

Statistics

Page: _____

- (1) Statistics is the discipline that concerns the collection, organization, analysis, interpretation, and presentation of data.

Types of Statistics

- (1) Descriptive stats
(2) Inferential stats

The data can be from population or sample

Descriptive stats

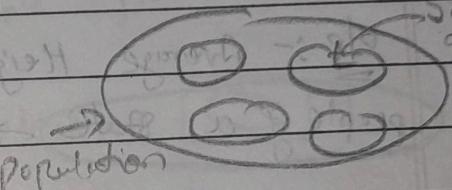
Data "facts or pieces of information"

e.g. Age = 20, height = 5.6 ft

Inferential stats

- (1) Analyzing data, summarizing data. Organizing data in the form of number & graph.

(1)



- (2) BAR plot, Histogram, Pie chart, PDF, CDF, Normal Distribution. From the entire population we are just taking a sample and we are inferring some information.

- (3) Measure of Central Tendency.
Mean, median, mode.

e.g.

- (1) Z-Test \rightarrow Hypothesis testing
• T-Test

- Chi square Test

- (4) Measure of Variance
SD, Variance

Definition:

- e.g. There are 20 classes at a university and you have collected the ages of students in one class

Ages = {20, 21, 18, 25, 26, 24, 22, 21}

What is the most common age of student in your class?

Ans = {21} mode

consists of using data you have measured to form conclusions / inferences

* Population Vs Sample

* Statistics

(u) (1) Population → Population Mean? or other are

(x) (2) Sample → Sample Mean (mean, median, mode)

③ Random Variable

Discrete Random Variable

Continuous Random Variable

* Population

e/p :- Average Height of all the people of a state → Karnataka

↓ considering

$\mu = \text{mean}$

x	x	x
x	x	x
x	x	x

Population

⇒ All → Average height of all the people

~~use in~~

$$\text{population } \mu = \frac{1}{M} \sum_{i=1}^M x_i$$

Note :-

When we have large population, it is very difficult to find each and every people's height.

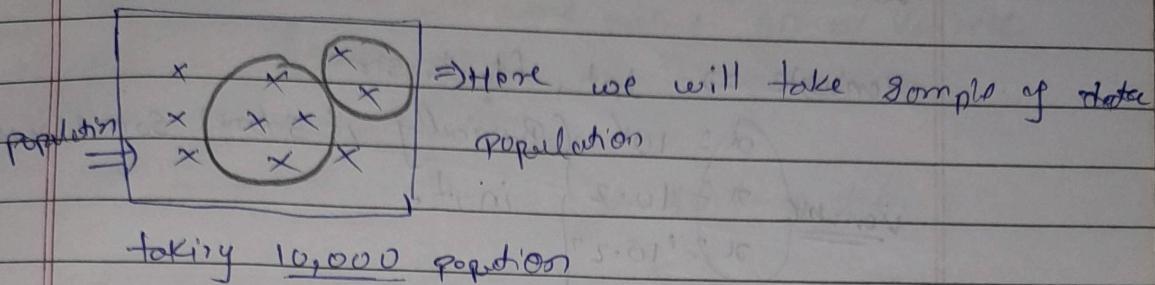
So, we are using Samples from it.

Sample mean is denoted by \bar{x}

Population count $\Rightarrow N$ | Population mean = μ
 Sample count = n | Sample mean = \bar{x}

Date: _____
 Page: _____

Sample



$$\text{Sample } \bar{x} = \frac{1}{n} \sum_{i=1}^{10,000} h_i$$

Sample mean $\frac{1}{n} \sum_{i=1}^{10,000} h_i \approx \text{height of all people}$

Note:- The population are v.v. Large values.

so we try to take sample from population.

