

Basic Stats

Probability

1. Introduction to Basic Terms

2. Variables

3. Random Variables

4. Population, Sample, Population Mean, Sample Mean

5. Pop Distribution, Sample Distribution, Sampling Distribution

6. Mean, Median, Mode - Measures of Central Tendency

7. Range

8. Measure of Dispersion

9. Variance

10. Std. Deviation

11. Gaussian / Normal Distribution

25. Q-Q Plot

26. Chebychev's Inequality

27. Discrete & Continuous Distribution

28. Bernoulli & Binomial Distribution

29. Log Normal Distribution

30. Power Law Distribution

31. Box Cox Transform

32. Poisson Distribution

33. Application of Non-Gaussian Distributions

## Intermediate Stats

12. Std. Normal Distribution

21. Covariance

13. Z score

22. Pearson Correlation Coefficient

14. Prob. Density Function

23. Spearman Rank Correlation

15. Cumulative Distribution Function

24. Hypothesis Testing

16. Hypothesis Testing

17. Many different plotting graphs

18. Kernel Density Estimation

19. Central Limit Theorem

20. Skewness of Data

~~DATA PROCESSING~~

①

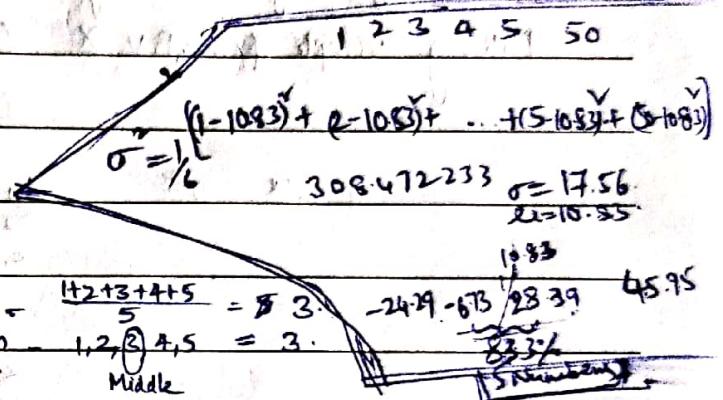
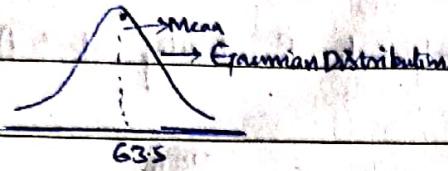
## V. Mean, Median, Mode - ~~Measures of~~

(i) Mean :- Measure of Central Tendency.

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

Sample of Height = {168, 170, 150, 160, 182, 146, 175}

$$\bar{x} = \frac{168 + 170 + 150 + 160 + 182 + 146 + 175}{7} = 163.5$$



(ii) Median :-

$$\{1, 2, 3, 4, 5\} \rightarrow \text{Mean} = \frac{1+2+3+4+5}{5} = 3. \quad \text{Median} = \frac{1, 2, 3, 4, 5}{5} = 3.$$

$$\text{Mean} = \frac{1+2+3+4+5+50}{6} = 10.83 \quad (65\%)$$

$$\text{Median} = \frac{1, 2, 3, 4, 5, 50}{6} = \frac{3+4}{2} = 3.5$$

Write in Ascending Order to find Median.

Size: Odd No - Middle

$$\text{Even No} - \frac{m_1 + m_2}{2}$$

Age

23

24

27

32

35

21

=

3 Data Mining.

If I remove the missing Date. I may lose some information.

In that case, I will fill Data using Mean, Median, Mode.

Mean - If no outliers are present - Prefer Mean.

Median/Mode - More outliers are present - Prefer Median/Mode.

(iii) Mode -

We touch your electricity everyday!

Statistics

I.  
1. Population :-

Question : Average Height of all people in a state

 $N \rightarrow$  Population

 $n \rightarrow$  Sample

Solution =

$$\mu = \frac{1}{N} \sum_{i=1}^{10^6} x_i \Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

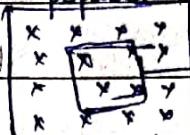
 $\mu \rightarrow$  Population Mean

 $\bar{x} \rightarrow$  Sample Mean

2. Sample :-

To get Data from all people is difficult (Population).

population



Sample

So, We consider Sample:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{10000} x_i \Rightarrow \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

 $n < N$ 

When Sample Size increases  $n \approx N$ 

## II.

3.1 Random Variables :- A random Variable which stores Data. Types < Numerical R.V.

Types of R.V. :- Numerical R.V. & Categorical R.V.

