

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. Measures of Dispersion.

9. Variance

10. Std. Deviation.

11. Gaussian / Normal Distribution

25. Q-Q Plot

26. Chebychev's Inequality

27. Discrete & Continuous Distributions

28. Bernoulli & Binomial Distribution

29. Log Normal Distribution

30. Power Law Distribution

31. Box Cox Transform.

32. Poisson Distribution

33. Application of Non-Gaussian Distribution

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

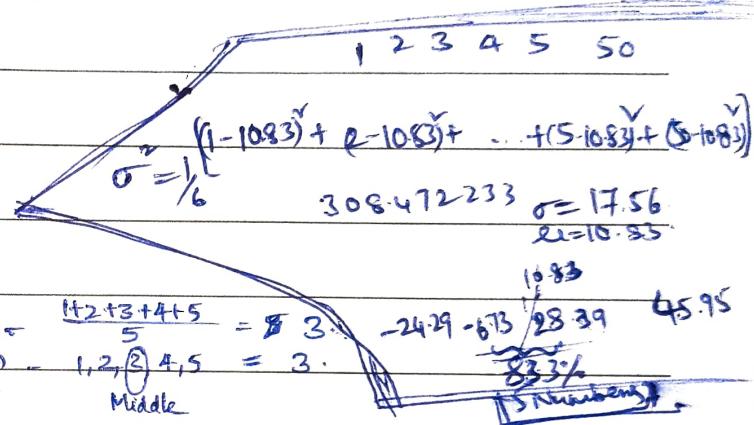
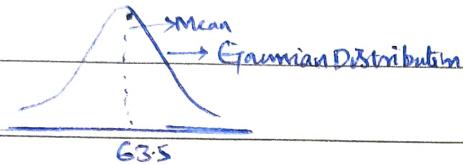
I. Mean, Median, Mode - Measures of

(i) Mean: Measure of Central Tendency.

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

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

$$\bar{x} = \frac{168 + \dots + 175}{7} = 17.5$$



(ii) Median:

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

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

Write in Ascending Order to find Median.

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

Size: Odd No - Middle.

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

Age

23
24

27

32

35

21

= ? Data Mining.

If I remove the missing Data. 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 -

I. Population :-

Question : Average Height of all people in a state

Solution :-

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

$N \rightarrow$ Population.

$n \rightarrow$ Sample.

$\mu \rightarrow$ Population Mean

$\bar{x} \rightarrow$ Sample Mean

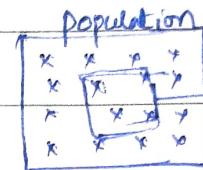
2. Sample:-

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

So, We consider Sample.

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

$n < N$



Sample

When Sample Size increases $n \approx N$

II.

3.1 Random Variables :- A random Variable which Stores Data. Types

Numerical R.V.

Categorical R.V.

Types of Numerical R.V. :-

(i) Discrete Variables - In whole Number.

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

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.: Height of a Person - 63Kg.

Interest Rate

63.2 Kg. (depends on scaling)
63.269 Kg.

Salary.

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
Numerical C.R.V. D.R.V.	Age group	Numerical D.R.V.	No. of People in House	Height of House Owner Numerical C.R.V.

Distributions

1. Gaussian / Normal Distribution

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

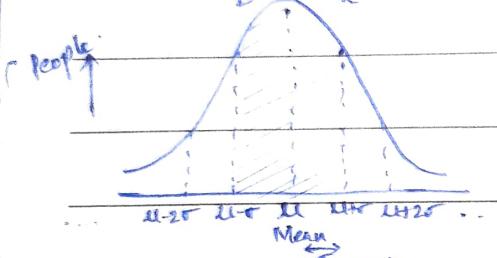
Mean S.D.
 $x \sim GD(\mu, \sigma)$ or $N(\mu, \sigma)$

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

Prob Density Fun / Histogram
Bell Curve

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

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



$$x = x_0$$

→ Height

Empirical formula :-

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

(Probability of Having Data points for $\mu \pm 1\sigma$ (in first deviation))
 $\text{Pr}(\mu - 2\sigma \leq x \leq \mu + 2\sigma) = 95\%$
 $\text{Pr}(\mu - 3\sigma \leq x \leq \mu + 3\sigma) = 99.7\%$

Bell Curve & Symmetrical.

Dataset

- eg:- (i) Distribution of Height
(ii) IRIS Dataset - Petal length.