Eg:- EXIT Polls (Election O)

Population (N) & Sample (n)

$$\text{Population mean } \mu = \frac{1}{N} \sum_{i=1}^N h_i$$

$$\text{Sample mean } \bar{x} = \frac{1}{n} \sum_{i=1}^n h_i$$

Population mean \gg Sample Mean

Note:- When the sample size increases Sample mean \approx Population mean equal mean

* Gaussian Distribution or Normal Distribution

1) Random Variable

- Discrete Random Variable
- Continuous Random Variable

④ Gaussian / Normal Distribution

Random Variable

variable \Rightarrow int, float, str

variable $\left\{ \begin{array}{l} x = 10 \\ x = 10.2 \\ x = "10.2" \end{array} \right\}$ Basically it stores some value in it.

\rightarrow Discrete \Rightarrow It will be whole number

It will not be floating number

e.g. - No. of Bank Account a person has

$\Rightarrow 2, 3, 4, 5, \dots$

e.g. - population of a state 1 million

\rightarrow continuous \Rightarrow within a range of values we can have any value

e.g. - 10 \dots to \dots 15

it can be $10, 10.2, 10.3, 14.9, 14.8, 15$

It can be whole number, Decimal.

height = 5.11 inch e.g. galaxy, moon orbit
5.6 inch

② Gaussian Distribution / Normal Distribution

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

(3 sigma limit)

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\underline{\underline{\sigma}} \rightarrow \sigma = \sqrt{\text{Var}}$$

Gaussian Distribution Model

mean

$\sigma = \sigma$



random variable

mean

σ

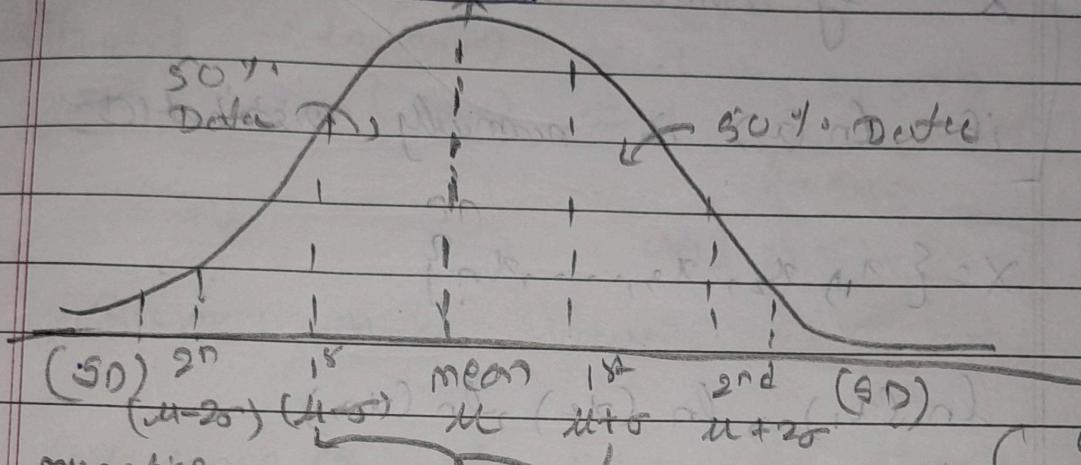
Date: _____
Page: _____

specifically
continuous
variable

$$x \sim G(\mu, \sigma)$$

with same value σ^2 or variance

Bell curve



Three properties

* Empirical formula in Gaussian Dist

68.26% of random variable will be falling in this area

Standard Deviation

$$\text{① } \Pr(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

x which is a part of X

The total of n elements from the X lies between $1^{\text{st}} \sigma$ to the left and $1^{\text{st}} \sigma$ to the right

$$\text{② } \Pr(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$\text{③ } \Pr(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$$

③ * Log Normal Distribution (Refer Line Class)

