

Machine Learning

⇒ what is the statistics And its application.

Statistics is a field that deals with collection organization, analysis, interpretation and presentation of the data

so it's application is to make the decision based on the given data.

Types of Statistics ⇒

There are two types of Statistics -

- ① Descriptive
- ② Inferential

① Descriptive ⇒

① measure of central tendency

① mean.

② median.

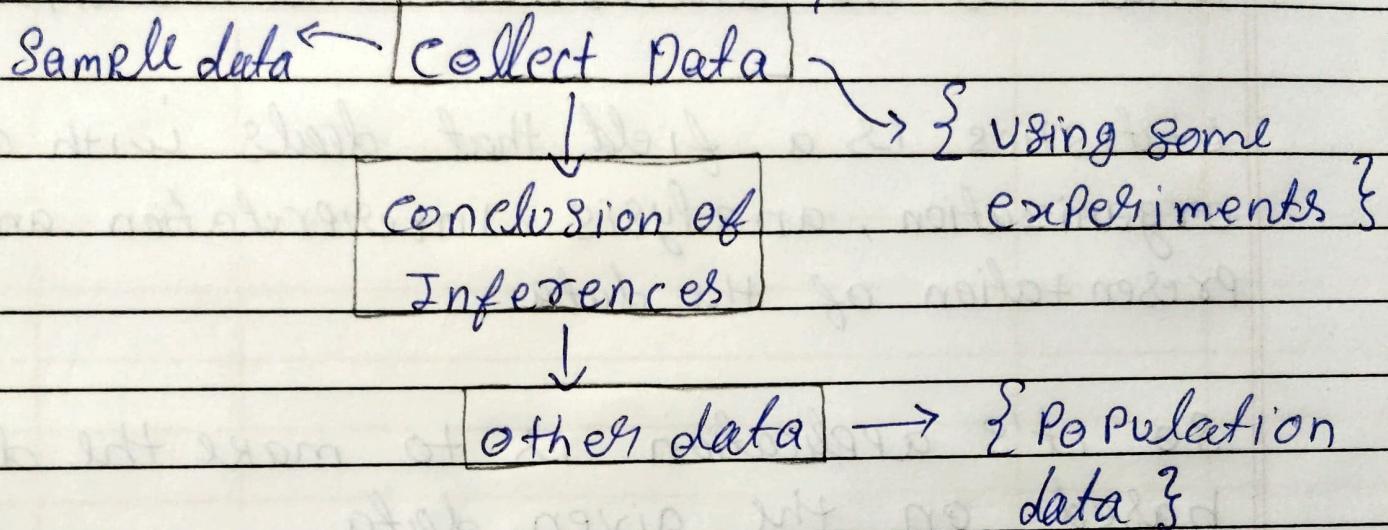
③ mode.

④ measure of Dispersion

① variance.

② Standard deviation.

① Inferential \Rightarrow



\Rightarrow population VIS Sample data \Rightarrow

A large circle represents the 'Population (100K)'. Inside it, four smaller circles represent 'Sample data (10K)'. An arrow points from the population circle to the sample data, with the label '(N)' above the arrow. Another arrow points from the sample data back to the population, with the label '(n)' below the arrow.

① → measure of central tendency

① mean

② Median

③ Mode

① Mean →

Population data

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

Sample data.

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

② age = {1, 3, 4, 5}

$$\mu = \frac{1+3+4+5}{4} = \frac{13}{4} = 3.25$$

③ median → age = {4, 3, 1, 5, 100}

→ first sort the numbers.

{1, 3, 4, 5, 100}

median → 4 → fight with outliers impact

(ii) mode \Rightarrow

$$\text{age} = \{4, 3, 2, 1, 1, 4, 4, 5, 2, 100\}$$

In mode we select the element which have maximum frequency.

$$\text{o/p} \Rightarrow 4$$

⑧ Measures of Dispersion:

i) variance

ii) Standard Deviation.