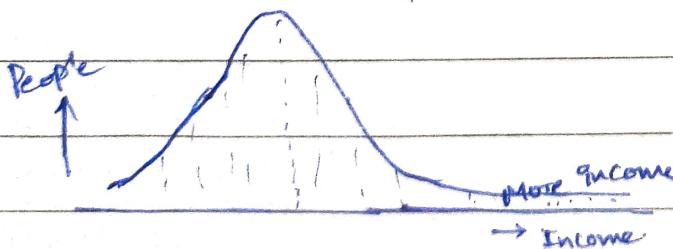
$$[Ex: N.D \Rightarrow \mu=0 \sigma=1]$$

2. Log Normal Distribution

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

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

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



→ Similar to Ex.D, but Right Hand side goes more flatten.
→ Right Skewed.

Dataset eg:- (i) Income of the people.
(ii) Amazon Product Review Description length

More Income people are very less

Why can we use all types of this distributions?

	R&D	Marketing	Profit	States	Company
1	Tobool-	100000/-	-	-	-
2	-	-	-	-	-
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~~ to Std. Normal Distribution ($\mu=0, \sigma=1$)

$$\frac{x_i - \mu}{\sigma}$$

Marketing \rightarrow

$$\begin{aligned} \text{Marketing} &\rightarrow \ln(\text{marketing}) \\ \text{Marketing} &\rightarrow \ln(\text{marketing}) \end{aligned}$$

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

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

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

It is called Log Normalization.

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

~~Confidence:~~

Covariance :- Useful in Data Preprocessing.

Size Price

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

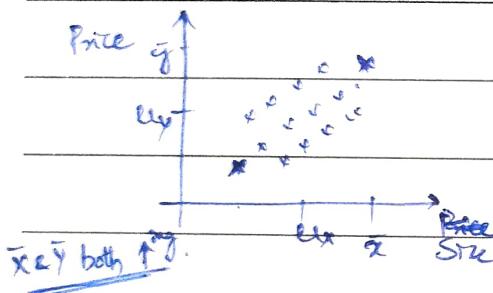
Relation b/w Size & Price

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

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

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

$x \uparrow \& y \uparrow \Rightarrow \text{Cov is +ve}$
 $x \uparrow \& y \downarrow \Rightarrow \text{Cov is -ve.}$

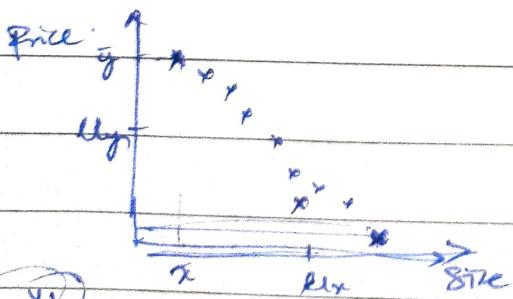


Here Covariance doesn't tell about how much Covariance b/w 2 Quantities.

That's why we use Pearson Correlation Coef.

$$\begin{aligned} \bar{x} > \bar{x} \\ \bar{y} > \bar{y} \end{aligned} \quad \text{Cor.} = \text{+ve, +ve} = \text{+ve.}$$

$$\begin{aligned} \bar{x} < \bar{x} \\ \bar{y} < \bar{y} \end{aligned} \quad \text{Cor.} = \text{-ve, -ve} = \text{-ve.}$$



$x \uparrow \& y \downarrow$
Vice versa.

$$\begin{aligned} \bar{x} > \bar{x} \\ \bar{y} < \bar{y} \end{aligned} \quad \text{Cor.} = \text{+ve, -ve} = \text{-ve.}$$

$$\begin{aligned} \bar{x} < \bar{x} \\ \bar{y} > \bar{y} \end{aligned} \quad \text{Cor.} = \text{-ve, +ve} = \text{-ve.}$$

IV. Central Limit Theorem:-

$$X \not\sim G.D(u, \sigma^2) \quad n \geq 30$$

(R.V.)

This R.V. may/may not belong to Gaussian Distribution

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

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

~~\bar{x}_{100}~~

According to Theorem ~~Mean~~ $\bar{x} \approx G.D(u, \frac{\sigma^2}{n})$ $n \geq 30$.

Mean follows G.D.

Chebyshew's Inequality

For suppose R.V. X following G.D

$$X \approx G.D(u, \sigma)$$

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

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

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

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

$$Y \not\sim G.D(u, \sigma)$$

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

Chebyshew's Inequality

$$\Pr(u - K\sigma \leq X \leq u + K\sigma) \geq 1 - \frac{1}{K^2}$$

$K=2$ 2nd Std deviation.