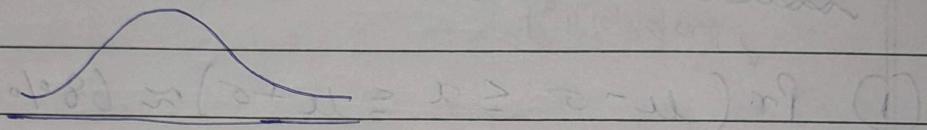
$X \sim \text{log normal distribution}$

If $\ln(x)$ is normally distributed

$$X = \{x_1, x_2, x_3, \dots, x_n\}$$

$$\ln(x_1), \ln(x_2), \ln(x_3)$$

↳ normally distributed



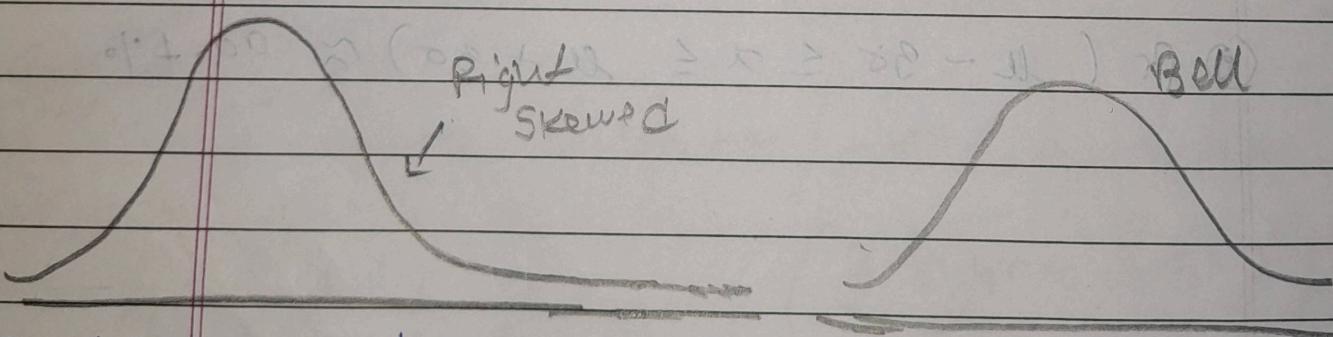
* log normal

$X \sim \text{log normal}$

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

Ex: Income of people

log normal \Rightarrow Gaussian Dist



Eg:- 1) Income of the people

Eg:- Distribution of

2) Amazon product reviews

height

• Sentiment Analysis usually follows
log normal Distribution

• Test score of 100

• Hospitalization Days

Marking

$1500\$ \rightarrow \ln(1500)$ will follow $\approx N(\mu, \sigma)$ \Rightarrow SAD

$2000\$ \rightarrow \ln(2000) \approx N \Rightarrow$ SAD

$3000\$ \rightarrow \ln(3000) \approx N$ then convert \Rightarrow SAD

\therefore Accuracy will be increased.

(*) Covariance Quantify Relation \rightarrow to find out some exact relationship through some numerical number we use concept of Covariance! - It is one of the very imp. topic when we consider Data Preprocessing, Data Analysis.

e.g. Consider two Random variable

<u>Size</u>	<u>Price</u>	Quantify relation
1200 Sqm	100K\$	size \longleftrightarrow Price
1500 Sqm	200K\$	
1800 Sqm	300K\$	{ Quantify relation is basically means finding relationship between size and Price}

\therefore If, $S \uparrow$ $P \uparrow$ (Increase) \therefore Size and Price
 $S \downarrow$ $P \downarrow$ (decrease)

Covariance

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

$\overbrace{x \quad y}$ \uparrow Sample mean of random Variable of x \uparrow random variable of y
 random variable

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

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

$$\therefore \text{cov}(x, x) = \text{Var}(x)$$

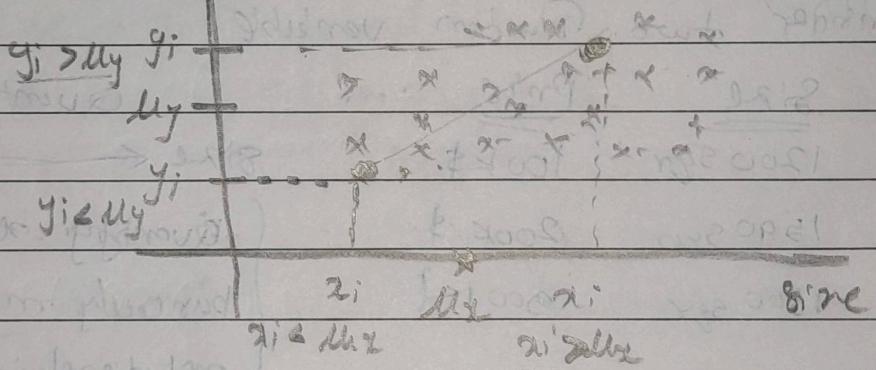
If we are trying to find out cov relation between $x \to x$ it is nothing but variance of x