{ Spread is more }

$$\text{age 1} = \{2, 2, 4, 4\}$$

$$\mu_1 = \frac{2+2+4+4}{4}$$

{ spreading less }

$$\Rightarrow 3$$



{ the data points are near to the mean }

$$\text{age 2} = \{1, 1, 5, 5\}$$

$$\mu_2 = \frac{1+1+5+5}{4}$$

$$\Rightarrow 3$$



{ the data points are far away from the mean }

①. variance

population data $\{N\}$

(points)

(population mean)

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



(Population Size)

$$\text{Age 1} = \{2, 2, 4, 4\}$$

$$\text{Age 2} = \{1, 1, 5, 5\}$$

$$\mu = \frac{2+2+4+4}{4}$$

$\Rightarrow 3$

x_i	μ	$(x_i - \mu)^2$
2	3	1
2	3	1
4	3	1
4	3	1
$\frac{4}{4}$		

$$\mu = \frac{1+1+5+5}{4}$$

$\Rightarrow 3$

x_i	μ	$(x_i - \mu)^2$
1	3	4
1	3	4
5	3	4
5	3	4
16		

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

$$\sigma^2 = \frac{4}{4} = \boxed{1}$$

{less dispersion}

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

$$\sigma^2 = \frac{16}{4} = \boxed{4}$$

{more dispersion}

(i) Sample data $\{x_i\}$

point

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

mean (Sample)
Sum of size

{ why $(n-1)$ \Rightarrow Imp invas. }so if we use the $\{s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}\}$.

Imp it is underestimating the true population variance.

and if we use the $\{s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}\}$ known as biased correction

it comes with Unbiased Underestimation and it is focusing the different - different data points.

(ii) Standard Deviation.

$$S.D = \sqrt{s^2} \rightarrow \text{population / sample standard deviation}$$

How far data points is away from the mean.

\Rightarrow Variable

variable is a property that can take up any value

eg: age = 25, Height = 7". 2'
Gender = male

Types of variable:

(i) Quantitative

Discrete Quantitative
a whole no
not any fraction
value
{eg:- {3, 5, 6, 18, 80}}

Continuous Quantitative

any value
can be whole
no and fraction
as well.
{eg: 2, 2.5, 3.3, 6.8}

(ii) Qualitative / Categorical

eg: Gender
male female

Colors:
red
blue
green

Random variable $\Rightarrow (X)$

it is a function whose values are derived from the "process or experiments"

$$X = \{0, 1\}$$

Tossing a coin.

