

COMPLETE STATISTICS FOR DATA SCIENCE

Data Scientist

Data Intelligence

Business Intelligence

Data Analyst

Descriptive Statistics

Inferential Statistics

Descriptive Statistics

Measure of central Tendency

Measure of Dispersion

Mean, Median, Mode, Variance, Standard Deviation

Related to Summarizing the Data

Histogram, pdf, cdf,

Probability, Permutation & Combination

Gaussian's Distribution

Log Normal Distribution

Transformation & Standardization

Binomial Distribution

Q-Q Plot

Bernoulli Distribution

Pareto or Power Law Distribution

Standard Normal Distribution

Inferential Statistics

Z Test

t Test

Annova Test or F test

Chi Square Test

Hypothesis Testing

Confidence Interval

By Python

P value

What is Statistics ?

Science of Collection, organizing and analyzing Data.



Used for **BETTER DECESSION MAKING**

What is Data?

Facts or Process of Information that can be measured.



Ex: Ages of Student of Class
30, 25, 20,28,20 (Data)

Types of Statistics ?

1. Descriptive Stats



It consists of organizing and summarizing of Data

2. Influential Stats



Technique where in we use the **DATA**, that we have measured to form **CONCLUSION**

Example of Descriptive stats

Class room of Math Students: 20 No.

Marks of Ist Semester

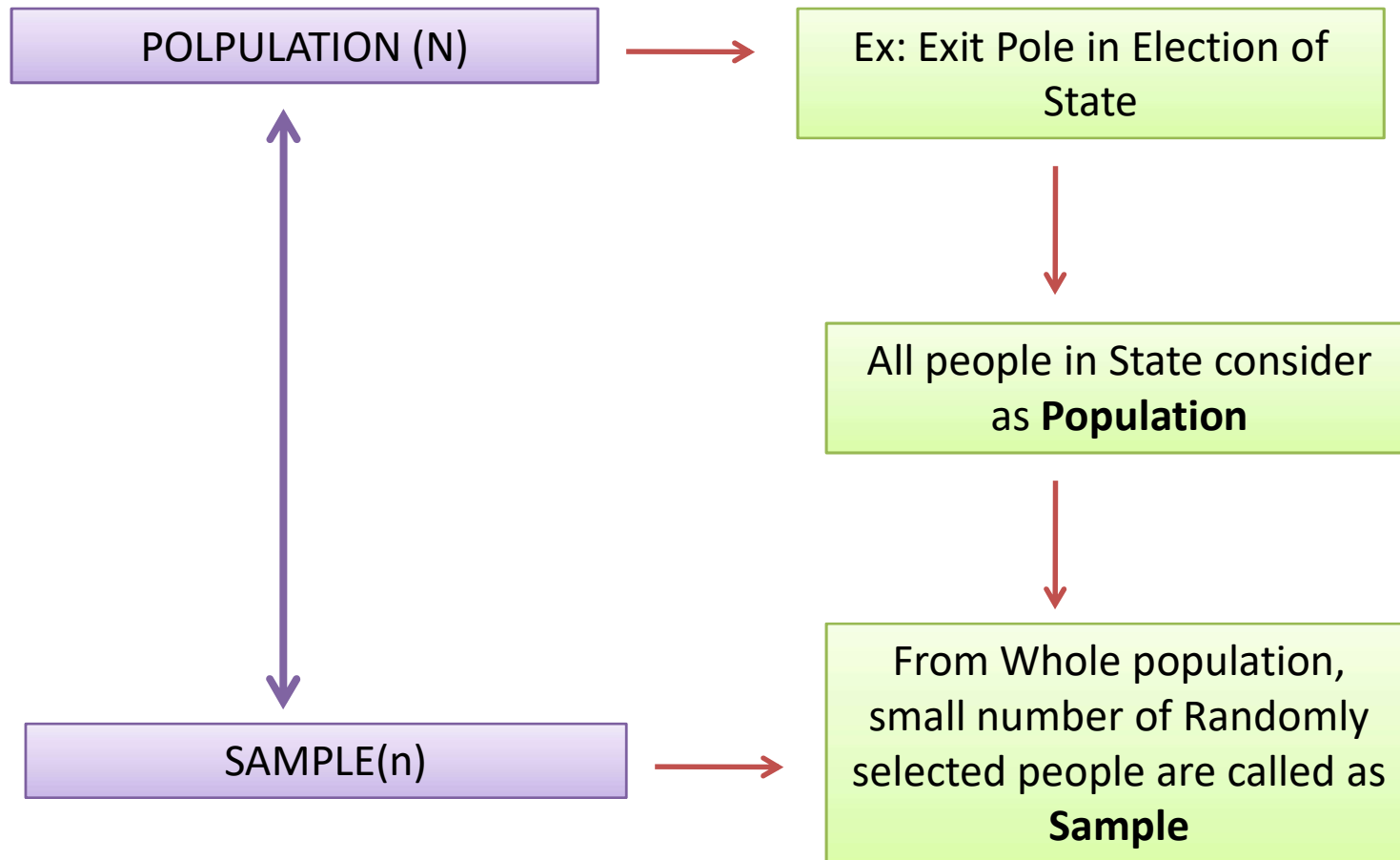
84,85,89,90,75,55, 84,85,89,90,75,55 84,85,89,90,75,55,35,40