(i) Discrete Variables - Whole Number.

e.g. No. of Bank Accounts a person has?

Ans: 2, 5, 7... Can't be 2.5 (decimal)

(ii) Continuous Random Variables - Within a range, Any Value.

10.1, 10.2, 10.368...

Depends on Scaling e.g. - Weight of a Person - 63kg.

Interest Rate

63.2 kg.

Salary.

depends on scaling

63.29 kg.

Categorical Variables for each record it is continuous, repeated.

Gender is Categorical R.V.

e.g.

Demographic Statistics	House Owner	Age	Gender	No. of People in House	Height of House Owner
Number C.R.V.	Atmospheric	Numerical D.R.V.	Gender	Numerical D.R.V.	Numerical C.R.V.

We touch your electricity everyday!

### Distributions:

#### 1. Gaussian / Normal Distribution:

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

Mean S.D.

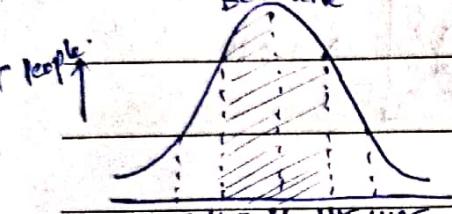
$$x \sim GD(\mu, \sigma) \text{ or } N(\mu, \sigma)$$

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

$$GD(\mu, \sigma^2) \text{ or } N(\mu, \sigma^2)$$

$$\text{Std. deviation } \sigma = \sqrt{\text{Var}}$$

Prob Density Fun / histogram  
Bell Curve



$$x = x_i$$

→ Height

Probability of  
Having  
Data point for  
(in 3σ range)

Mean

2σ distance  
2σ deviation away

Empirical formula:

$$Pr(\mu - \sigma \leq x \leq \mu + \sigma) = 68\%$$

$$Pr(\mu - 2\sigma \leq x \leq \mu + 2\sigma) = 95\%$$

$$Pr(\mu - 3\sigma \leq x \leq \mu + 3\sigma) = 99.7\%$$

Bell Curve & Symmetrical.

Dataset

e.g. Br. (i) Distribution of Height

(ii) IRIS Dataset - Petal length-

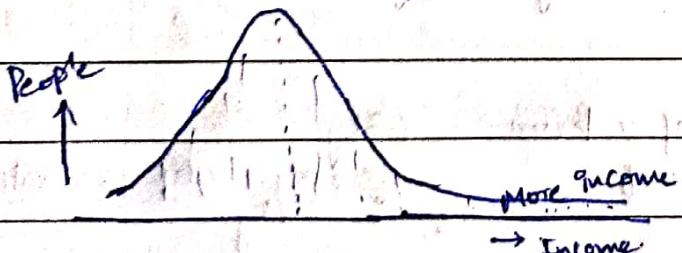
$$GD \Rightarrow \mu=0 \sigma=1$$

#### 2. Log Normal Distribution

$x \sim \text{Log Normal Distribution}$  if  $\ln(x)$  is Normally Distributed.

$$x = \{x_1, x_2, x_3, \dots, x_n\}, \ln(x_1), \ln(x_2), \dots, \ln(x_n)$$

$$\ln(x) \sim N(\mu, \sigma)$$



At Similar to G.D, but Right Hand side goes more flatten.

→ Right Skewed.

Dataset e.g.: (i) Income of the people

(ii) Amazon Product Review Description Length

More income people are very less

Why do we use all types of distributions?

Money Spent on:

	R&D	Marketing	Profit	States	Company
1	100000	1000000	-	-	-
2	1-	-	-	-	-
3	-	-	-	-	-
4	-	-	-	-	-
5	-	-	-	-	-
6	Gaussian	Log Normal	-	-	-

As per Domain Knowledge, If we know R&D is Gaussian Distribution. & Marketing is Log Normal Distribution.

