

## Descriptive Statistics

### 1) → Measure of Central tendency

- Mean
- Median
- Mode

#### Population mean ( $\mu$ )

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

#### Sample mean ( $\bar{x}$ )

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

#### Median

\* Sort the values either asc or desc order

\* create the mid value

\* If you get mid 2 values take avg of those two values

→ outlier (far away from 5)

Q. [1, 2, 2, 3, 4, 5]

[1, 2, 2, 3, 4, 5, 100]

mean = 2.8

Mean = 16.7

Median = 2.5

Median = 3

\* Mean will be affected by outlier where as Median won't affect by outlier

\* These are calculated for null values imputation

#### Mode

Most repetitive value.

## Measure of Dispersion

Eg:- Person 1 Person 2

Monday	7:30 AM	8 am
Tuesday	7:45 AM	11 am
Wednesday	8 AM	9 am
Thursday	7:15 AM	7 am
Friday	7 am	10 am
Prediction	7-8 15 min 30 min	9-10 2 hr 1 hr

Variance high

Prediction accuracy low

Variance less

Prediction accuracy high

→ Variance

→ standard deviation

→ Range

#### Population Variance ( $\sigma^2$ )

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

#### Sample Variance ( $s^2$ )

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

[n-1 = degree of freedom]  
Imp. (in interview questions)

a. calculate  $\sigma^2$  &  $s^2$

{1, 2, 2, 3, 4, 5}

$$\sigma^2 = \frac{(1-2.8)^2 + (2-2.8)^2 + (2-2.8)^2 + (3-2.8)^2 + (4-2.8)^2 + (5-2.8)^2}{6}$$

$$= \frac{3.36 + 0.64 + 0.64 + 1.36 + 4.69}{6} = \frac{10.69}{6} = 1.78 \text{ (approx)}$$

$$\sqrt{\sigma^2} = \sqrt{1.8 \text{ km}^2}$$

$$\sigma = \underline{1.34 \text{ km}} \text{ (standard deviation)}$$

$$k = 2.8$$

Standard deviation is square root of  $\sigma^2$  ( $\sqrt{\text{variance}}$ )

$$\text{Population standard deviation } (\sigma) = \sqrt{\sigma^2}$$

$$\sigma = \sqrt{1.8}$$

$$\sigma = 1.34$$

Sample variance ( $s^2$ )

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{10.83}{5} = \underline{2.16}$$

Sample SD ( $s$ )

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

$$s = \sqrt{2.16}$$

$$s = 1.46$$

Range (maximum - minimum)

$$= 5 - 1 = \underline{4}$$

Percentile and Quartile

Percentile is the value below which a certain percent of observations will come

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

How much % of data will come below 6?

$$\text{Percentile rank of } x = \frac{\text{No of value below } x}{N} \times 100$$

$$= \frac{7}{11} \times 100$$

$$= \underline{63\%} \text{ observation data value } x < 6$$

Quartile

Quartile helps to find the value which is present at the given percentile rank.

\* Which value is present at 25%?

90%?

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

$$= \frac{90}{100} \times 12$$

$$= \frac{25}{100} \times 12 = 3 \rightarrow \text{index}$$

$$\text{value} = 2$$

$$= 10.8 \rightarrow \text{index}$$

10  $\rightarrow$  2nd value before dec

$$\text{value} = \underline{7}$$

## Five Number Summary

- 1) Minimum
- 2) First Quartile ( $Q_1$ ) 25%
- 3) Median ( $Q_2$ ) 50%
- 4) Third Quartile ( $Q_3$ ) 75%
- 5) Maximum

Note: choose these 5 numbers after removing the outliers from the data by first boundary values.