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

$$\Rightarrow \Pr(u - 2\sigma \leq X \leq u + 2\sigma) \geq 75\%$$

$K=3$ 3rd Std deviation

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

Pearson Correlation Coef.

Useful for feature selection

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

$$\begin{array}{ll} x \uparrow y \uparrow & = + \\ x \uparrow y \downarrow & = - \end{array}$$

Here Covariance gives Only direction of Relationship.

but doesn't know about Strength.

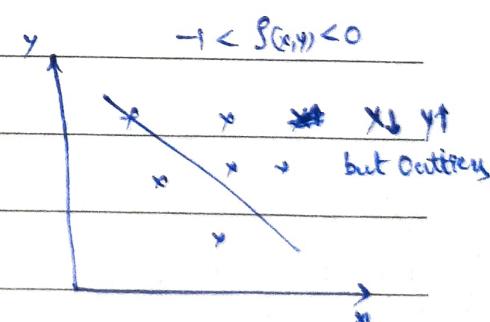
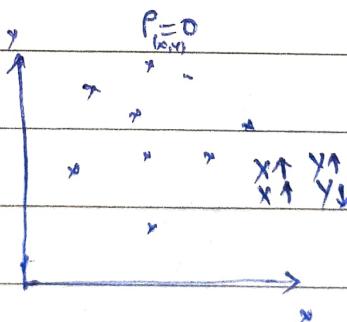
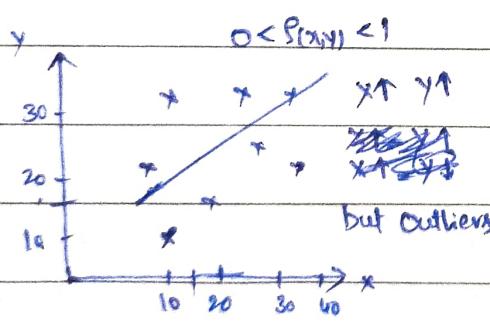
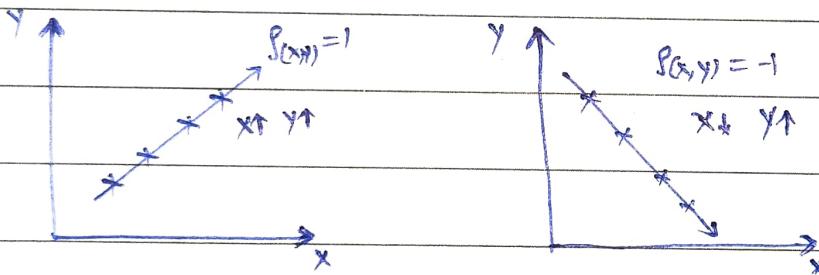


$$\text{Pearson C.C.} = P(x,y) = \frac{\text{cov}(x,y)}{\sqrt{x} \cdot \sqrt{y}}$$

Value Range - $P(x,y) \in [-1 \leq P \leq 1]$

We get How much positivity/Negativity is Strength

x y
Height Height
 $x \uparrow \Rightarrow y \uparrow$



$$\begin{aligned} x &= 10, 20, 30, 40 \\ x &= 10, 20, 30 \\ z &= 10, 20, 30 \end{aligned}$$

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

B.

Spearman's Rank Correlation Coef.

Wikipedia

* Heat Map use this Spearman's Correlation.

$$\rho_s = \text{Cov}(r_{gx}, r_{gy}) = \frac{\text{cov}(rg_x, rg_y)}{\text{Rank}_x \times \text{Rank}_y}$$

Here we Considering Ranks.

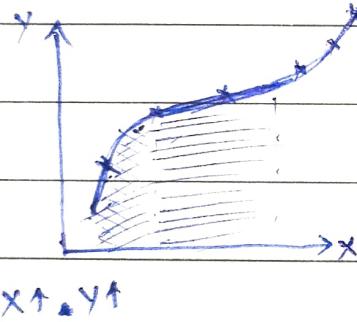
P → Pearson Correlation Coef., but applied to Rank Variables.

$\text{cov}(rg_x, rg_y) \rightarrow$ Covariance of Rank Variables

$\text{Org}_x, \text{Org}_y \rightarrow$ 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~~ \Rightarrow Small increase in Y.
Large increase in X \Rightarrow Large increase in Y.

So NonLinear We have consider this difference also.

