

Descriptive Statistics

- Statistics is the science of learning from Data
- Data is essentially numbers (text/Symbols) which represent some information.
- It helps to think of data as value of quantitative and qualitative variables.
- What are the variables types:
 - * Numerical or Quantitative (Continuous & Discrete)
 - * Categorical or Qualitative (always Discrete)
 - * Nominal: It is ^{not} possible to arrange data in any order.
 - * Ordinal: It is possible to arrange data in order.

represent characteristics like Gender (Male / Female / Transgender)

Example of Nominal data } Marital Status — (Married / Unmarried)
 } Home State — (States of India)

Example of Ordinal Data { Color of Terror alert possible value
 { Green, Yellow, Orange & Red
 ↓
 Very low Risk Very high Risk.
 ordered Categorical Variable

Player Quality: A class, B class, C class, D class

Very high skill player

Significantly low skill Player

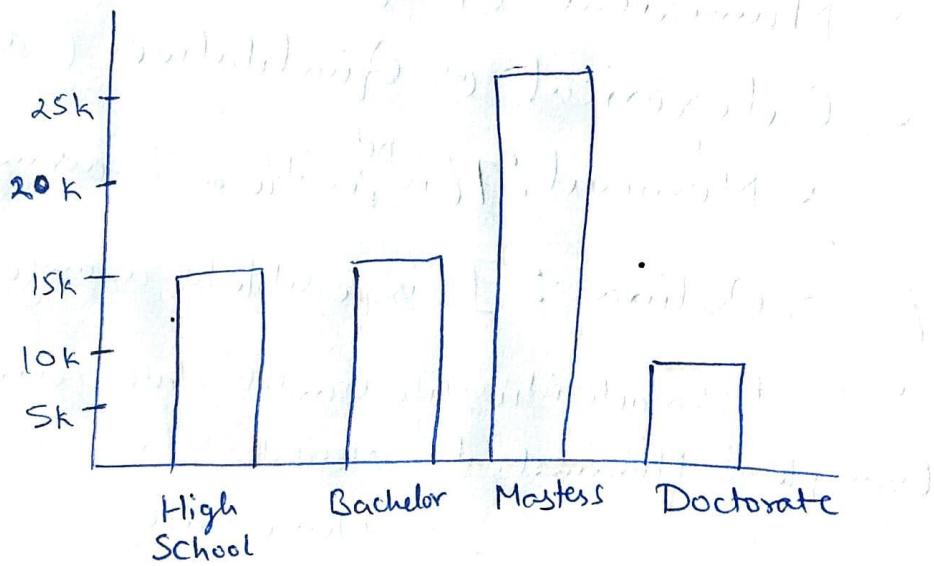
Descriptive Statistics is method to Quantitatively
describing the data

- ↳ Graphical Representation
- ↳ Tabular Representation
- ↳ Summary Statistics

Graphical Representation

Categorical Variable

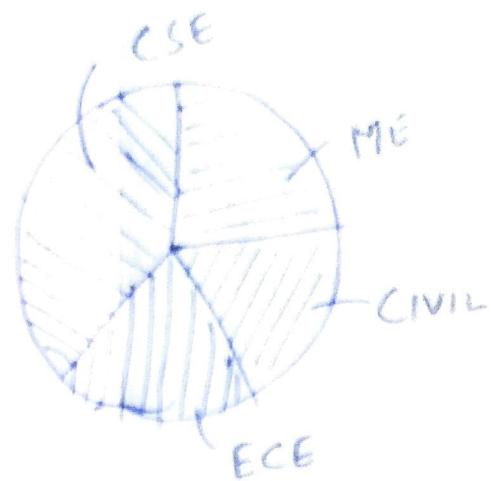
- ↳ BAR CHART



Ordinal Variable (Data)