What is the Average Marks of Students of class room?

Example of Influential stats

Are the Average marks of the students of Maths classroom similar to the marks of Math Class in college?

POPULATION & SAMPLE



METHOD OF SAMPLING TECHNIQUES

1. Simple Random

Mostly used; Ex. Exit Pole opinion of Election

2. Stratified Sampling

Population (N) is split into
Non Overlapping Groups (Strata)
Example:
a. Gender (Male/ Female)
b. Survey Based on Specific age group

3. Systematic Sampling

Pick up every nth (specific number) individual from population

4. Convenient Sampling

- Survey from specific domains
- Related to specific Topic

What is Variable?

Variable is the property that can have different value.

Ex: Height, Age, Body Weight

Two Type

Quantitative

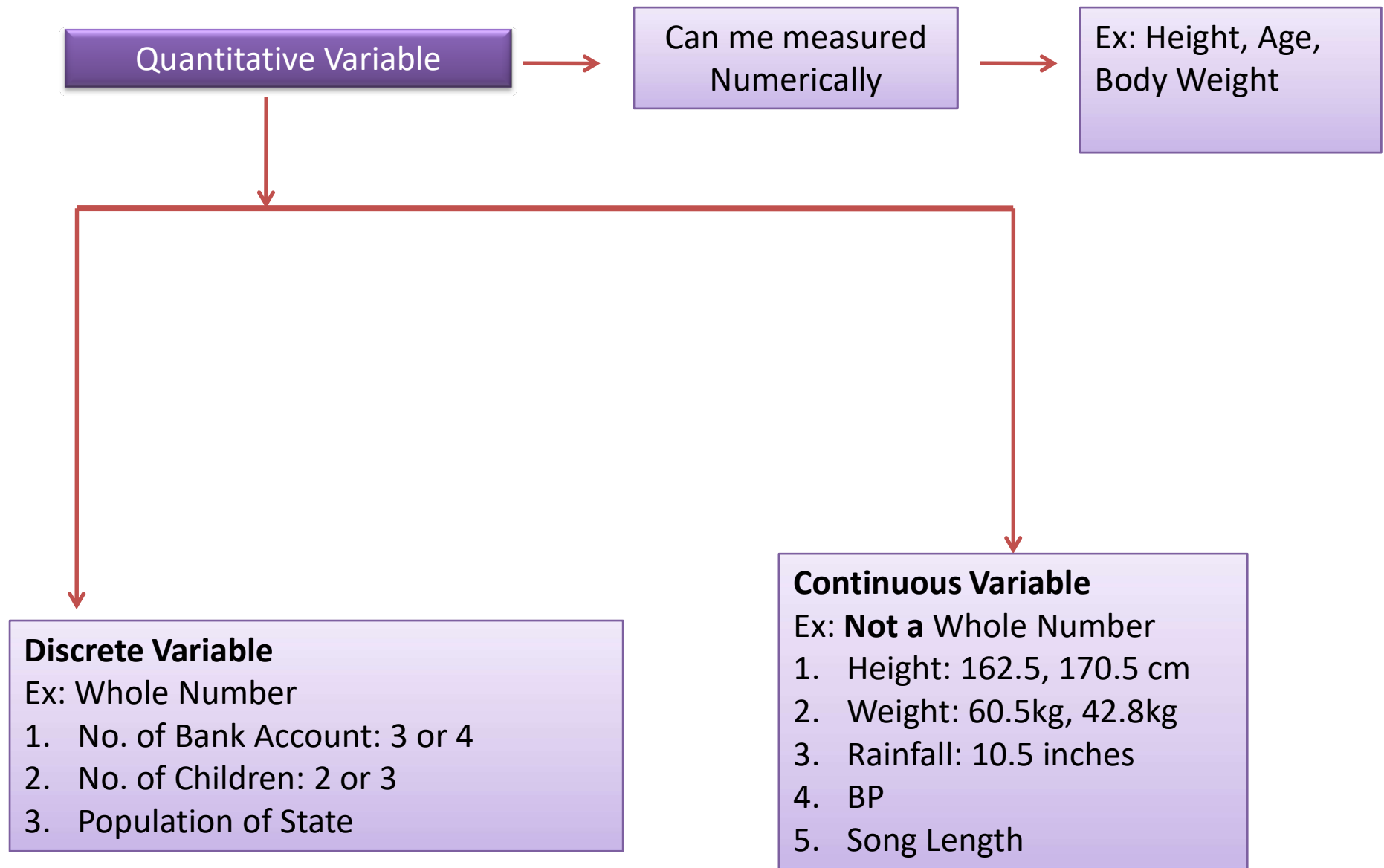