X↑, Y↑

Pearson Correlation = 0.88

Spearman correlation = 1

Question

Answer
Sort it

Step 1 :- Sort the Order (X_i)

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

X _i	Y _i	rank X _i	rank Y _i	d _i <small>rank(X_i) - rank(Y_i)</small>	d _i ²
86	2	1	1	0	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	29	6	9	-3	9
106	7	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

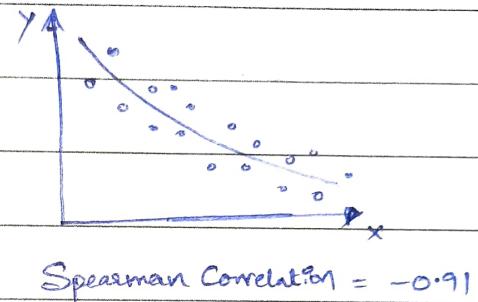
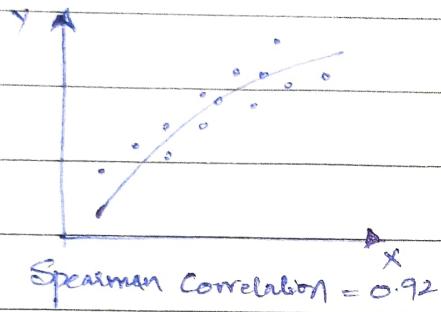
$\left\{ \begin{array}{l} \text{Pearson Correlation:} \\ \text{Spearman Correlation:} \end{array} \right.$
 Cov(x,y) \rightarrow diff. of x & diff. of y
 Rank(x,y) \rightarrow diff. of x & y
 independent 9
 dependent

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

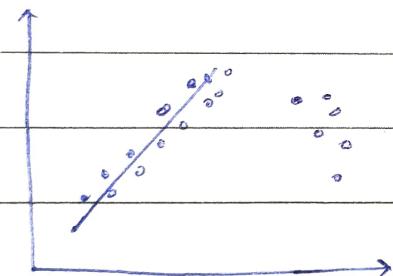
P Coef. is very low.

No ~~relation~~ b/w IQ & 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.

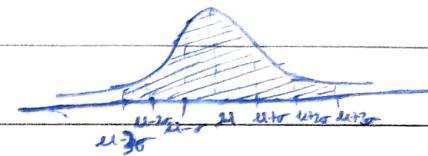
Z-score:

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

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

$$Z = \frac{x - \mu}{\sigma} > 3 \quad [\text{Outlier}]$$

Out of 3rd Std Deviation.



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

Z > 3.
Outlier.

I.Q.R InterQuartile Range - Eat

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

for ~~Suppose~~

1st Quartile: 3/10th 50% 6/10th 7/10th 90% 100%.

Here Calculation Wrong.

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

1 → ~~1 Number~~ 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. $2/10 \times 100 = 20\%$

4 → 30%

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

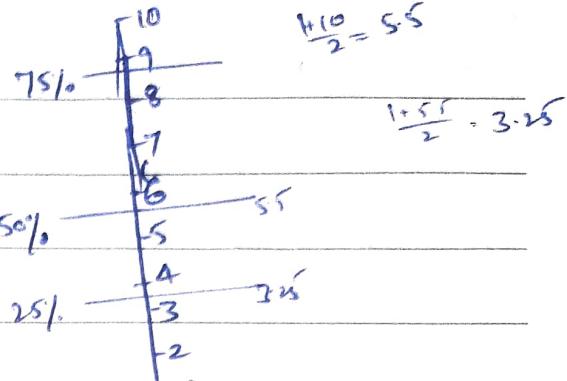
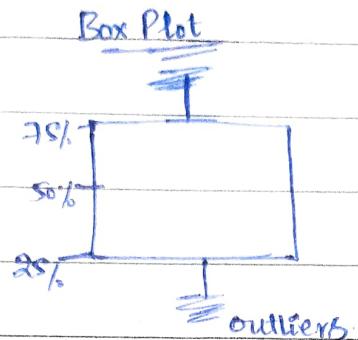
$$\text{I.Q.R} = \underline{75\%} - \underline{25\%} \text{ value}$$

(11)

Scatter Plot



Box Plot



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

~~10, 10, 10, 10, 10, 11, 11, 12, 12, 12, 12, 12, 12, 13, 13, 13, 13, 13,~~
~~14, 14, 14, 14, 14, 14, 14, 15, 15, 15, 15, 15, 17, 17, 19, 102, 107, 108~~
32