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

~~$x \uparrow y \downarrow$~~ \Rightarrow $\boxed{+ve}$

~~$x \uparrow y \uparrow$~~ \Rightarrow $\boxed{-ve}$

\uparrow mix

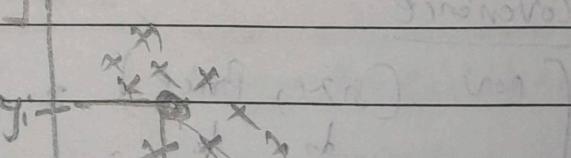


$$\text{cov}(x) = (+ve) \times (+ve) \Rightarrow +ve$$

$$\text{cov}(x) = (-ve) \times (-ve) \Rightarrow +ve$$

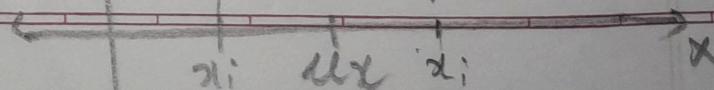
* $x \uparrow y \downarrow$

$y \uparrow$



$$\Rightarrow (+ve) \times (-ve) = -ve$$

$$(-ve) \times (+ve) = -ve$$



In covariance

$x \uparrow y \uparrow \Rightarrow$ +ve \Rightarrow How much +ve (It doesn't tell)

$x \uparrow y \downarrow \Rightarrow$ -ve \Rightarrow How much -ve (It doesn't tell)

To overcome this Disadvantages we use

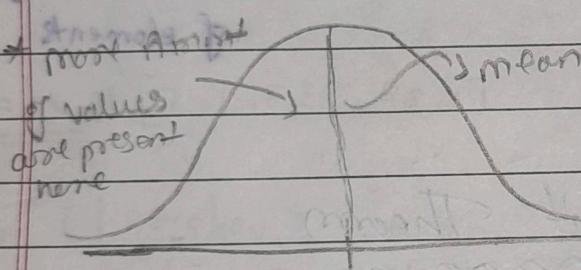
- Because of these we use Pearson correlation coefficient

(5) * Mean, Median and Mode

\Rightarrow Sample of height $\{168, 170, 150, 160, 182, 140, 175\}$ cm

$$\mu = \frac{\sum_{i=1}^n x_i}{n} = \frac{[]}{7} = 163.5$$

$n = \text{no. of sample}$



\Rightarrow This mean specifies measures of central Tendency

• Median

$$\Rightarrow [1, 2, 3, 4, 5] \quad \therefore \frac{5}{2} = 3$$

$\mu = 3$

{ If i add outlier

$$\Rightarrow [1, 2, 3, 4, 5, 50]$$

$$\mu = 13$$

= After i added outlier

$$\therefore \mu = 3 \quad \mu = 13$$

This very harmful for doing any statistical analysis

Note:- for Median the number should be in shaded order

Solved [1, 2, 3, 4, 5, 50]

Solved $\Rightarrow \frac{3+4}{2} = 3.5$

$m = 3.5$

If i have some outlier i can basically use median and mode

Mode:-

Mode: Basically find out the maximum number of numbers in the this particular Distribution

[repeated many times]

[1, 2, 3, 3, 3, 4, 5, 50] # Max no. of occurrence of elements

$\therefore m = 3$

[1, 2, 3, 3, 3, 4, 5] $\Rightarrow m = 2.83$

⑥ * Central Limit Theorem

may or may not belong to G.D

* $X \stackrel{D}{\sim} G.D(m, \sigma^2)$

Suppose i know my mean and variance value

m, σ^2

Central Limit theorem specifies two different point, suppose ; take sample from this (X) random variable X.