Can be measured Numerically

Ex: Height, Age, Body Weight

Qualitative/ Categorical

Based on characteristic, we can derive categorical variable

Ex: Gender (Male/Female)



Variable Measurement Scale

1. Nominal

Categorical Data

Ex: Gender- Male/Female, Flower's Type: Lilly, Rose

2. Ordinal

Order of Data Matters, value doesn't matter.

Ex: In case of Marks & Ranks of students, we can analyze performance through Ranks (Order) only

3. Interval

Order & Value matter, Natural Zero not exist

Ex. Interval of Temperature (70-80F; 80-90F) or Distance like 10-20km

4. Ratio

Quantitative scale, true zero exist, equal interval between point.

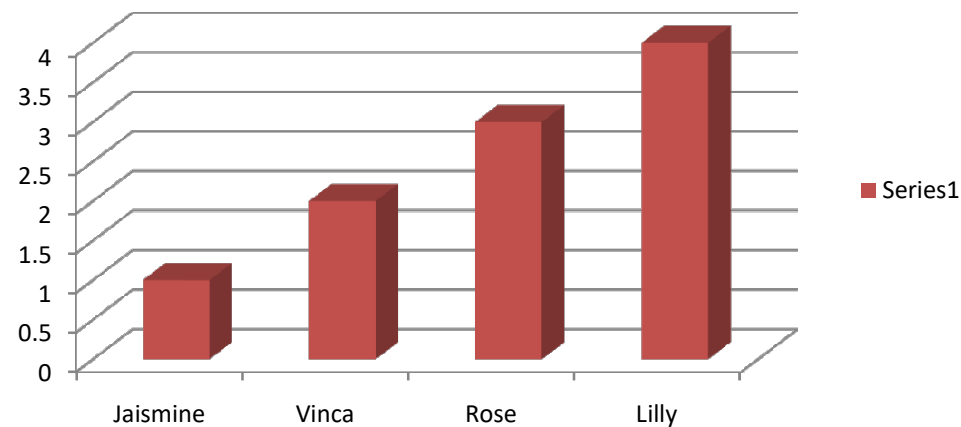
Ex. Length, Area

Frequency Distribution

Ex: Flowers- Rose, Lilly, Rose, Rose, Vinca, Jaismine, Vinca, Lilly, Lilly, Lilly