Sorted the List -

~~0% 2.9% 5.8%~~
~~8.3% 11.1% 14.7% 17.6% 20.5% 23.5% 26.4% 29.4% 32.5% 35.2% 38.2% 41.1% 44% 47% 50%~~
~~S = [10, 10, 10, 10, 10, 11, 11, 12, 12, 12, 12, 12, 12, 13, 13, 13, 13, 13,~~
~~14, 14, 14, 14, 14, 14, 14, 15, 15, 15, 15, 15, 17, 17, 19, 102, 107, 108]~~
~~18 9 20 21 22 23 24 25 26 27 28 29 30 31 32 33~~
~~52.9 55 58.8 61.1 64.7 67.6 70.5 73.5 76.4 79 82 85 88 91 94 97~~

Total Elements = 34

Percentile:

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

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

5th Element 5 Index '5' - $\frac{5}{34} \times 100\% = 14.7\%$ (5 elements are before 5)

Upper

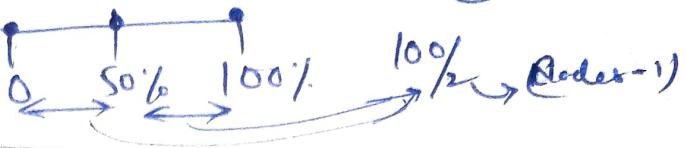
$$\frac{27}{34} \times 100 = 79\%$$

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

$$y = 25.5$$

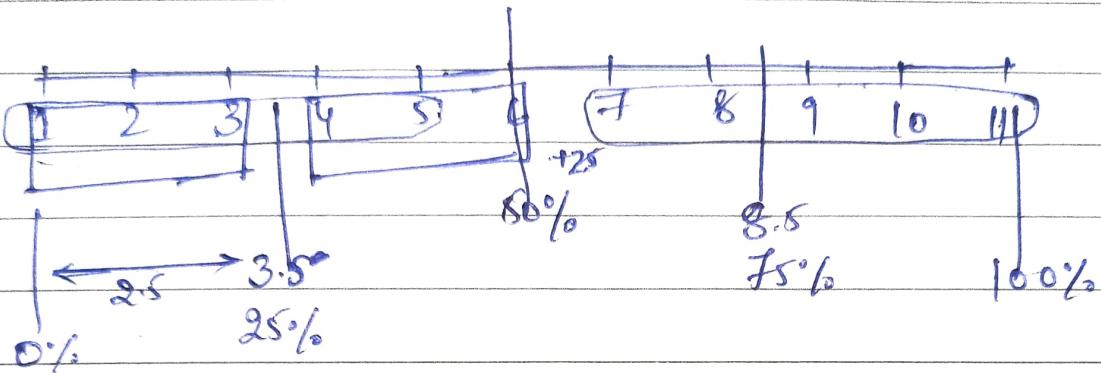
$x = 8.5$
 Number of elements from 1 to 9
 25
 24

(12)



IQR Practice

11 Numbers:



$$25\% \times \frac{(\text{Total No. Count} + 1)}{(1+1)25\%} = \frac{(11+1)}{4} \times \frac{1}{25\%} = 3 \text{ Per } \% \text{ Value.}$$

~~100% = 25%.~~

11 nodes = Total No. Count

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

$$1 \rightarrow 0\%$$

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

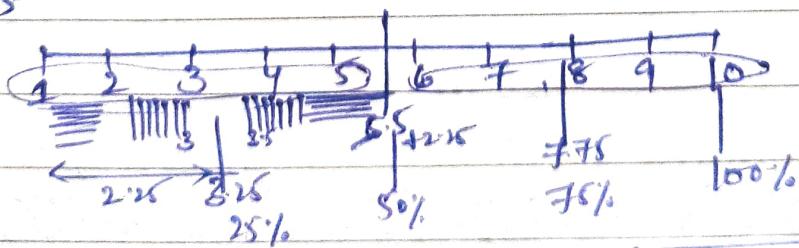
$$3 \rightarrow \frac{2}{11-1} \times 100\% = 20\%$$

