Q_2)$$

19.00

20.00

$$\bullet \text{ IQR} \Rightarrow Q_2 - Q_1$$

$$\Rightarrow 7 - 3$$

$$\Rightarrow 4$$

$$\bullet \text{ Lower fence} = Q_1 - 1.5(\text{IQR})$$

$$\Rightarrow 3 - 1.5(4)$$

$$\Rightarrow 3 - 6 \Rightarrow \boxed{-3} = \text{lower fence}$$

$$\bullet \text{ Higher fence} \Rightarrow Q_3 + 1.5(\text{IQR})$$

$$\Rightarrow 7 + 1.5(4)$$

$$\Rightarrow \boxed{13} = \text{high fence}$$

24 Sunday

So, our lower to higher fence range is:-

$$[\text{lower} \longleftrightarrow \text{higher fence}] \text{ range}$$

$$\boxed{-3 \longleftrightarrow 13}$$

\* here, every no.  $> 13$  is outlier & no  $< -3$  is also considered as outlier.



Here 27 should be removed as it is greater than 13, belongs to higher fence.

remaining data

{ 1, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, ~~27~~ }

minimum value = 1

$Q_1 = 3$

$Q_3 = 7$

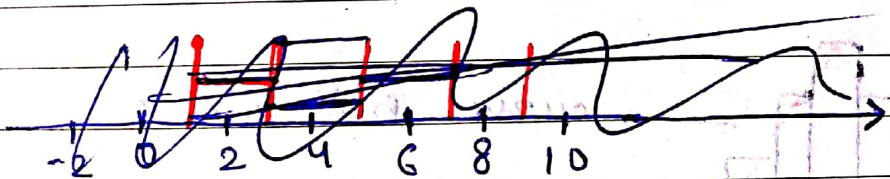
maximum = 9

median = 5 ( $\frac{5+5}{2}$ )

(5 no summary)

→ making box plot from the resultant of 5 no. sum

\* "removing an outlier is done by lower fence & higher fence and IQR"



Boxplot basically use to see 'outlier'