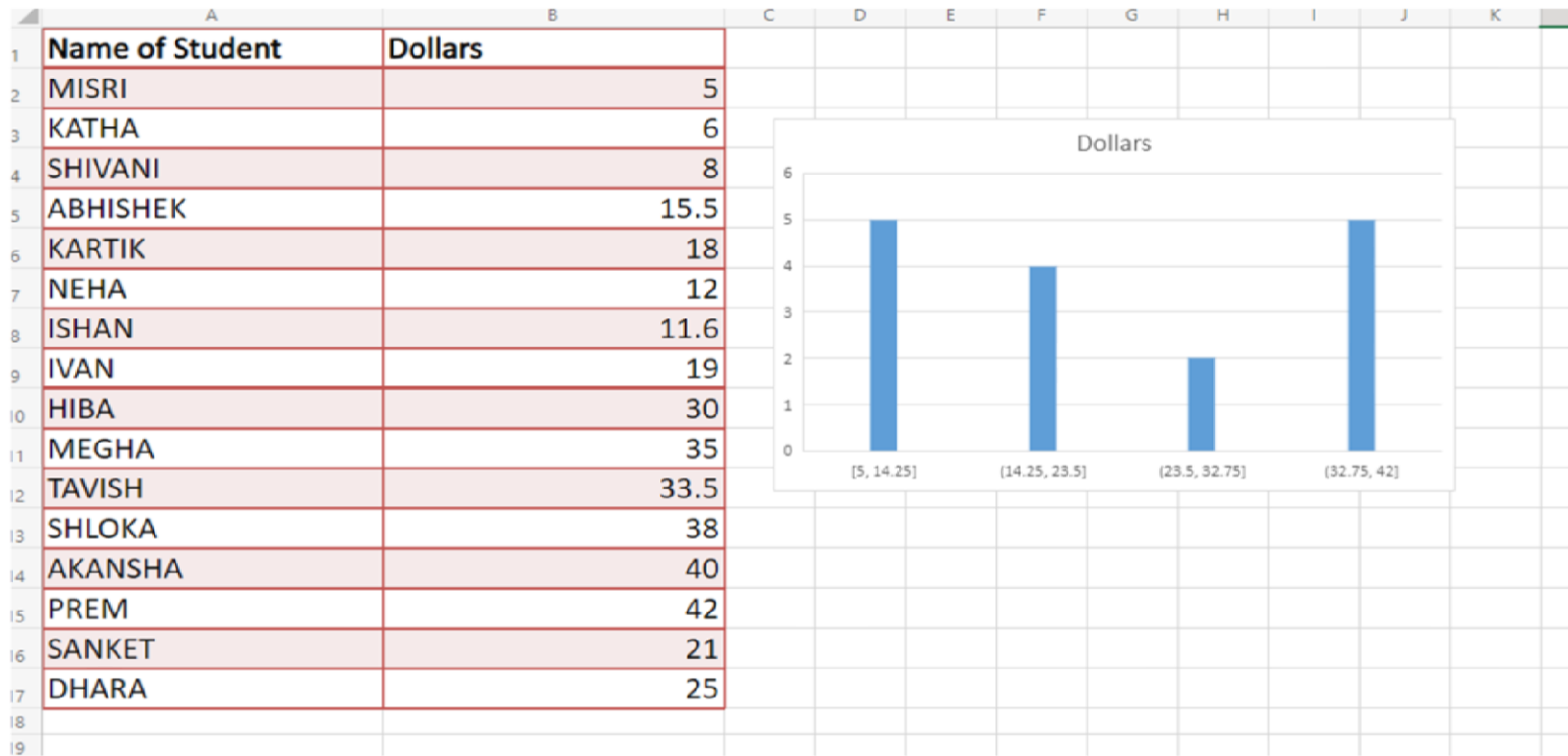
Flowers	Frequency	Cumulative Frequency
Rose	3	3
Vinca	2	5
Lilly	4	9
Jaismine	1	10

Bar Chart



Value is
Discrete
Variables

Histogram: Data Should be continuous Variable

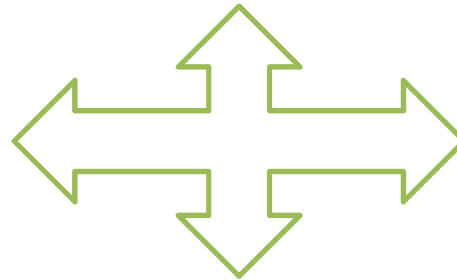


Note:
In above Histogram, we have selected 4 bins

MEASURE OF CENTRAL TENDENCY

Refer to the measure used to determine the central of distribution of data

Arithmetic **MEAN** for
Population(N) and
Sample(n)



Arithmetic **MODE** for
Population(N) and
Sample(n)

Arithmetic **MEDIAN** for
Population(N) and
Sample(n)



Arithmetic **MEAN** for Population(N) = Calculate Average only

Arithmetic **MEDIAN** for Population(N) and Sample(n)= Select Central or middle element of data

Arithmetic **MODE** for Population(N) and Sample(n)= Select most frequent Element of data

MEASURE OF CENTRAL TENDENCY

Example of Data: MODE

Suppose we have data with some missing elements

Type of Flower	Petal Length
Rose	
Vinca	
Lilly	
Rose	
Vinca	
Rose	
?	
?	