R&D  $\rightarrow$  Scaling to Std. Normal Distribution ( $\mu=0, \sigma=1$ )

$$\frac{\bar{x}_i - \mu}{\sigma}$$

Marketing  $\rightarrow$

$$\text{Marketing} \xrightarrow{\log(\text{marketing})} \text{will follow Normal Distribution.}$$

$$1000 \xrightarrow{\ln(1000)} \ln(x)$$

$$\xrightarrow{\ln(x)} \ln(x)$$

$\ln \sim \text{Gaussian } (\mu=0, \sigma=1)$

$$\downarrow \quad S.N.D = \frac{\bar{x}_i - \mu}{\sigma} \quad \text{Scaling done.}$$

It is called Log Normalization.

Comparing

Now after Scaling both R&D, Marketing  $\rightarrow$  Analysis with More Accuracy.

Conclusion:

Q. Covariance :- Useful in Data Preprocessing.

Size      Price

1200 sqm	100/-
1500 sqm	200/-
1800 sqm	300/-

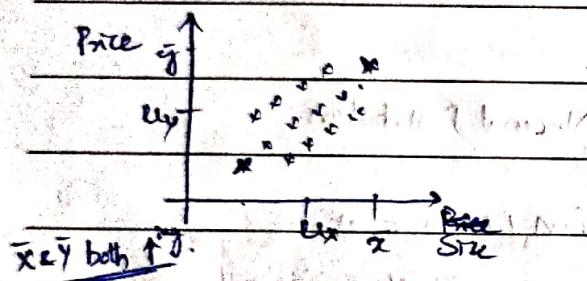
Relation b/w Size & Price

$S \uparrow \rightarrow P \uparrow$        $S \downarrow \rightarrow P \downarrow$

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Cov}(X, X) = \text{Var}(X)$$

$X \uparrow \& Y \uparrow \Rightarrow \text{Cov is +ve}$   
 $X \uparrow \& Y \downarrow \Rightarrow \text{Cov is -ve}$



Here Covariance doesn't tell about how much covariance b/w 2 quantities.

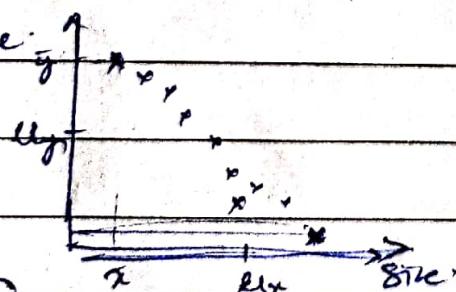
That's why we use Pearson Correlation Coef.

$\bar{x} > \bar{y}$

$\text{Cor.} = \text{+ve. +ve} = \text{+ve.}$

$\bar{x} < \bar{y}$

$\text{Cor.} = \text{-ve -ve} = \text{-ve.}$



$x \uparrow y \downarrow$   
Decreasing

$\bar{x} > \bar{y}$        $\text{Cor.} = \text{+ve. -ve} = \text{-ve.}$

$\bar{x} < \bar{y}$        $\text{Cor.} = \text{-ve. +ve} = \text{-ve.}$

We touch your electricity everyday!

~~Population & Sample~~

~~Population Mean~~

## VI. Central Limit Theorem:-

$$\bar{X} \xrightarrow{\text{(R.V.)}} \text{G.D. } (\mu, \sigma^2) \quad n \geq 30$$

This R.V. may/maynot belong to ~~Gaussian Distribution~~

Taking sample from R.V.  $X$  (30 data points. Randomly Selected)

$$\begin{aligned}
 u &= \sigma^2 \\
 S_1 &= x_1, x_2, \dots, x_{30} = \underline{\bar{x}_1} \\
 S_2 &= \dots = \underline{\bar{x}_2} \quad (\text{some other Random datapts}) \\
 S_3 &= \dots = \underline{\bar{x}_3} \\
 &\vdots \\
 S_{100} &= \dots = \underline{\bar{x}_{100}}
 \end{aligned}$$

According to Theorem ~~Then 30~~  $\bar{X} \approx \text{G.D. } (\mu, \frac{\sigma^2}{n}) \quad n \geq 30$

Mean follows G.D.