before B Node 3 \rightarrow There are 2 Numbers

$$x = \frac{25\% \times 10}{100\%} = 2.5 \quad \text{Each } 25\% \text{ is.}$$

~~25% Value = 25% of total~~

10 Numbers:



$$\therefore 25\% (10-1) = 9/4 = 2.25$$

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

Interquartile range

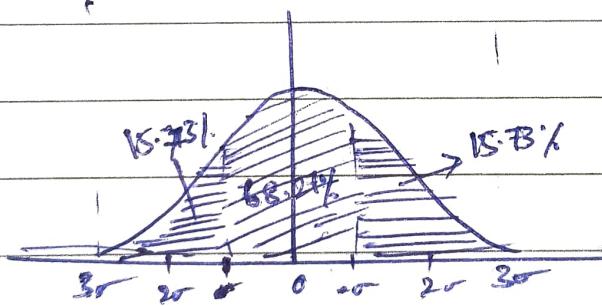
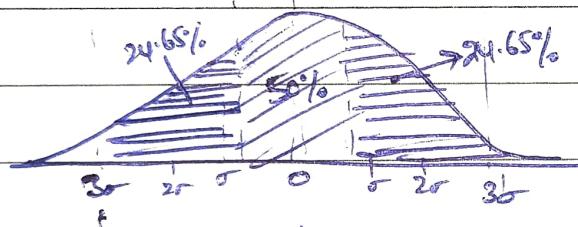
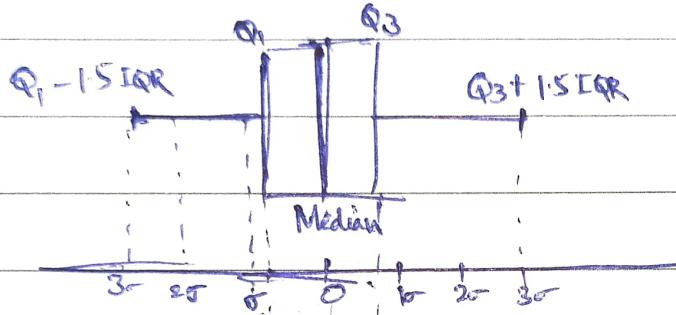
H-Spread, Q-Q, middle 50%

Measure of Statistical dispersion, diff b/w 75th and 25th percentiles

$$IQR = Q_3 - Q_1$$

Here Median is Corresponding Measure of Central Tendency

IQR - Used for Outliers Identifying



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

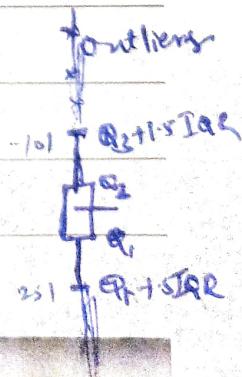
$\xrightarrow{\text{Avg}}$
 $31 \downarrow Q_1$
 Median of Upper Half
 (Median of Lower Table)

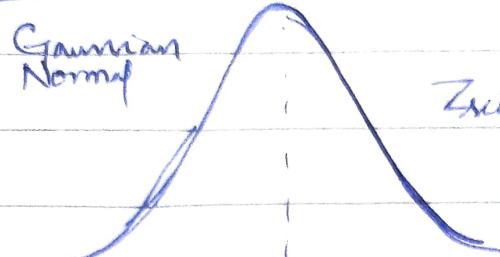
$\downarrow \text{Avg}$
 $119 \downarrow Q_3$

Median of Last Half

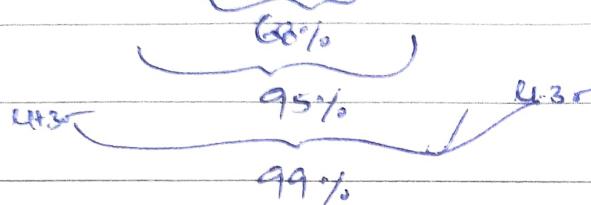
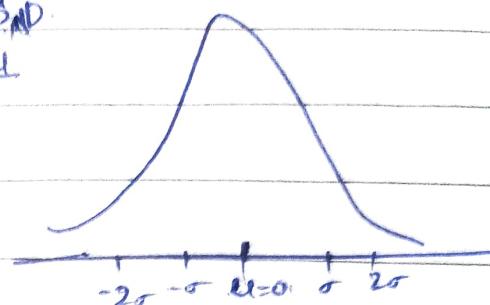
$$IQR = Q_3 - Q_1 = 119 - 31 = 88$$

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



Z-score :-

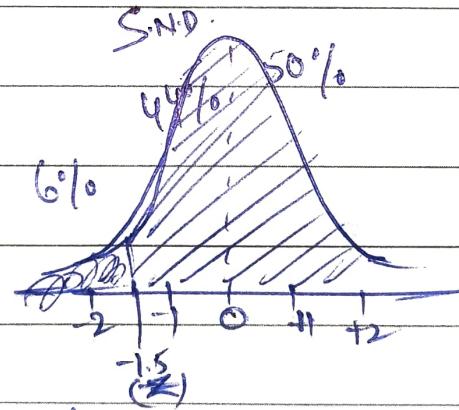
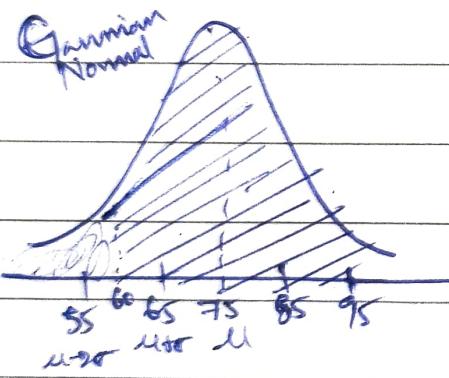
$$Z_{\text{score}} = \frac{x_i - \mu}{\sigma}$$