Missing elements: can be determined
with modes



Arithmetic **MODE** for Population(N) and
Sample(n)= Select most frequent
Element of data

**Mode will work with Categorical
Variable**

MEASURE OF CENTRAL TENDENCY

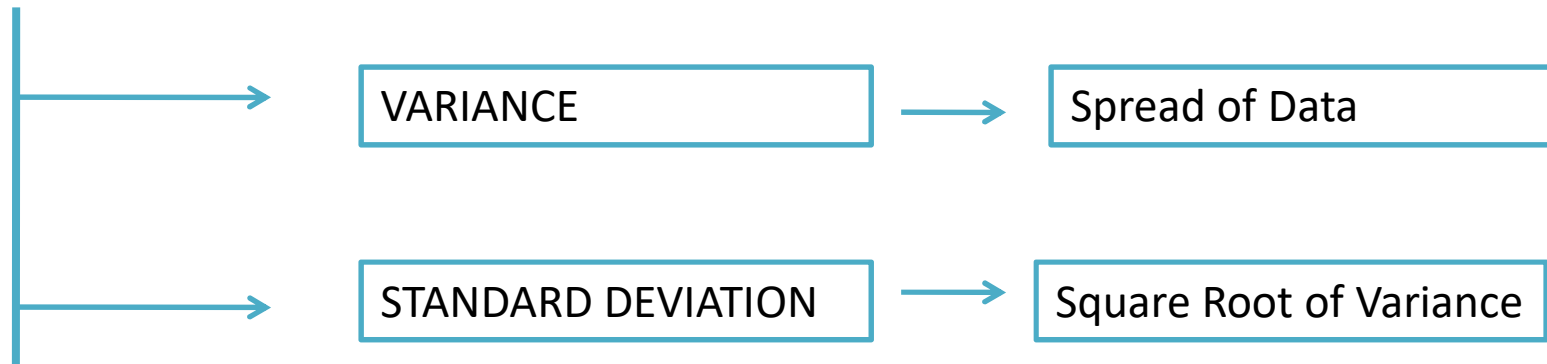
Example of MEAN

Suppose we have data (Quantitative variable) with some missing elements

Name of Student	Age
MISRI	30
KATHA	15
SHIVANI	35
ABHISHEK	35
KARTIK	?
NEHA	30
ISHAN	?

→ MEAN WILL WORK
BETTER TO
CALCULATE

MEASURE OF DISPERSION [SPREAD]



EXAMPLE OF TWO DATASET

DATASET: 1

1	1	2	2	4
---	---	---	---	---

DATASET: 2

2	2	2	2	2
---	---	---	---	---

MEAN OF BOTH DATASET IS
SAME that is **5**

It can be distinguish
with the help of
VARIANCE

So, How we can
distinguish Distribution
????

PERCENTILE & QUARTILE

Use to find Outliers or Odd numbers in Data

PERCENTILE is the value, below which a certain percentage of observation lie.

Ex: Data Set

Q1 . = What is the Percentile Ranking of 10?

2	2	3	4	5	5	5	6	7	8	8	8	8	8	9	9	10	11	11	12
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	----	----	----

Percentile value of 10 = $\frac{\text{No. of value below 10}}{n} \times 100$
n (sample size)

Percentile value of 10 = $\frac{16}{20} \times 100 = 80\%$

PERCENTILE & QUARTILE

Ex: Data Set

Q 2 . = What value exist at Percentile Ranking of 25%?
(Reverse of Question 1)

2	2	3	4	5	5	5	6	7	8	8	8	8	8	9	9	10	11	11	12
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	----	----	----

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

Solution: Value =
 $\frac{25 \times (20+1)}{100} = 5.25$ (This is INDEX value)

2	2	3	4	5	5	5	6	7	8	8	8	8	8	9	9	10	11	11	12
1	2	2	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