## VII. Chebyshchev Inequality

For suppose R.V.  $X$  following G.D.

$$X \approx \text{G.D. } (\mu, \sigma)$$

$$\Pr(\mu - \sigma \leq X \leq \mu + \sigma) \approx 68\%$$

$$\Pr(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%$$

$$\Pr(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99.7\%$$

For suppose R.V.  $Y$  not following G.D.

$$Y \notin \text{G.D. } (\mu, \sigma)$$

What % of # datapoints in  $Y$  falling within range of first Std. deviation?

Chebyshchev Inequality.

$$\Pr(\mu - k\sigma < X < \mu + k\sigma) \geq 1 - \frac{1}{k^2}$$

$k=2$  2nd Std. deviation.

$$\Pr(\mu - 2\sigma < X < \mu + 2\sigma) \geq 1 - \frac{1}{4} \Rightarrow \Pr(\mu - 2\sigma < X < \mu + 2\sigma) \geq \frac{3}{4}$$

$$\Rightarrow \Pr(\mu - 2\sigma < X < \mu + 2\sigma) \geq 75\%$$

$k=3$  3rd Std. deviation

$$\Pr(\mu - 3\sigma < X < \mu + 3\sigma) \geq 1 - \frac{1}{9} \Rightarrow \frac{8}{9} \times 100 \geq 88.9\%$$

We touch your electricity everyday!

## Pearson Correlation Coef.

Useful for feature selection

$$\text{Covariance} = \text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\begin{array}{ll} x_1 y_1 = +ve \\ x_1 y_2 = -ve \end{array}$$

Here Covariance gives Only direction of Relationship

but doesn't know about strength.

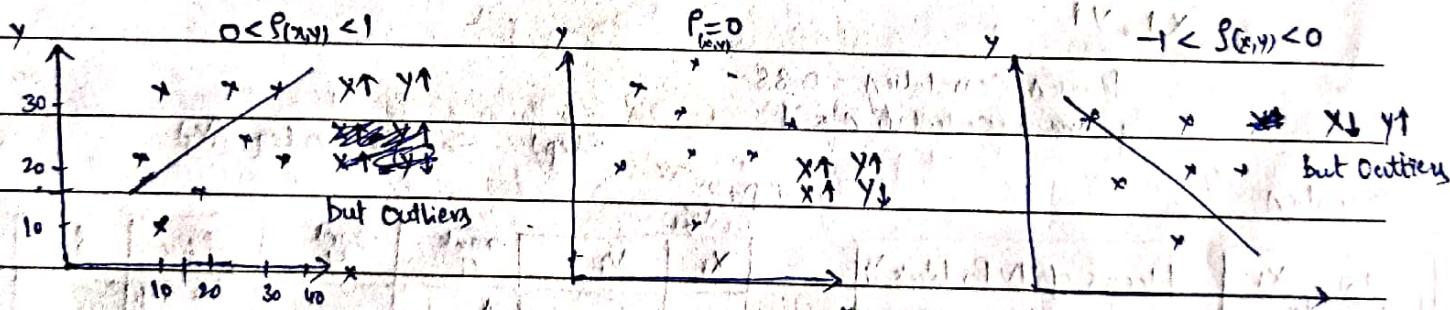
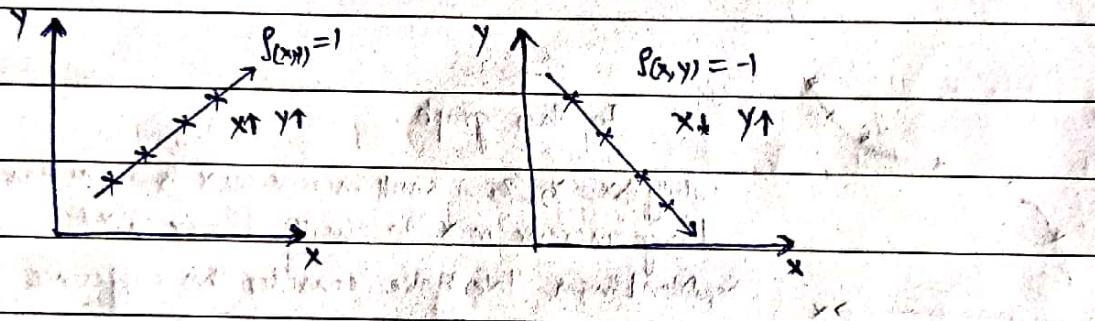


$$\begin{array}{ll} x & y \\ \text{height} & \text{weight} \\ x \uparrow \Rightarrow y \uparrow & \end{array}$$

$$\text{Pearson C.C.} = P(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

We get How much positivity/Negativity is Strength

$$\text{Value Range} : P(x, y) = -1 \leq P \leq 1$$



$$\begin{aligned} x &= 10, y = 10 \\ x &= 10, y = 30 \\ x &= 20, y = 10 \end{aligned}$$