Random variable

Discrete Random

eg: Tossing a coin
rolling a dice

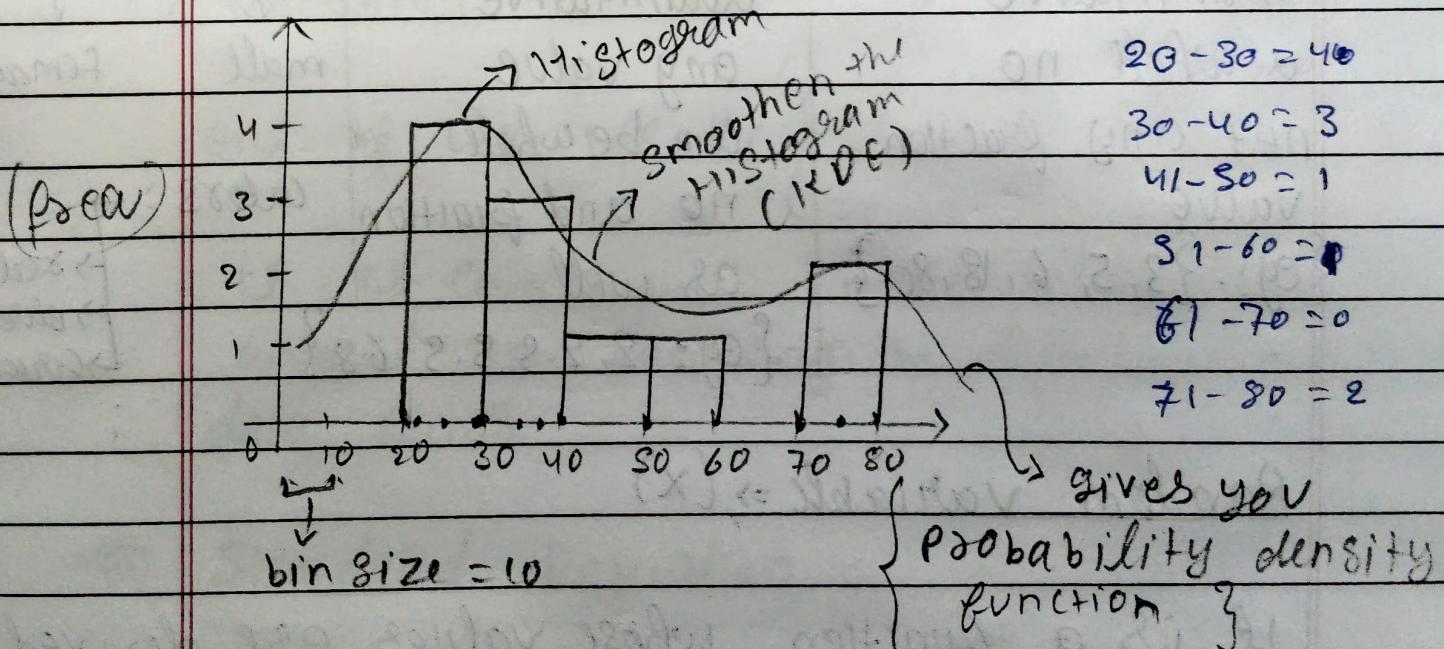
Continuous

Random variable.

eg: Tomorrow how
many inches it is
going to rain.

Histograms \rightarrow It's count the frequency of element

$$X = \{23, 24, 25, 30, 34, 36, 40, 50, 60, 75, 80\}$$



KDE \rightarrow Kernel density Estimation

and it gives you the probability density function

Percentiles and Quartiles.

Percentage:

$$\{1, 2, 3, 4, 5, 6\}$$

odd number = 3

$$\text{Percentage of odd number} = \frac{3}{6} \times 100 = 50\%$$

Percentiles:

A percentile is a value below which certain percentage of observations lie.

$$\{2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, 9, 9, 10\}, n=14$$

$$\text{Percentile of } x = \frac{\# \text{ value below } x}{n} \times 100$$

below x values are = 11

$$= \frac{11}{14} \times 100 = 78.57\% \text{ of value } 9$$

Percentile ranking.

So basically it means that the 78.57% of the entire distribution is less than 9

if Percentile is given and we find out the 25 percentile of the given sample value is.

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

$$\begin{aligned}\text{Value} &= \frac{25}{100} \times (14+1) \\ &= \frac{25 \times 15}{100} = 3.75\end{aligned}$$

but we don't have 3.75, so in this case take the avg of lower and higher value of the range like in the data points 3 & 4 are present.

$$80 \quad \frac{3+4}{2} = 3.5$$

$$\Rightarrow 3.75 \approx 3.50$$

② Quartiles \rightarrow

$$\left. \begin{array}{l} 25 \text{ Percentile} = 1^{\text{st}} \text{ Quartile} \\ 50 \text{ Percentile} = 2^{\text{nd}} \text{ Quartile} \\ 75 \text{ Percentile} = 3^{\text{rd}} \text{ Quartile} \end{array} \right\}$$

5 Number Summary.

outlier.

$$\text{eg: } \{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27 \}$$

- ① minimum
- ② 1st Quartile (25%)
- ③ median
- ④ 3rd Quartile (75%)
- ⑤ maximum.

IQR \Rightarrow Inter Quartile Range

$$Q_3 - Q_1 \rightarrow 1^{\text{st}} \text{ Quartile}$$

↓

3rd Quartile

Removing the outlier. \Rightarrow So in case of removing the outlier we have to calculate the lower fence and higher fence

- ① { if any value is less than lower fence then it is removed as outlier }
- ② { if any value is greater than higher fence then it is removed as outlier }

#

27

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR})$$

$$\text{Higher fence} = Q_3 + 1.5(\text{IQR})$$

$$\Rightarrow Q_1 \Rightarrow \frac{\text{Percentile}}{100} \times (n+1)$$

$$= \frac{25}{100} \times (20) = 5^{\text{th}} \text{ position} = \underline{3} \text{ value}$$

$$\Rightarrow Q_3 \Rightarrow \frac{75}{100} \times 20 = 15^{\text{th}} \text{ position} = \underline{7} \text{ value}$$

$$\text{IQR} = 7 - 3 = \underline{4}$$

$$\Rightarrow \text{Lower fence} \Rightarrow Q_1 - 1.5(\text{IQR})$$

$$\Rightarrow 3 - 1.5 \times 4$$

$$\Rightarrow 3 - 6 = \underline{-3}$$

$$\Rightarrow \text{Higher fence} \Rightarrow Q_3 + 1.5(\text{IQR})$$

$$\Rightarrow 7 + 1.5(4)$$

$$= \underline{13}$$

so the lower fence is -3 and Higher fence is 13.

the values which is less than -3 and greater than 13 is considered as outliers.

