

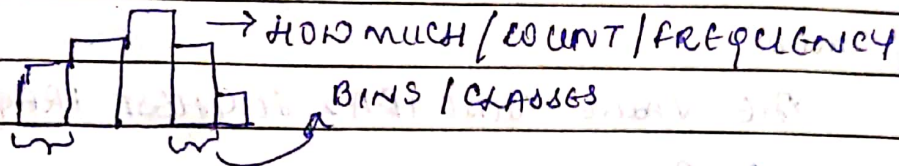
STATS!

RAW DATA → INFORMATION → KNOWLEDGE → WISDOM

↓  
DESCRIPTIVE  
STATISTICS

FREQUENCY DISTRIBUTION:

RAW DATA → CLASSES & FREQUENCY  
↓  
HISTOGRAM



ASSUME THERE ARE RANDOM NUMBERS B/W 8-23.

THESE NUMBERS ARE CLASSIFIED:

8-11	11-14	14-17	17-20	20-23	→ CLASSES
↓	↓	↓	↓	↓	
2	7	12	3	1	→ FREQUENCY

CUMULATIVE DISTRIBUTION =  $2 + 7 + 12 + 3 + 1 = 25$

CUMULATIVE DISTRIBUTION CURVE OR OGIVE CURVE:

CENTRAL TENDENCY:

WHERE IS THE DATA AROUND → AVG (LIKE)

1] ARITHMETIC MEAN:

✓  $\bar{x} = \frac{\sum x}{n}$  → SUM OF NUMBERS  
→ NO. OF VALUES / COUNT / SAMPLE SIZE  
MEAN

BUT EXTREME VALUES AFFECT THE MEAN HENCE  
DO NOT USE IF THE DATA SET HAS EXTREME VALUES.

2] MEDIAN  $\rightarrow$  VALUE @ 50th PERCENTILE  
MIDDLE MOST OBSERVATION AFTER  
ARRANGING THE DATA SET FROM LOWEST TO  
HIGHEST.

THIS TO A LARGE EXTENT REDUCES THE EFFECT OF  
EXTREME VALUE.

GIVEN NUMBER OF VALUES  $\rightarrow$  MIDDLE VALUE  
ODD NUMBER OF VALUES  $\rightarrow$  SAME AVERAGE OF  
MIDDLE VALUE

3] MODE

THE VALUE THAT HAS HIGHEST FREQUENCY.  
LIKELY OUTPUT.

DRAW BACK  $\rightarrow$  MULTIPLE VALUES OF SAME FREQUENCY

HOW IS DATA DISTRIBUTED AROUND CENTRAL  
TENDENCY  $\rightarrow$  MEASURE OF DISPERSION

1] RANGE OF DISPERSION

$$R = X_{\text{MAX}} - X_{\text{MIN}}$$

IF  $X_{\text{MAX}} = X_{\text{MIN}} \rightarrow$  ZERO DISPERSION

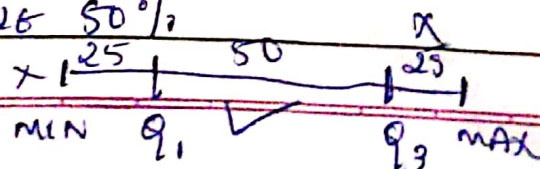
HOW IS THE SPREAD OF DATA.

RANGE SHOULD NOT BE USED IF DATA SET HAS  
EXTREME VALUES.

2] INTER-QUARTILE RANGE (IQR)

SOLUTION TO PROBLEM USING RANGE.

HERE LOWEST 25% & HIGHEST 25% ARE IGNORED  
AND RANGE IS CALCULATED USING ONLY THE  
MIDDLE 50%





## STANDARD DEVIATION:

AVERAGE DEVIATION FROM THE MODE OF THE DATA-

- CALCULATE MEAN  $\bar{x}$
- FROM EACH VALUE (OBSERVATION) SUBTRACT  $\bar{x}$
- SQUARE THE OUTPUT [EACH VALUE]
- FIND THE SUM
- DIVIDE THE SUM BY (NUMBER OF OBSERVATIONS - 1)
- $\hookrightarrow$  VARIANCE
- $\sqrt{\text{VARIANCE}} = \text{STANDARD DEVIATION}$

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

$$\text{STD DEV} = \sqrt{\text{VARIANCE}}$$

$(n-1) \rightarrow$  BECAUSE  $\bar{x}$  IS ESTIMATE

## COEFFICIENT OF VARIATION

$\hookrightarrow$  RATIO OF S.D TO MEAN

$$CV = \frac{s}{\bar{x}}$$

50 UNITS WITH SD = 5  $\rightarrow$  SALES MAN 1 } BETTER??  
 75 UNITS WITH SD = 25  $\rightarrow$  SALES MAN 2 }

$$\begin{aligned} CV_1 &= \frac{5}{50} & CV_2 &= \frac{25}{75} \\ &= \frac{1}{10} & &= \frac{1}{3} \\ &= 10\% & &= 33\% \end{aligned}$$