$$\mu = 75$$

$$\sigma = 10$$

$$P(x > 60)$$



Z-score Table gives ^{Area of} left hand side (Φ) ~~1500~~
-1.500

$$= -0.06681 \text{ i.e. } 6.681\%$$

Remaining Area 93.319%
 (Probability) $\leq 94\%$

$$\frac{100.689}{93.319}$$

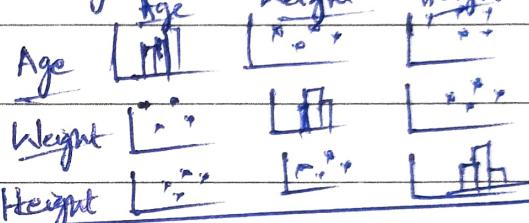
* Univariate, Bivariate & Multivariate Analysis

Height cm	Weight kg	O/p
180	90	Obesity
160	50	Slim
170	78	Fit
190	90	Fit
175	85	Slim.
:	:	

Multivariate

More than 3 features - Age, Weight, Height

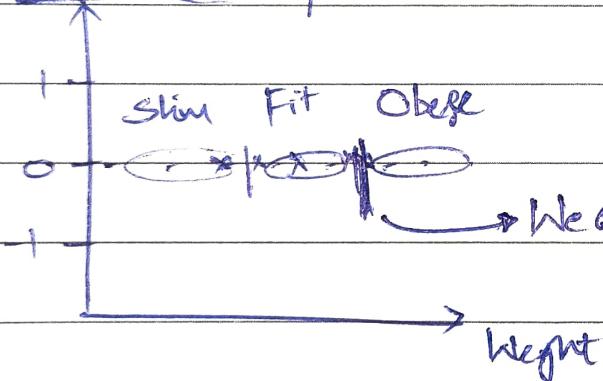
We go with Pairplot in Seaborn.



Gives concept of Correlation.

A

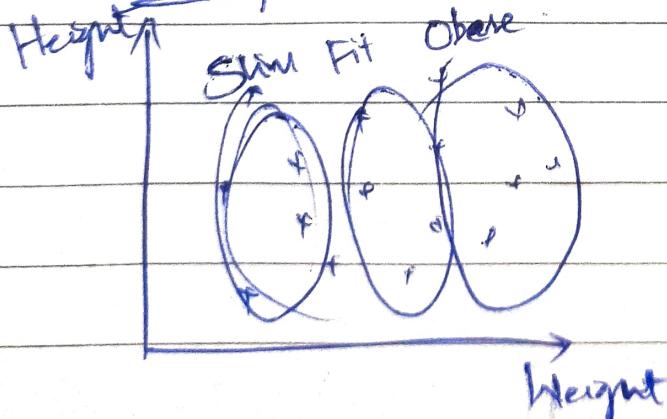
Univariate Analysis



We ~~are~~ maybe 100% Correct. Overlap point will be there

We go for Bivariate/Multivariate

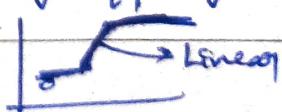
Bivariate Analysis



Later

If we are able to clarify point property, we can go for Logistic Regression

Logistic Regression :-



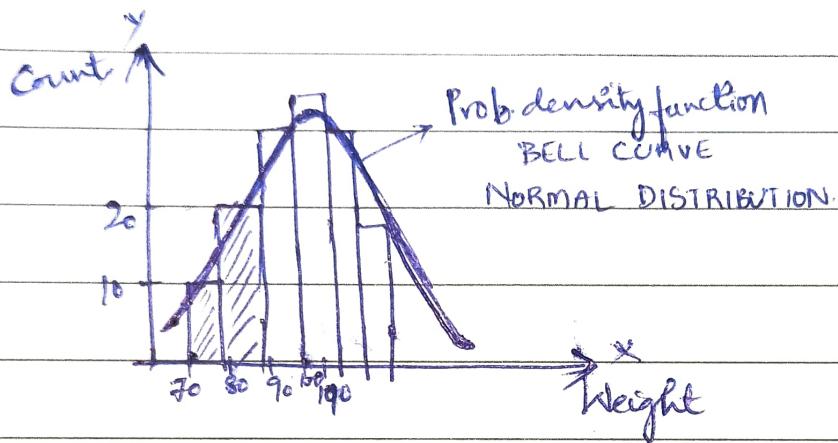
If Overlapping Occurs, we go with Decision Tree / XG Boost / Random forest (Nonlinear) SVM / AdaBoost / kNN Algo.