Suppose  $P(x, y) = 1$ , re  $x \uparrow y \uparrow$  we can drop one variable either  $x$  or  $y$  for feature selection in ML as both features are same.

Spearman's Rank Correlation Coef.
[Wikipedia](#)

\* HeatMap Case this Spearman's Correlation.

$$\text{Cov}(r_{gx}, r_{gy})$$

$$* \quad \rho_s = \text{Cov}(r_{gx}, r_{gy}) = \frac{\text{Cov}(r_{gx}, r_{gy})}{\sigma_{rg_x} \cdot \sigma_{rg_y}} \quad \text{Here We Considering Ranks}$$

Rank of x Rank of y

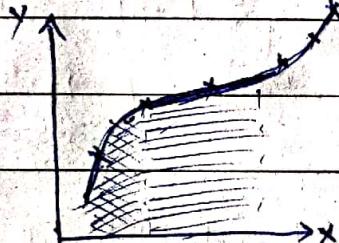
P → Pearson Correlation Coef., but applied to Rank Variables.

$\text{Cov}(r_{gx}, r_{gy}) \rightarrow$  Covariance of Rank Variables.

$\sigma_{rg_x}, \sigma_{rg_y} \rightarrow$  Std. deviations of Rank Variables

$$\rho_s = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$$

$d_i = rg(x_i) - rg(y_i) \Rightarrow$  Diff. b/w 2 Ranks of Each Observation.



In this graph,  
 Initially ~~a small increase in x~~ a small increase in x  $\Rightarrow$  Small increase in y.  
 Large increase in x  $\Rightarrow$  Small increase in y.  
 So NonLinear We have consider this difference also.

Pearson Correlation = 0.88  
 Spearman Correlation = 1

Steps :- Sort the order ( $X_i$ )

Question

Answer:  
Sort

	IQ. $X_i$	Hours of TV Per Week $Y_i$
1	106	7
2	100	27
3	86	2
4	101	50
5	99	28
6	103	29
7	97	20
8	101	50
9	113	12
10	112	6
	110	17

	$X_i$	$Y_i$	rank $x_i$	rank $y_i$	$d_i$	$d_i^2$
1	86	2	1	1	0	0
2	97	20	2	6	-4	16
3	99	28	3	8	-5	25
4	100	27	4	7	-3	9
5	101	50	5	10	-5	25
6	103	29	6	9	-3	9
7	106	7	7	3	4	16
8	110	17	8	5	3	9
9	112	6	9	2	7	49
10	113	12	10	4	6	36

We touch your electricity everyday!

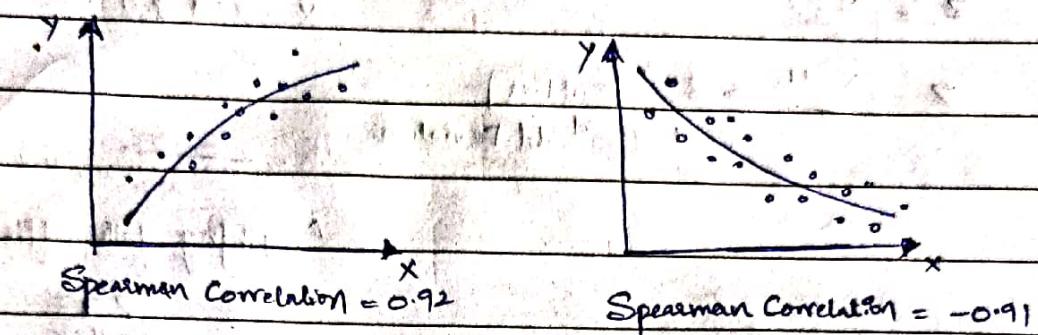
Pearson Correlation =  
 Spearman Correlation =  
 Cov(x,y) → diff. of x & diff. of y  
 Rank(x,y)  
 → diff. of x & y  
 dependent

$$\rho = 1 - \frac{6 \times 194}{10(10-1)} = -\frac{29}{165} = -0.17575757\dots$$

∴ Coef. is very low.