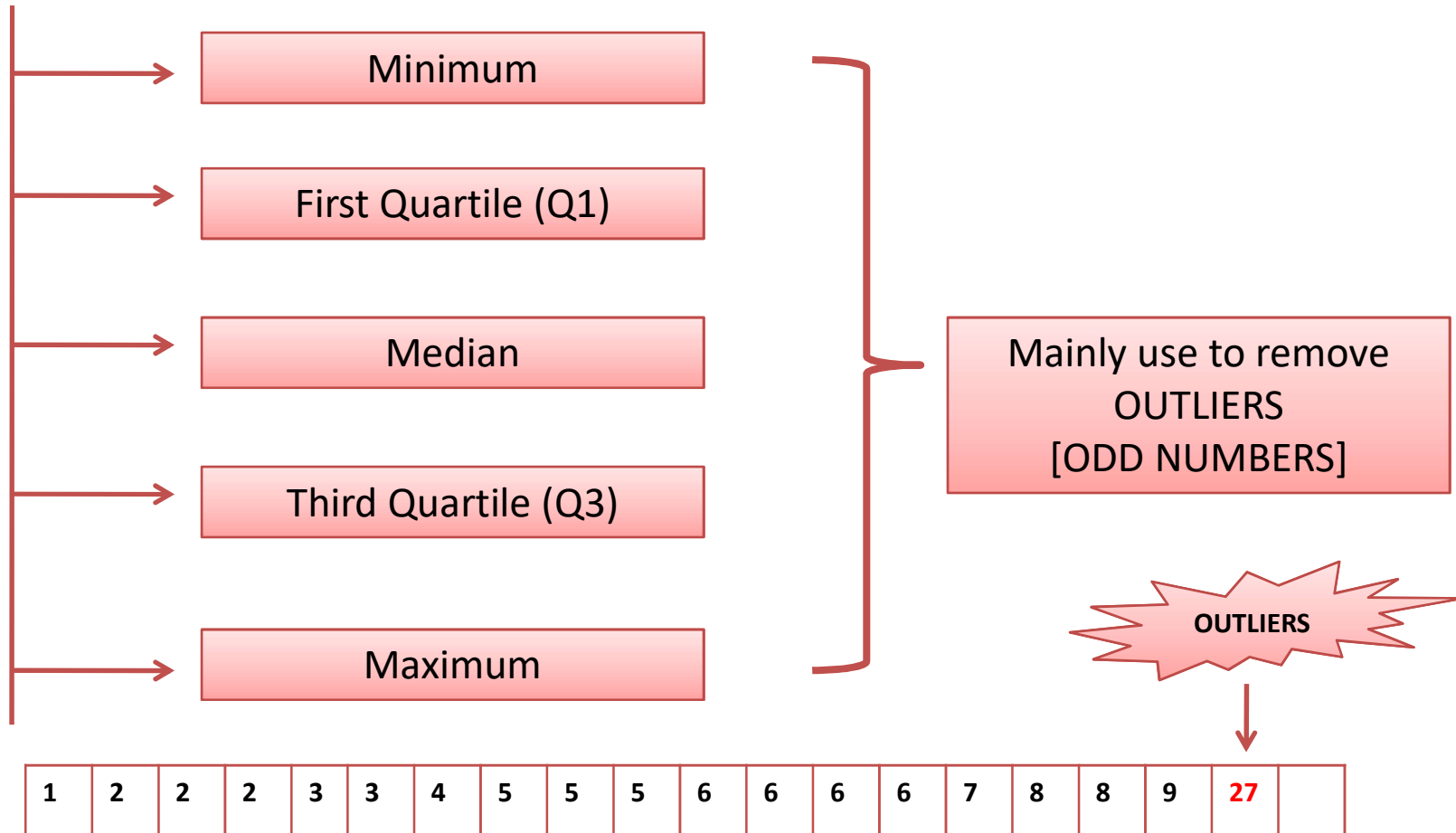
Index value

Index value 5.25 comes between 5 & 6 index

Will take
Average of
Index value of 5
& 6

In this case: Index value of 5 & 6 is 5,5: Average of $5+5/2=5$

FIVE NUMBER SUMMARY



FIVE NUMBER SUMMARY

How to Remove Outliers from following Data Set?

1	2	2	2	3	3	4	5	5	5	6	6	6	6	7	8	8	9	27	
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	--

Answer:

We need to define LOWER FENCE & HIGHER FENCE

Lower Fence: Below the lower fence
all number will be treated as Outliers

Lower Fence= $Q1 - 1.5 (IQR)$

Upper Fence: Above the Upper/Higher
fence all number will be treated as
Outliers

Upper Fence= $Q3 + 1.5 (IQR)$

Where:

IQR: Inter Quartile Range

IQR: $Q3 - Q1$

Q1: 25 Percentile

Q3: 75 Percentile

FIVE NUMBER SUMMARY: REMOVAL OF OUTLIERS

Solution:

Data Set

1	2	2	2	3	3	4	5	5	5	6	6	6	6	7	8	8	9	27	
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	--