→ PIE CHART

①
②

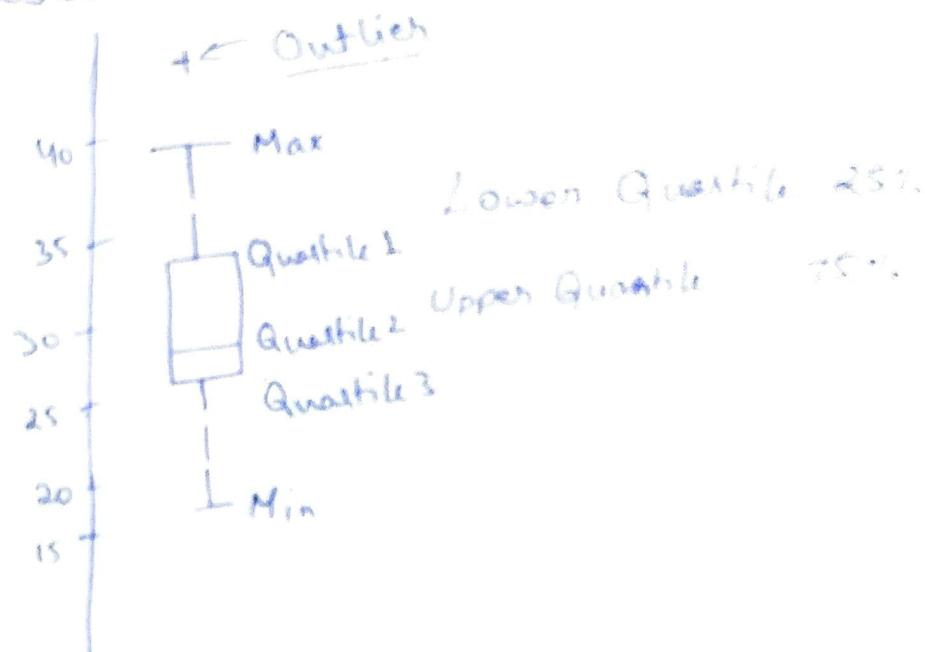


↳ Quantitative Variable

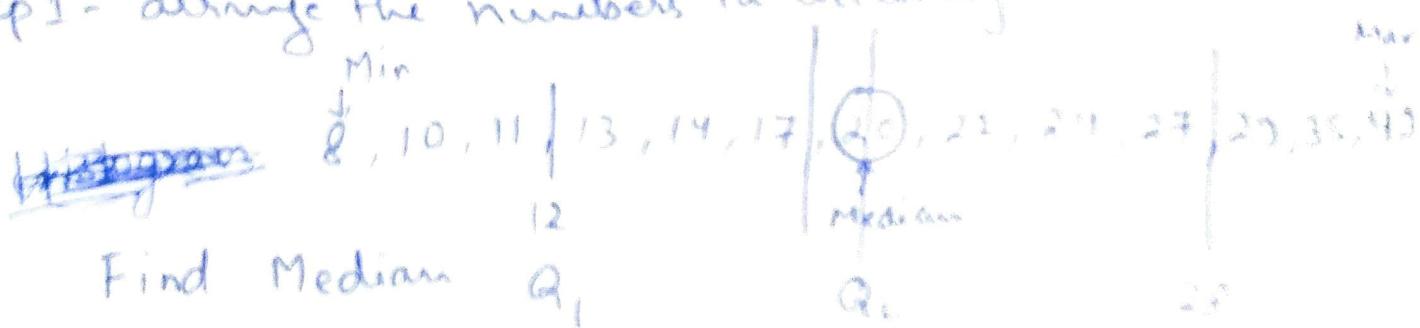
Box & Whisker Plot

Box Plot

Ex 11, 22, 20, 14, 29, 8, 35
27, 13, 43, 10, 24, 17



Step 1 - arrange the numbers in ascending order



Find Median Q_2

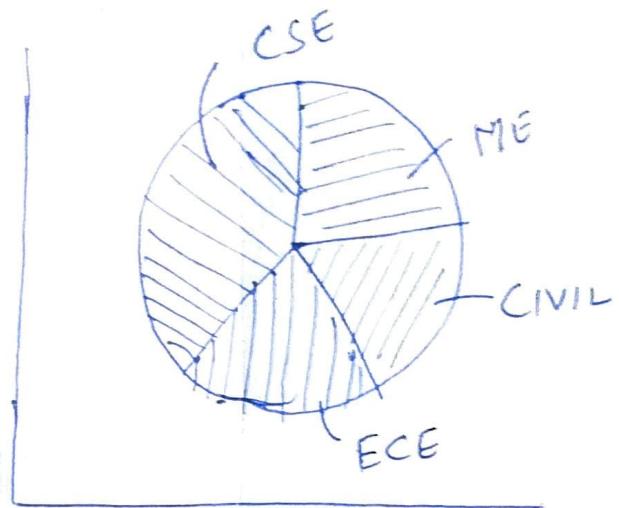
Q_2

Now to check the min & max are not outliers

Now we need to determine the range of numbers in which the outlier could be

$[Q_1 - 1.5 \text{ IQR}, Q_3 + 1.5 \text{ IQR}]$ Any number outside this range is an outlier