No ~~relation~~ b/w TO & Hours per TV

(p-value 0.627  
t-distribution)



### Outliers



Spearman Correlation = 0.84

Pearson Correlation = 0.67

Pearson Correlation focus on Linear

but Spearman helps on Nonlinear also.

### \* Finding Outliers in Dataset using Z-score and IQR.

#### Outlier:

An Outlier is a data point in a dataset that is distant from all other observations.

A data point that lies outside the overall distribution of the dataset.

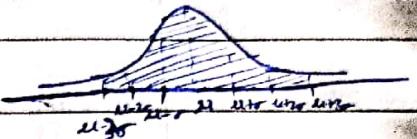
#### Z-score

$$Z = \frac{x - \mu}{\sigma}$$

$$x = \mu + Z\sigma$$

$$Z = \frac{x - \mu}{\sigma} > 3 \quad (\text{outlier})$$

Out of 3rd Std. Deviation.



Anything falls after  $\mu + 3\sigma$  is an Outlier

$Z \geq 3$   
Outlier.

### I.Q.R Interquartile Range - Est

$$\text{dataset} = [5, 6, 7, 1, 2, 8, 10, 3, 4, 9]$$

for suppose

1st. quartile 37.5% 50% 60% 75.5% 90% 100% 2

Here Calculation  
Wrong,

Ans Sort the dataset = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

1 → ~~These are~~ 1 → 0 numbers are less than 1 → 0%.

2 → 10% of Total Numbers are less than 2 ie. 1

$$(\text{before } 2 = 1 \text{ Number only}) \quad \text{P. } \frac{1 \text{ Number}}{\text{Total Numbers}} \times 100 \% = \frac{1}{10} \times 100 \\ = 10\%$$

10%

\*\*\* 3 →

1, 2, 3, 4...10

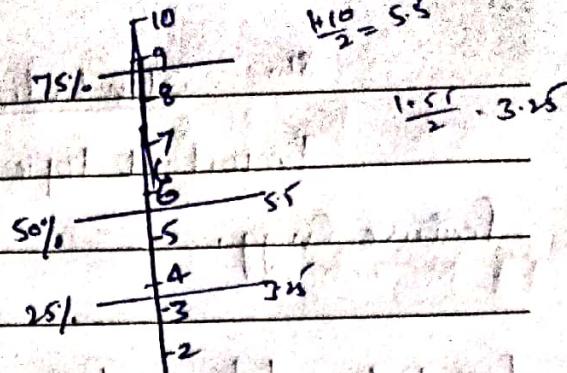
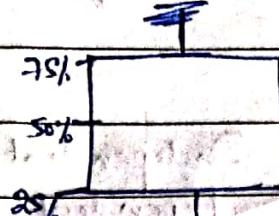
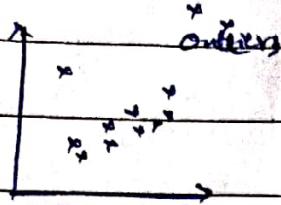
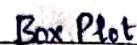
2 Numbers are less than 3. ie.  $\frac{2}{10} \times 100 = 20\%$

4 → 30%

In InterQuartile Range, Most Numbers are Concentrated at 25% - 75%

$$\text{I.Q.R} = \underline{75\%} - \underline{25\%} \text{ values}$$

We touch your electricity everyday!