1	2	2	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----

Q1=25 Percentile

Q3=75 Percentile

Index Number

$$\frac{25}{100} * (n+1)$$

$$\frac{75}{100} * (n+1)$$

Upper Fence= $Q3 + 1.5 (IQR)$

Lower Fence= $Q1 - 1.5 (IQR)$

$$\frac{25}{100} * (19+1)$$

$$\frac{75}{100} * (19+1)$$

$IQR = Q3 - Q1$

$IQR = 7 - 3 = 4$

Result: 5 or
Index No. 5

Result: 15 or
Index No. 15

Data in index 5 is "3"

Data in index 15 is "7"

FIVE NUMBER SUMMARY: REMOVAL OF OUTLIERS

Solution: Continue.....

Upper Fence= $Q3 + 1.5 (IQR)$

→ Upper Fence= $7 + 1.5 (4)$

→ Upper Fence= 13

Lower Fence= $Q1 - 1.5 (IQR)$

→ Lower Fence= $3 - 1.5 (4)$

→ Lower Fence= -3

Data Set

1	2	2	2	3	3	4	5	5	5	6	6	6	6	7	8	8	9	27	
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	--

As per definition of OUTLIER, Numbers that exist below the lower fence and those exist above Upper fence in Data Set are consider as OUTLIERS.

OUTLIERS

Data Set after removal of OUTLIER

1	2	2	2	3	3	4	5	5	5	6	6	6	6	7	8	8	9
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

FIVE NUMBER SUMMARY: CONCLUSION

Solution: From above solutions, now we can figure out following points

Data set: after removal of Outliers

1	2	2	2	3	3	4	5	5	5	6	6	6	6	7	8	8	9
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Minimum



1

First Quartile (Q1)



3

Median



5

Third Quartile (Q3)

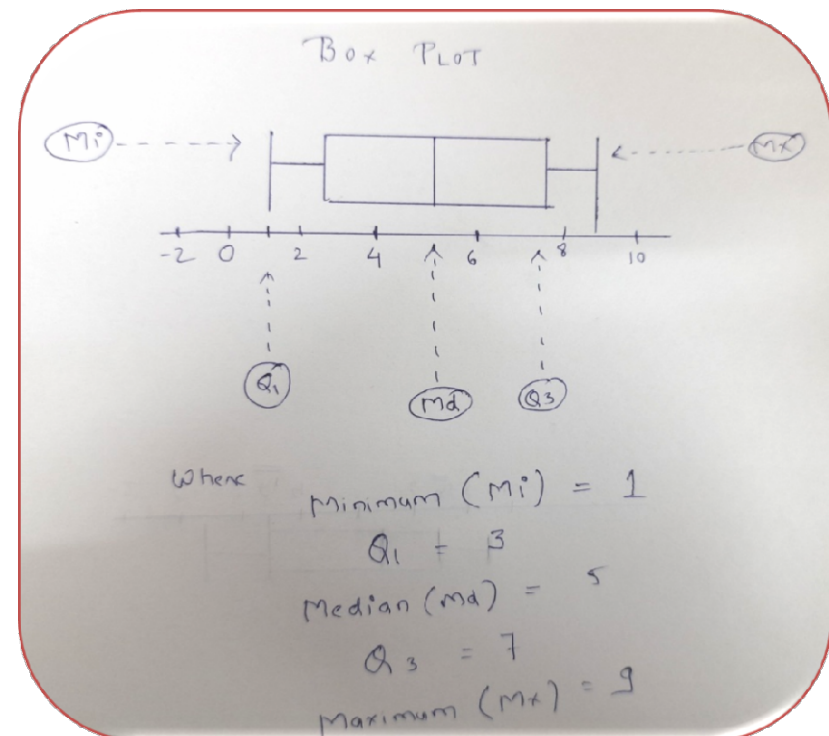


7

Maximum



9



Note: In above diagram, M_i , M_x and M_d is used for illustrative purpose only .

INTERVIEW QUESTION: WHAT IS THE USE OF BOX PLOT?

Box Plot: It is used for the determination of Outlier.
: It gives u the visualization way, where the outlier is actually presents

INTERVIEW QUESTION: How BOX PLOT uses for the removal of Outlier?

With the help of Five Number Summary, like Minimum, Q1, Median, Q3 and maximum, we can identify Upper fence and lower Fence of Data Set and then with the help of Box Plot, we can remove Outliers