Histograms:-

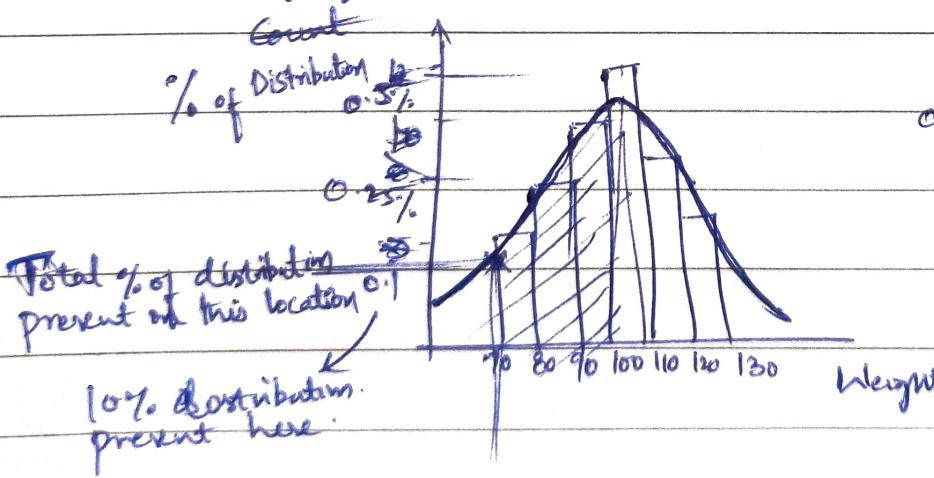
The given range

Gives Information of How many No. of points in a Given range
(box)

i.e. \Rightarrow Builds Bar graph.

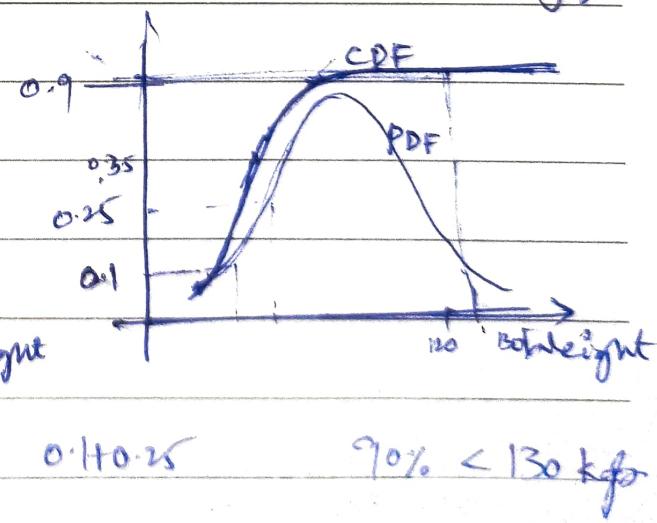


Prob. density function PDF



CDF

Cumulative Density function



$$0.1 + 0.25$$

$$70\% < 130 \text{ kgs}$$

E. 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)
Units (No. of Years)	25	60	162
Magnitude	25	No. of Years	Kgs

Two scaling Techniques are Standardization & Normalization.

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

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

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

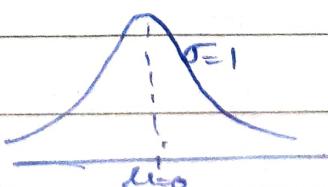
(ii) Standardization :- (Z-score Normalization)

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

With Mean $\mu = 0$
Std. Deviation $\sigma = 1$

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

~~Answe~~ is



Normalization (Min Max)

(StandardScalar)

- * M.L. Algorithm 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-means Neighbour
K-means Clustering

ANN Artificial Neural Net
CNN Convolutional Neural Net } Scale down
Linear Regression Images (0 to 255)
Logistic Regression.

* 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 Input

b/w 0 to 1 to analyse heights quickly.