[lower fence upper fence]

$$LF = Q_1 - 1.5(IQR)$$

$$UF = Q_3 + 1.5(IQR)$$

IQR = Inter Quartile Range

$$IQR = Q_3 - Q_1$$

$$\# \{1, 1, 1, 2, 3, 4, 4, 4, 5, 5, 6, 7, 7, 8, 8, 9, \boxed{28, 36}\} \text{ outliers}$$

$$\begin{aligned} LF &= IQR = \frac{75}{100} \times 17 - \frac{25}{100} \times 17 \\ &= 13.5 - 4.25 \quad [\text{inolen}] \\ &= 9.25 \\ &= 9 \end{aligned}$$

$$LF = 9 - 1.5 \times 5 = 4.5$$

$$UF = 9 + 1.5 \times 5 = 15.5$$

Anything  $< 4.5$  &  $> 15.5$  is a outlier. (remove them).

$$\text{Minimum} = 1$$

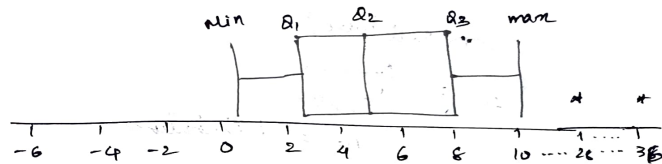
$$(Q_2) \text{ Median} = 5$$

$$\text{Maximum} = 9$$

$$Q_1 = 3$$

$$Q_3 = 8$$

Box plot



This graph is used to the outliers.

Different types of Distributions

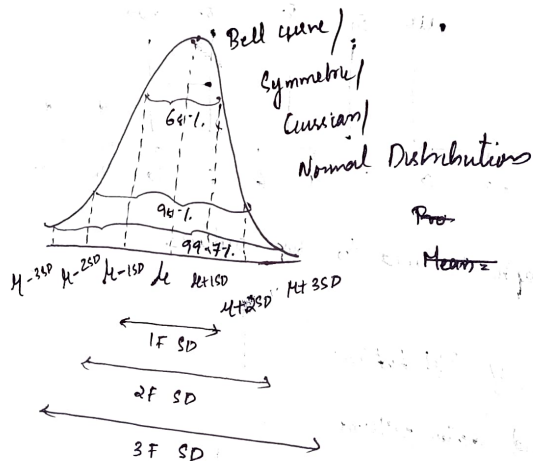
- 1) To understand data patterns
- 2) To summarize the data easily
- 3) To calculate Probability
- 4) To make prediction and Decision.
- 5) To choose right statistical test.

There are 2 category of distributions :-

- Continuous distribution (Numbers, Numerical Distributions)
- Discrete distribution (Categorical Distributions)

- |                                 |               |
|---------------------------------|---------------|
| 1) Normal Distribution          | } CD          |
| 2) standard normal Distribution |               |
| 3) Bernoulli Distribution       | } categorical |
| 4) Binomial Distribution        |               |
| 5) Poisson Distribution         |               |

## 1) Normal Distributions



### Property

Mean = Median = Mode

68% - 95% - 99.7% is called as confidence interval

### Empirical Rules :-

68% of data will present in 1SD.

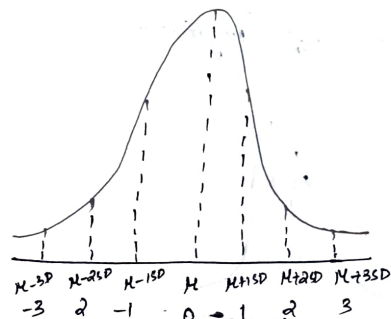
95% of data will present in 2SD.

99.7% of data will present in 3SD.

### 2) Standard Normal Distribution

$$\mu = 0$$

$$\sigma = 1$$



$$Z \text{ score} = \frac{x_i - \mu}{\sigma}$$

### Normal Dist Data

2
7
5
4
1
3
5
μ = 3.86
σ = 2.00

### Std Normal Dist Data

-0.93
1.57
0.57
0.07
-1.43
-0.43
0.57
μ = 0
σ = 0.94 = 1