→ PIE CHART



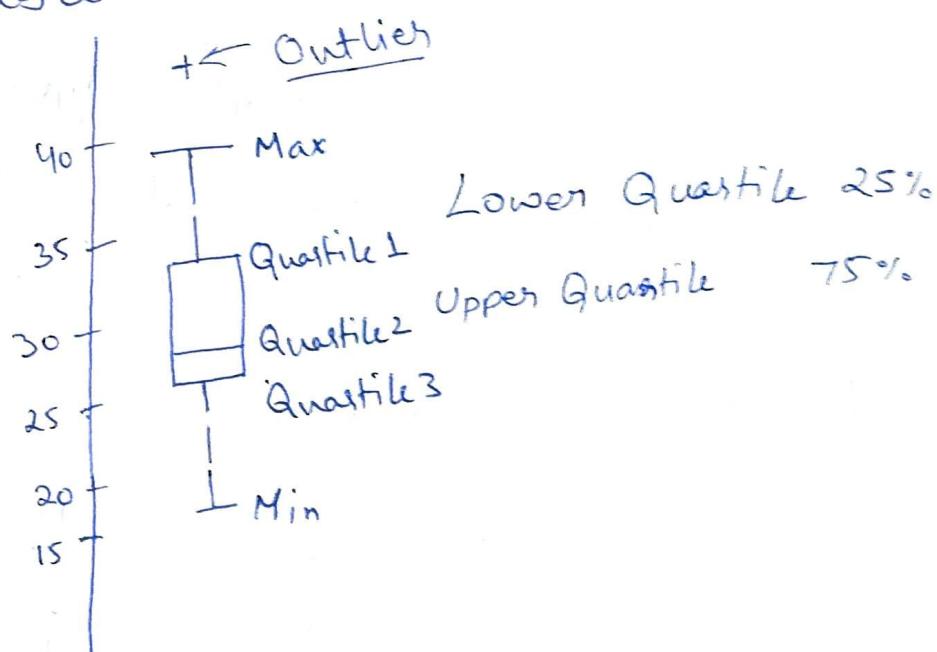
(2) (2)

→ Quantitative Variable

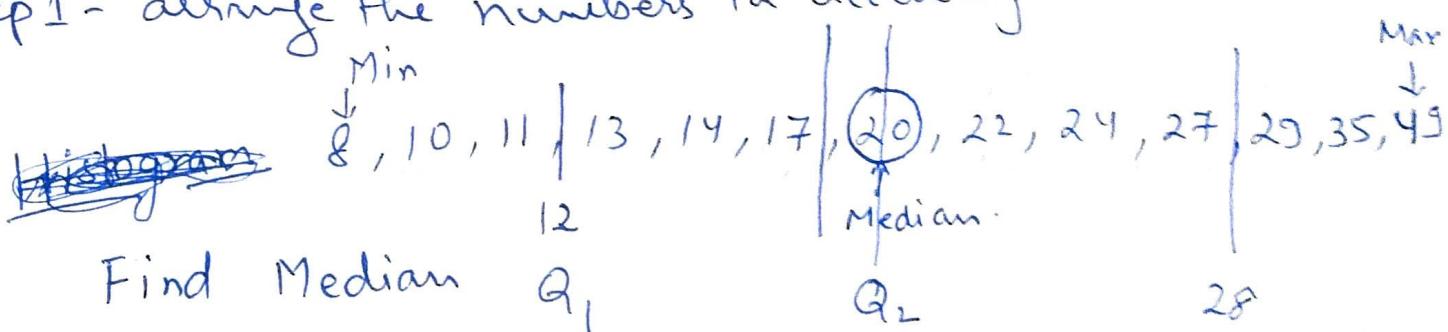
Box & Whisker Plot

Box Plot

Ex 11, 22, 20, 14, 29, 8, 35
27, 13, 49, 10, 24, 17



Step 1 - arrange the numbers in ascending order.



Find Median

Q_1

Q_2

28

Now to check the min & max are not outliers

Now we need to determine the range of numbers in which the outlier could be

$[Q_1 - 1.5 \text{ IQR}, Q_3 + 1.5 \text{ IQR}]$ Any number outside this range is an outlier.