## THE EMPIRICAL RULE

IF  $\mu \pm 1\sigma$  HAS 68% OF DATA  
 $\mu \pm 2\sigma$  HAS 95% OF DATA  
 $\mu \pm 3\sigma$  HAS 99.7% OF DATA

SP1

AVG 50 UNITS

SD 5 UNITS

(40, 60)

$\mu \pm 2\sigma$

SP1 SELL AT LEAST

40 UNITS

SP2

AVG 75 UNITS

SD 25 UNITS

(25, 125)

$\mu \pm 2\sigma$

SP2 SELL AT LEAST

25 UNITS

## CHEBYSHEV RULE

IF DATA NOT BELL SHAPED OR BELL SHAPED

AT LEAST  $\left(1 - \frac{1}{K^2}\right) \%$  DATA FALL WITHIN

K STANDARD DEVIATIONS.

IF  $K = 2$

75% OF DATA WILL FALL BETWEEN 2 SD OF MEAN

## FIVE NUMBER SUMMARY:

$X_{\text{smallest}}$

FIRST QUARTILE ( $Q_1$ )

MEDIAN ( $Q_2$ )

THIRD QUARTILE ( $Q_3$ )

$X_{\text{largest}}$

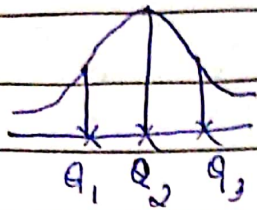
DIFF b/w EACH

IS MEASURE OF

DISPERSION

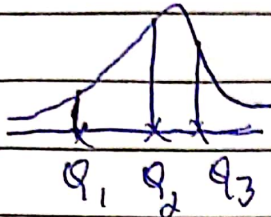


## DISTRIBUTION SHAPE



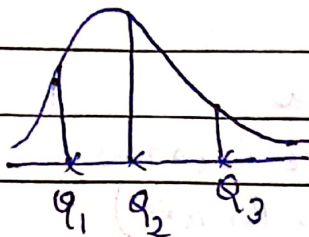
$Q_2 \rightarrow$  middle of  $Q_1$  &  $Q_3$

~~DETERMINED~~  
SYMMETRIC



$Q_2$  &  $Q_3$  closer  
 $Q_1$  larger observations

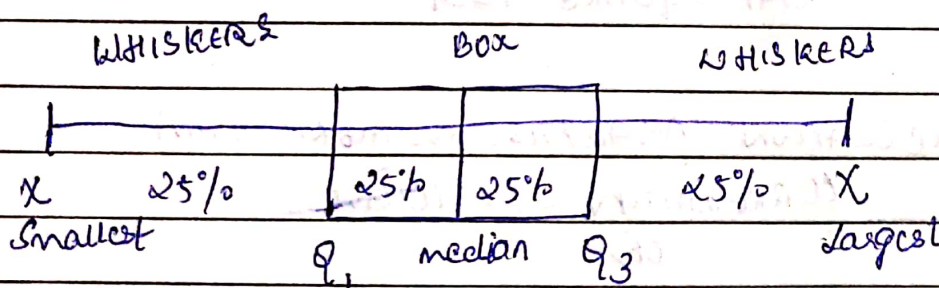
~~DETERMINED~~  
LEFT SKEW



$Q_1$  &  $Q_2$  closer  
 $Q_3$  larger observations

RIGHT SKEW

## BOX PLOT: GRAPHICAL REP. OF FIVE NUMBER SUMMARY



Box  $\rightarrow Q_1$  to  $Q_3$

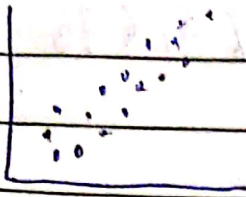
Box Plot  $\rightarrow$  HORIZONTAL OR VERTICAL

$\rightarrow$  SHOWS THE SKEW

$\rightarrow$  GOES TILL  $1.5 \times IQR$

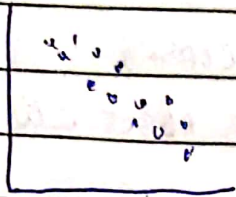
CORRELATION: VARIABLE 1  
V1

VARIABLE 2  
V2



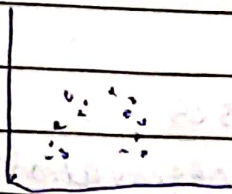
+VE CORRED

$V1 \propto V2$



-VE CORRED

$V1 \propto \frac{1}{V2}$



NO CORRED

(MULTIPLE RELATION  
OR NO PATTERN)

CORRELATION ANALYSIS: NOMINAL DATA  
CHI SQUARE TEST

CORRELATION ANALYSIS: NUMERIC DATA  
CORRELATION COEFFICIENT  
OR

PEARSON'S COEFFICIENT.

SUMMARY:

HISTOGRAM

CENTRAL TENDENCY

DISPERSION

CORRELATION ANALYSIS.