Then,

Taking Sample as 5,

from this random variable (x). And considering
 $n = 30$.

$\therefore [n \geq 30]$ (n is greater than or equal to 30 note points)

Every time selecting 30 random sample values

$$S_1, x_1, x_2, \dots, x_{30} = \bar{x}_1$$

$$S_2, x_1, x_2, \dots, x_{30} = \bar{x}_2$$

$$S_3, x_1, x_2, \dots, x_{30} = \bar{x}_3$$

$$S_4, x_1, x_2, \dots, x_{30} = \bar{x}_4$$

$$S_5, x_1, x_2, \dots, x_{30} = \bar{x}_5$$

$$S_6, x_1, x_2, \dots, x_{30} = \bar{x}_6$$

$$\vdots$$

$$S_{100}, x_1, x_2, \dots, x_{30} = \bar{x}_{100}$$

$$\boxed{\bar{x} \approx \text{GD}(u, \frac{\sigma^2}{n})}$$

$$-1 < (\bar{x} + 3\sigma) - (\bar{x} - 3\sigma) > 8\sigma$$

(f) Chebyshov's Inequality

$$\boxed{nx \approx \text{GD}(u, \sigma)}$$

x , random variable belongs to GD with some value of mean (u) and standard deviation or variance (σ). If we consider variance it will be sigma square (σ^2) then we know that based on the Empirical formula.

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

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

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

- * If my $y \neq \text{GID}$, in this we use Chebyshew's Inequality

$$y \neq \text{GID}$$

\Rightarrow Chebyshew's Inequality, says that if you want to find out the probability of a random variable falling within the range of 1σ ~~SD~~ will always be greater than or equal to

$$\left| 1 - \frac{1}{K^2} \right|$$

K = specifies for which range of SD, I have to find out this particular value, i.e., is my percentage of Data points

If my random variable does not belong to GID, I am basically able to find out

$$\Pr(\mu - K\sigma < x < \mu + K\sigma) > 1 - \frac{1}{K^2}$$

$\approx 68\%$ of the data will be: $K=1$

$$\Pr(\mu - 2\sigma < x < \mu + 2\sigma) > 1 - \frac{1}{K^2} = 1 - \frac{1}{4} = 75\%$$

(8) Pearson Correlation Coefficient

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

e.g.: x y
Height weight

$$\text{cov}(x, y) : \begin{array}{|c|c|} \hline & 1 & 2 & 3 & 4 & 5 \\ \hline 1 & & & & & \\ \hline \end{array}$$

$x \uparrow y \uparrow = +ve$

$x \uparrow y \downarrow = -ve$

But, in covariance it ~~does not~~ tell how much +ve and how much ~~not~~ -ve may be.

② Pearson Correlation Coefficient.

→ It says that based on ~~the~~ your variance of x and y ~~#~~ will be able to tell you that this parameter i.e. \Rightarrow strength (How strong this is) ~~correlated~~