Below 25% & Above 75%

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33

Sorted the List -

Total Elements = 34

Percentile:

~~Percentile:~~ To 1st Element '0' Index: '0' -  $\frac{0}{34} \times 100\% = 0\%$

$$2^{\text{nd}} \text{ Element } 1 \text{ Index: '10'} - \frac{1}{34} \times 100\% = 2.9\%$$

6<sup>th</sup> element 5 index 'ii' -  $\frac{5}{34} \times 100\% = 14.7\%$ . (5 elements are there before 11)

$$\text{lower } \frac{x}{34} \times 100 = 25\%$$

Oppos

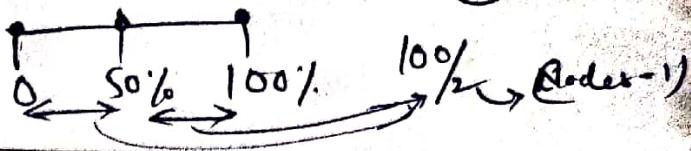
$$\frac{24}{32} \times 100 = 75\%$$

$$x = 8.5$$

~~Number 9~~ ~~Card 10 from 12~~

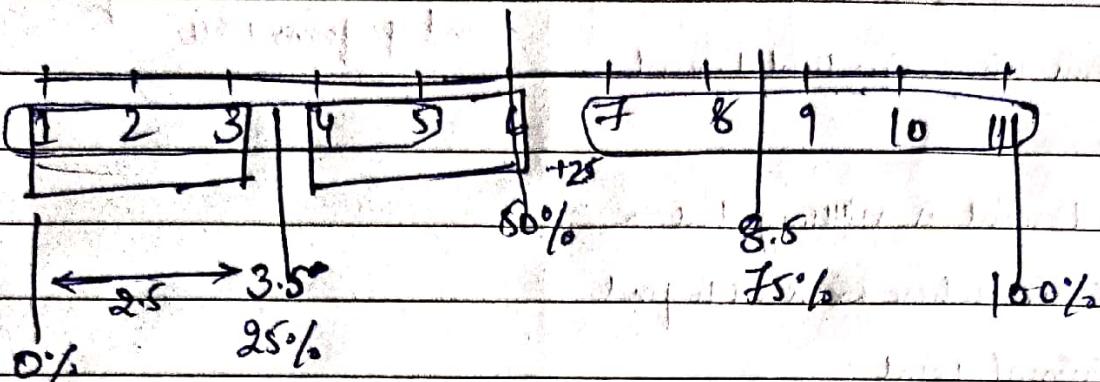
$$y = 25.5$$

We touch your electricity everyday!



## TQR Practice

11 Numbers:



~~(Total No. Count + 1) / 100%~~

~~$25\% \cdot (11+1) 25\% = (11+1) \times \frac{1}{11+1} = 3$~~

11 nodes = Total No. Count

Each is  $\frac{100}{11-1} = 10\%$   
Segment

$$1 \rightarrow 0\%$$

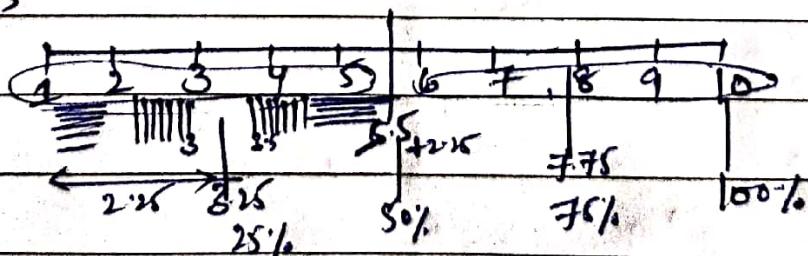
$$2 \rightarrow \frac{1}{11} \times 100\% = 10\%$$

$$3 \rightarrow \frac{2}{11} \times 100\% = 20\%$$

before B Node 3  $\rightarrow$  There are 2 Numbers

$$x = 25\% \times 10\% = 2.5\% \quad \text{Each } 25\% \text{ is.}$$

10 Numbers:



$$\therefore 25\% (10-1) = \frac{9}{4} = 2.25$$

$$\phi\% \text{ Value} = \phi\% \cdot \frac{1}{(\text{Total Count}-1)}$$

We touch your electricity everyday!