Box & Whisker Plot

2'

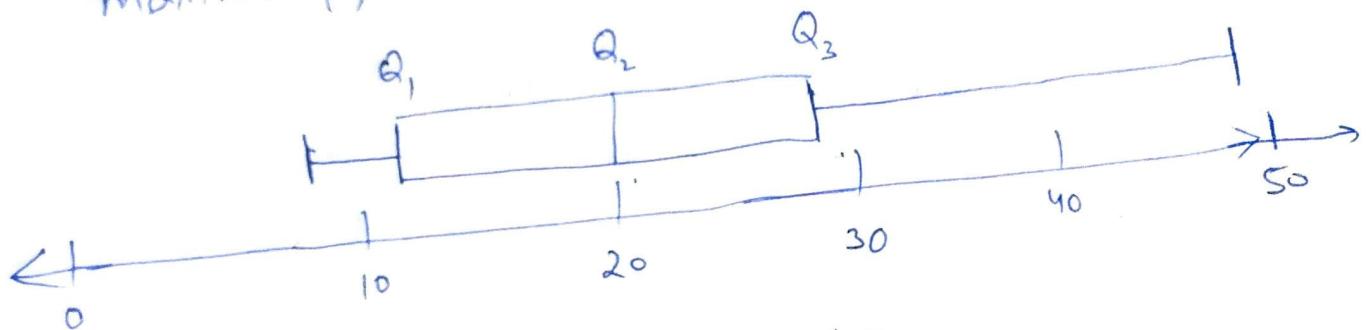
IQR : Inter Quartile Range is the difference between Q_3 & Q_1 .

$$IQR = Q_3 - Q_1 = 28 - 12 = 16$$

$$Q_1 - 1.5 IQR = 12 - 1.5 \times 16 = -12$$

$$Q_3 + 1.5 IQR = 28 + 1.5 \times 16 = 52$$

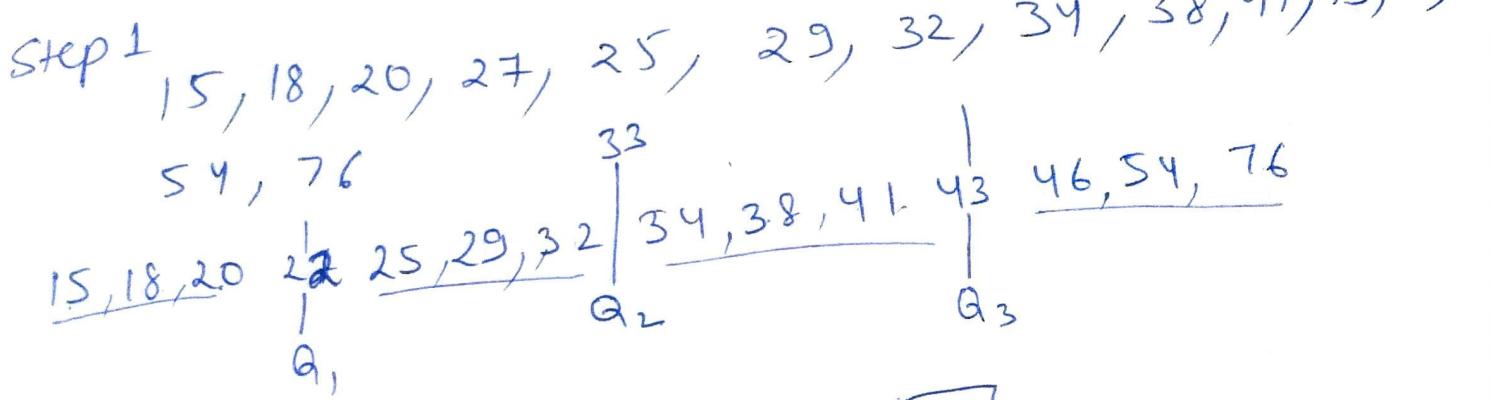
$\Rightarrow [-12, 52]$ hence minimum(8) & maximum(43) is not an outlier.