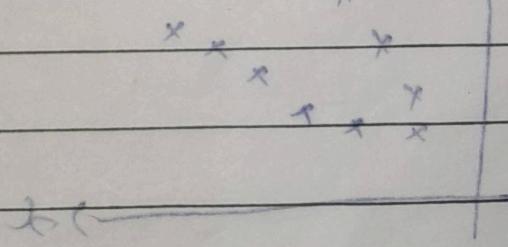
(2) Direction of Relationship

• Here, it is combining both the things

① Strength and ② Direction of Relationship

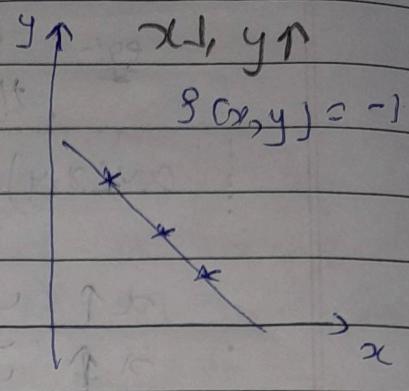
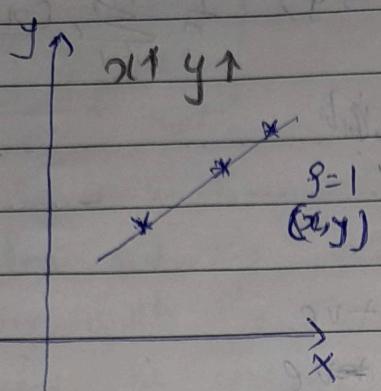
$$\text{Pearson CL} = \rho(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$\left[\rho(x, y) \in -1 \leq \rho \leq 1 \right]$ (the value ranges between -1 to +1) in Pearson

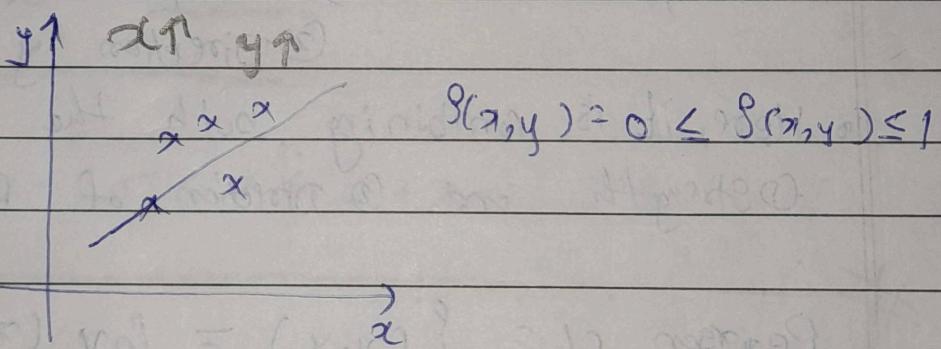
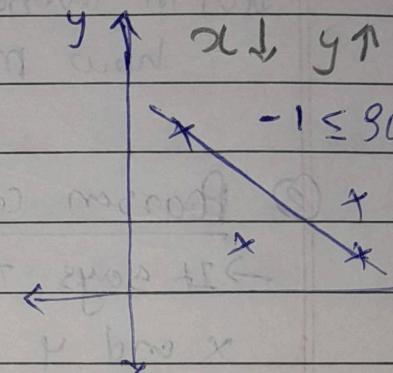
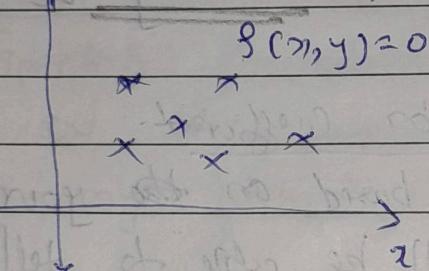


$$\boxed{g(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}} \quad [-1 \leq g \leq 1]$$

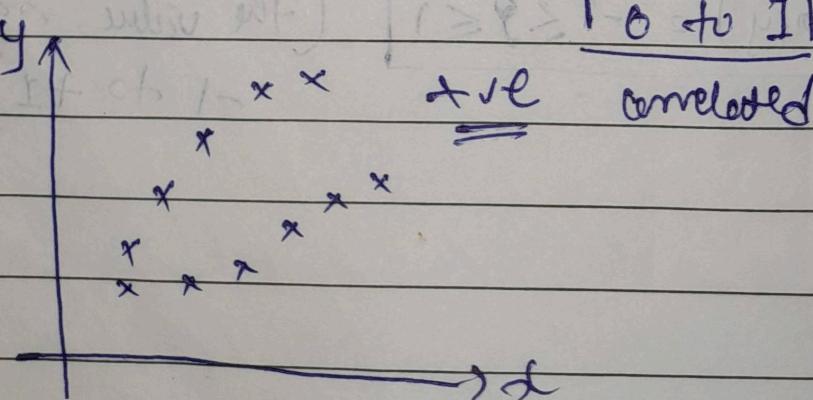
Correlation



No Relationship



Note: This technique is basically used in feature selection.