## Interquartile range

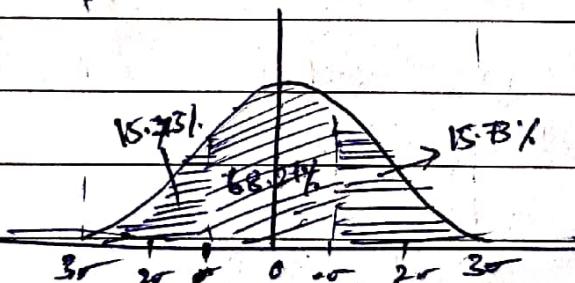
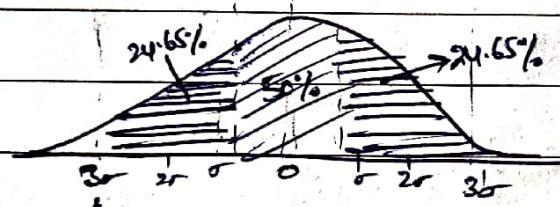
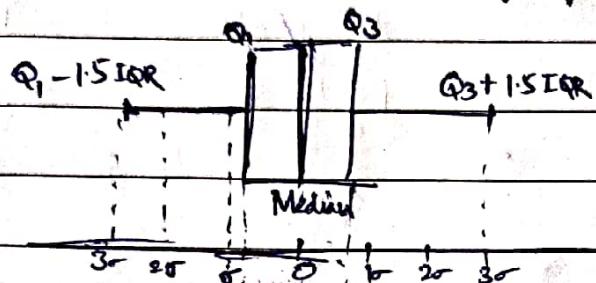
H-Spread,  $Q_3 - Q_1$ , middle 50%.

Measure of Statistical dispersion, diff. b/w 75<sup>th</sup> and 25<sup>th</sup> percentiles

$$IQR = Q_3 - Q_1$$

Here Median is Corresponding Measure of Central Tendency

IQR - Used for Outliers Identifying



Order  $7, 7, 31, 31, 47, 75, 87, 115, 116, 119, 119, 155, 177$ .

$\downarrow$   $\downarrow$   $\downarrow$   
31       $Q_1$       119  
Median of First Half      (Median of Whole Table)

$\downarrow$  Avg  $\downarrow$   
119       $Q_3$   
Median of Last Half.      Frontiers  
 $-101 \quad Q_1 - 1.5 \text{ IQR}$   
 $251 \quad Q_3 + 1.5 \text{ IQR}$

$$IQR = Q_3 - Q_1 = 119 - 31 = 88.$$

$$\begin{aligned} \text{Range (L.B, U.B)} &= (Q_1 - 1.5 \text{ IQR}, Q_3 + 1.5 \text{ IQR}) \\ &= (-101, 251) \end{aligned}$$

We touch your electricity everyday!

## II. Standardization Vs Normalization

Important topic for Feature Scaling which is an integral part of feature Engg.

For Data, We analyse features

Independent  
Dependent

\* Using Independent Variables, We analyse the Dependent Variables.

Features → We Consider ~~&~~ Magnitude & Units.

Age	Weight (kg)	Height (cm)
25	60	1620
Magnitude	No. of hrs	Cms

Two scaling Techniques are Standardization & Normalization.

(i) Normalization :- (min-max Normalization)

Helps to scale down the feature b/w 0 to 1.

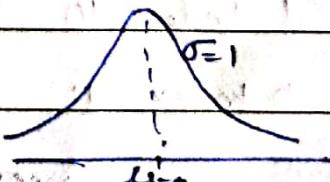
$$x_{\text{norm}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

(ii) Standardization :- (Z-score Normalization)

Helps to scale down the feature based on Std. Normal Distribution.

Mean  $\mu = 0$   
with Std Deviation  $\sigma = 1$

$$z = \frac{x - \mu}{\sigma}$$



## Normalization (Min Max)

\* M.L. Algorithms which involves Euclidean distance.

\* D.L. Algorithms where Gradient Descent is involved.

→ Gradient Descent is nothing but a Parabola Curve where we need to find Global Minimal point.

\* We use Min-Max Normalization

\* Algorithms like KNN

K-nearest Neighbour

K-means Clustering

ANN Artificial Neural Net  
CNN Convolutional Neural Net } Scale down

Linear Regression

0 to 1.  
Images [0 to 255]

Logistic Regression.

## Standardization (StandardScalar)

\* Many of the use cases

Standardization Technique is used and performs well.

\* Some Algorithms, we won't perform Scaling: like Decision Tree, Random Forest, XGBoost and all boosting techniques.

No use of Scaling Techniques here.

ANN - Tensor Flow, Keras require inputs

b/w 0 to 1 to analyse weights quickly.