In our example the 27 is present which is greater than 13 so it is considered as an outlier. and remove all the outliers.

Values = { 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9 }

- ⇒ minimum = 1
- ⇒ 1st Quartile = 3
- ⇒ median = 5
- ⇒ 3rd Quartile = 7
- ⇒ maximum value = 9

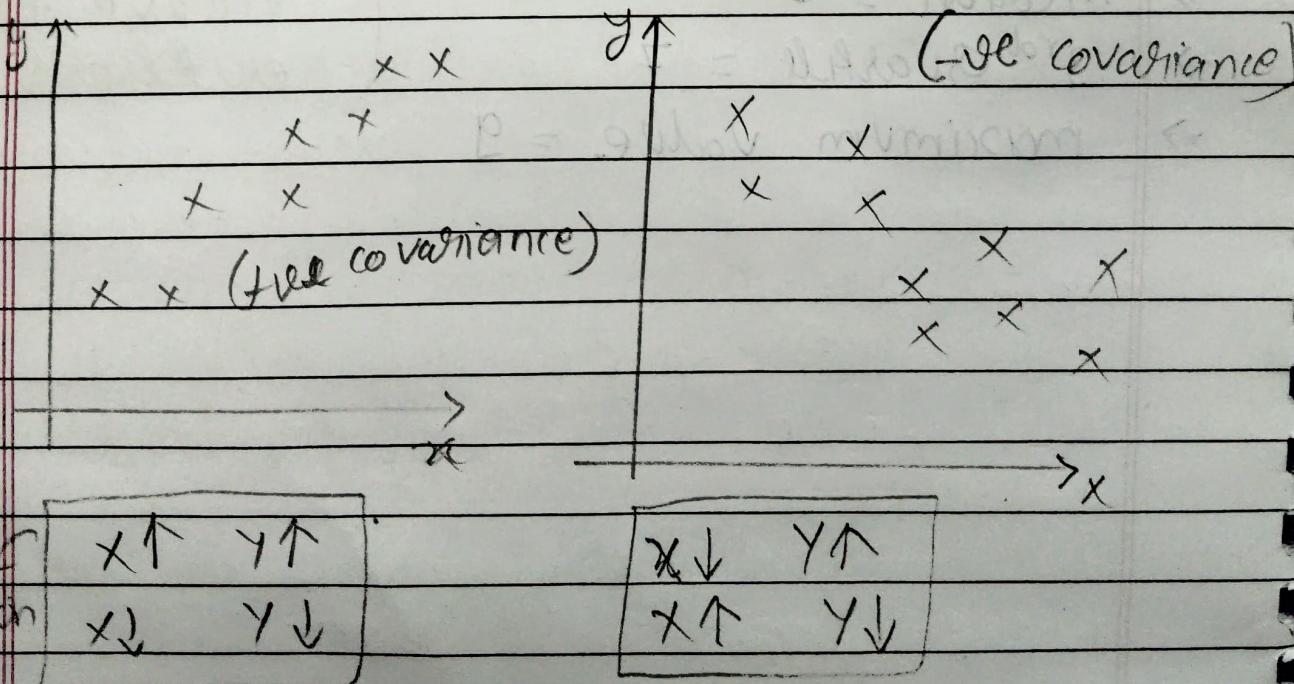
} we use the
Box Plot to
remove the
outliers.

Covariance and correlation.

Covariance and correlation are two statistical measures used to determine the relationship between two variables. Both are used to understand how changes in one variable are associated with changes in another variable.

Covariance \Rightarrow

Covariance is a measure of how much two random variables change together. If the variable tends to increase and decrease together, the covariance is (+ve). If one tends to increase when the other decreases, the covariance is (-ve).