If the data have outlier then

Let us take an example

15, 18, 20, 27, 25, 29, 32, 34, 38, 41, 43, 46, 54, 76, 22, 18, 34, 76, 29, 15, 41, 46, 25, 54, 38, 20, 32, 43, 45, 22



$$IQR = Q_3 - Q_1 = 43 - 22 = 21$$

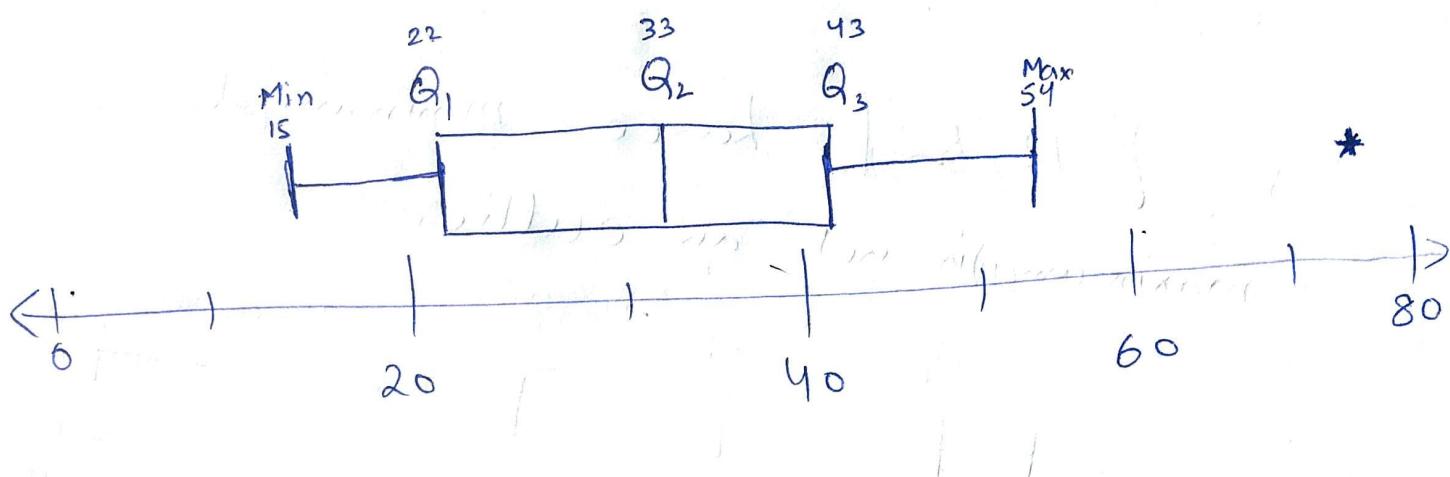
$$[Q_1 - 1.5 \text{ IQR}, Q_3 + 1.5 \text{ IQR}]$$

$$[22 - 1.5 \times 21, 43 + 1.5 \times 21]$$

$$[-9.5, 74.5]$$

Min 15 is in range

Max 76 is outside the range



Histogram

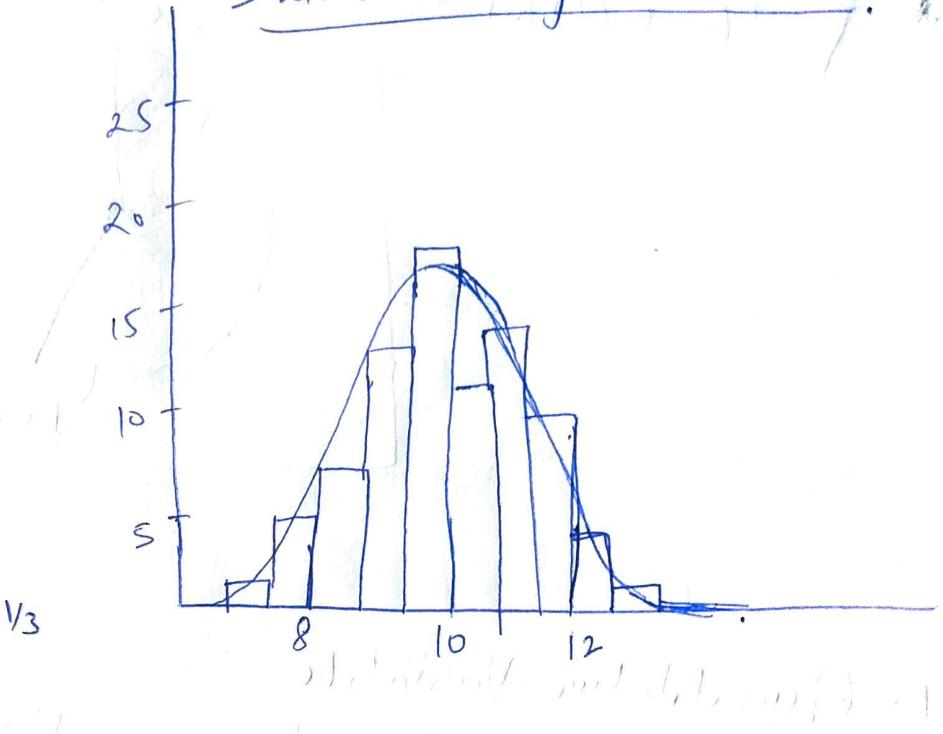
How to Choose Number
of bins:

Scott's rule

$$\hat{\Delta}_b = 3.49 \hat{\sigma} n^{-1/3}$$

$\hat{\sigma}$: estimate of Standard deviation

n : Number of data points



Freedman-Diaconis Rule

$$\hat{\Delta}_b = 2(Q_{75} - Q_{25}) n^{-1/3}$$

where $Q_{75} - Q_{25} = \text{IQR}$ (Inter Quartile Range)

"Optimal Data-Based Binning of Histogram"

by Kelvin H. Knuth}, arXiv: physics/0605197.

① Number of bins k can be assigned directly by choosing bin width h as

$$k = \left\lceil \frac{\max - \min}{h} \right\rceil$$

② Square - Root Rule

$$k = \lceil \sqrt{n} \rceil$$

③ Sturges formula

$$k = \lceil \log_2 n \rceil + 1$$

④ Rice Rule

$$k = \lceil 2^{\sqrt[3]{n}} \rceil$$

⑤ Doane's Formula: Doane's modifies Sturges rule to improve the performance for non-Gaussian data

$$k = 1 + \log_2(n) + \log_2 \left(1 + \frac{|g_1|}{\sigma_{g_1}} \right)$$

g_1 : 3rd Moment (Skewness) of the

$$\sigma_{g_1} = \sqrt{\frac{6(n-2)}{(n+1)(n+3)}}$$

87

Shimazaki & Shinomoto's Rule

$$\arg \min_h \frac{2\bar{m} - v}{h^2}$$

where \bar{m} & v are mean & biased variance

of a histogram with bin width h

$$\bar{m} = \frac{1}{K} \sum_{i=1}^K m_i \quad v = \frac{1}{K} \sum_{i=1}^K (m_i - \bar{m})^2$$

Variable bin widths

where number of samples in each bin is expected to be approximately equal

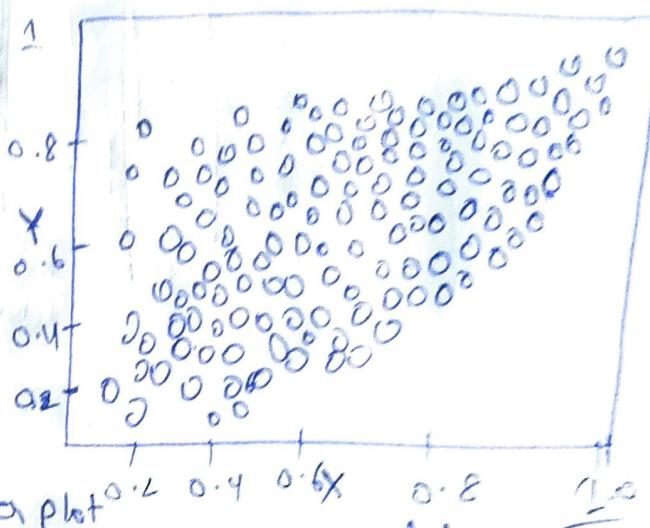
②

Graphical Representation: Multiple Variable

Scatter Plots: Two Quantitative Variables (correlation)

It is used to Show a relationship between two variables.