(a) Spearman's Rank Correlation Coefficient

The Spearman correlation coefficient is defined as the Pearson correlation coefficient between the rank variable

here we are trying to find out person of correlation of rank of X and rank of Y

$$\rho_s = \frac{\rho_{r_x, r_y}}{\sigma_{r_x} \sigma_{r_y}} = \frac{\text{Cov}(r_{x_i}, r_{y_i})}{\sigma_{r_x} \sigma_{r_y}}$$

where,

ρ denotes the usual Pearson correlation coefficient, but here applied to the rank Variable

$\text{Cov}(r_{x_i}, r_{y_i})$ is the covariance deviations of the rank variable

$\sigma_{r_x}, \sigma_{r_y}$ are the standard deviations to the rank variable

only if all n ranks are distinct integers, it can be computed using the popular formula

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

$$d_i = r_g(x_i) - r_g(y_i)$$

n = no. of observations

eg Calculate the correlation between IQ of a person with the number of hours spent in front of TV per week.

TQ, x_i Hours of TV per week, y_i

106 \dagger

86 2

86 2

101 50

99 28

103 29

97 20

113 12

112 6

110 \dagger

sort in ascending order x_i

rank y_i
di = $y_i - x_i$

TQ x_i	Hours of TV per week, y_i	rank x_i	rank y_i	d_i	d_i^2
86	2	1	1	0	0
91	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	9	81
112	6	9	2	7	49
113	12	10	4	6	36

$$\sum d_i^2 = 194$$

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

$$= 1 - \frac{6 \times 194}{10(10^2-1)}$$

$$P = 1 - \frac{29}{100} = 0.145757 \dots$$

Conclusion ⇒ The value is close to 0, the correlation between IQ and TV hour is very low.

⑩ * finding outliers in Dataset Using Z-score and IQR

Definition: An outlier is a data point in a dataset that is distant from all other observation. A data point that lies outside the overall distribution of the dataset.

① Z-score

Consider Data point,

$$f = [1, 2, 3, 4, 5] \quad \begin{matrix} \text{if we want to invert} \\ \text{this distribution} \end{matrix} \quad \text{SND} \Rightarrow \mu = 0, \sigma = 1$$

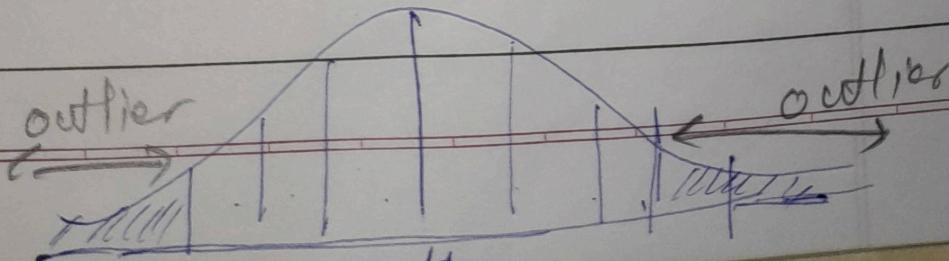
$\mu = 3 \quad \sigma = 1$ in SND

∴ Hence, in order to convert this, we usually apply Z-score observation

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

Out fall after the 3rd SD are

• Z-score, tell that in SND whether the ~~far~~ considered as outliers



② IQR

75%, 25% values in a dataset

Steps:-

- 1) Arrange the data point in Increasing order
- 2) Calculate first (Q_1) and third quartile (Q_3)
- 3) find Interquartile range ($Q_3 - Q_1$)
- 4) find lower bound $Q_1 - 1.5$
- 5) find upper bound $Q_3 + 1.5$

Anything that lies outside of lower and upper bound is an outlier.

• Various ways of finding outliers

- 1) Scatter plot or looking which is outlier or not
- 2) Box plot
- 3) Z-score
- 4) Using IQR