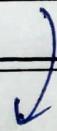
#

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

#

$$\begin{aligned} \text{Cov}(x, x) &= \sum_{i=1}^n \frac{(x_i - \bar{x})(x_i - \bar{x})}{n-1} \\ &= \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1} \end{aligned}$$

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



↳ variance of x

Covariance of x

with respect to x

$x_i \rightarrow$ Data points of random variable x

$\bar{x} \rightarrow$ Sample mean of x

$y_i \rightarrow$ Data points of random variable y

$\bar{y} \rightarrow$ Sample mean of y

Advantage:

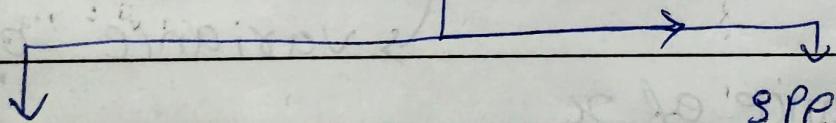
- ① Quantify the relationship between x and y

Disadvantage:

- ① covariance does not have a specific limit value

$$\text{cov}(x, y) \Rightarrow -\infty \text{ to } \infty$$

Correlation



Spearman Rank

Pearson correlation

coefficient

correlation

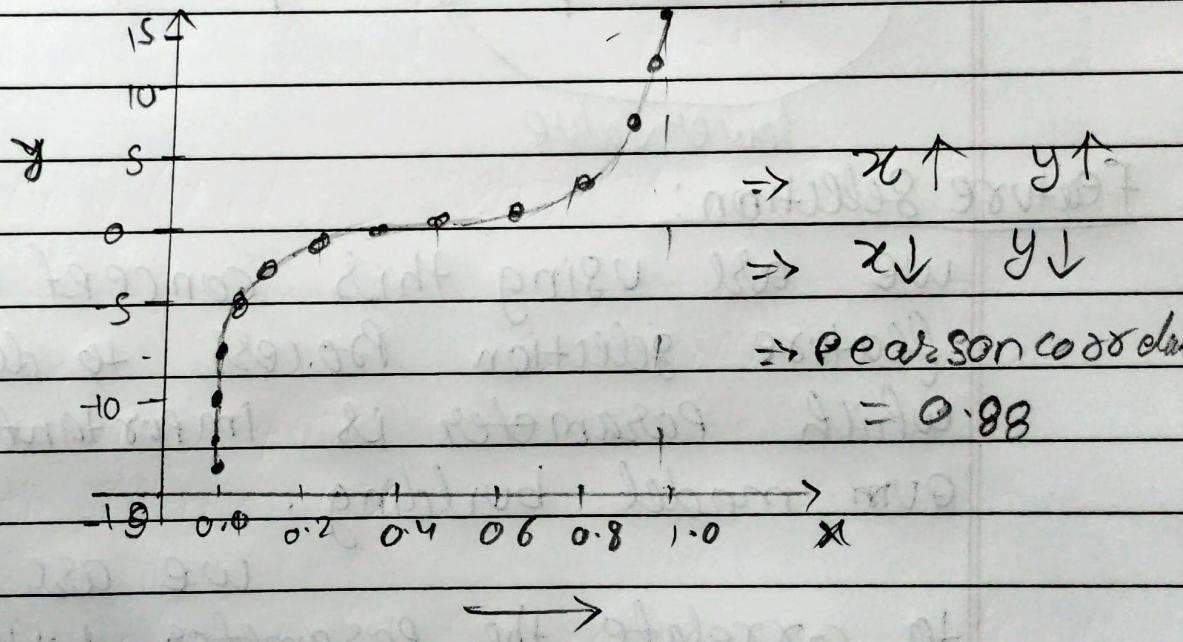
- ① Pearson correlation coefficient:

it's limit the value $[-1, 1]$

$$\$ P_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

standard deviation.

- ① The more the value towards (+1) the more (+ve) correlated x & y is.
- ② The more the value towards (-1) the more (-ve) correlated it is (x, y)
- ③ Spearman Rank correlation:



} So the Pearson correlation is not correlating the non-linear data points.

$$\rho_s = \frac{\text{cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

rank

Rank \rightarrow So basically Rank is considered the lower value is lower rank.

x	y	$R(x)$	$R(y)$
1	2	2	1
3	4	3	2
5	6	4	3
7	8	5	5
0	7	1	4

lower valve

Feature Selection:

we are using this concept in our feature selection process to determine which parameter is important for the our model building.

we are trying to correlate the parameters with each other with respect to the output.

eg:

(+ve)

~~remove / drop this feature~~

Size of House	No of rooms	location	No. of people	Haunt-ed	price
↑	↑	↑	people living	↓	↑
(10)	• (10)	• (10)	(~0)	-ve corresp	locality