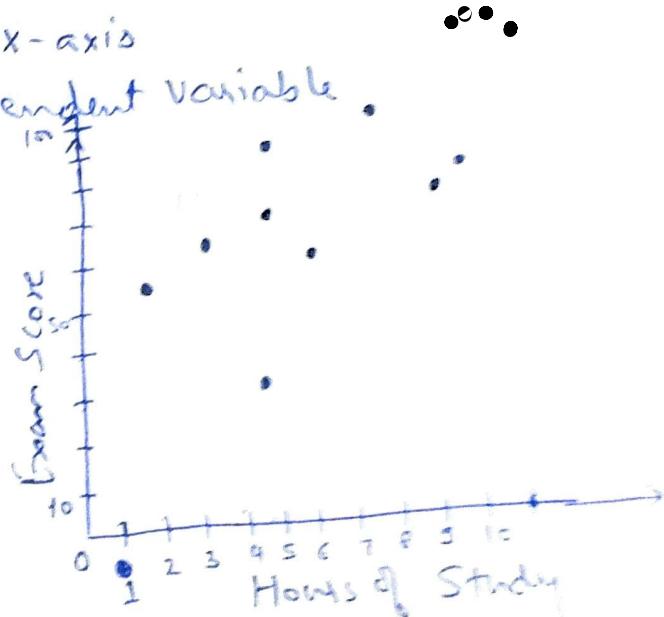
A symbol, usually a dot is used to show a data pair.



- When to use Scatter plot**
- * When you have paired Numerical data
 - * Dependent variable have multiple value for each value of independent variable
 - * To determine whether the two variables are related
- x-axis
Independent Variable

Ex

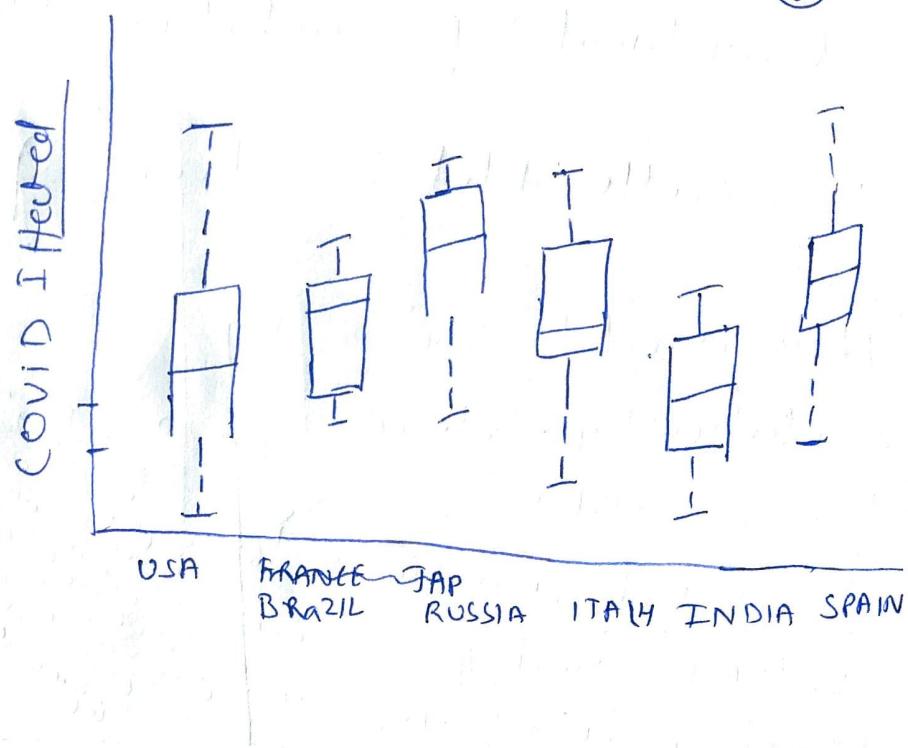
Hours Studying	Exam Score
2	53
4.5	35
5	91
5	72
6	60
3	62
10	85
9.5	78
8	99



Find the Relationship between two variables if exist!

Matrix Scatter Plot: To represent Correlation between more than two variables.

Box Plots



Contingency Table: 2 Categorical Variables with frequency of occurrence as the theme (4)

		WORK EX	
		Y	N
MBA	Y	22	7
	N	32	17

How many of Employees have MBA before joining the organization

How many of Employees have work experience before joining the organization.

Ex Smoker/Non Smoker. in Org. joint dist.

		Smoker	Non-Smoker	Total
		Male	Female	
Gender	Male	71.55.47%	44.44.57%	116.57.11%
	Female	34% - 75.1%	55.25.11%	88.75%
Total	100% - 22.1%	97.75% - 78.1%	20.3	100%

Relative frequency
 $P(\text{smoking} | \text{Male})$, Rankin's hist.