

Stat

Statistics

- ① Population → population mean, median mean } measure of central tendency
- ② Sample → sample mean mode $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- ③ Random Variable
 - Discrete Random Variable
 - Continuous Random Variable

Sample → just taking an amount of data as a sample and create all the prediction start from here.

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

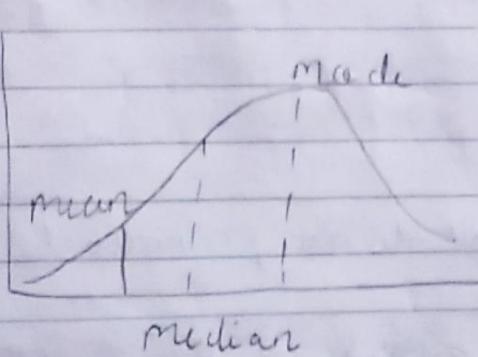
Random Variable

- discrete → whole no., cannot be a floating no.
-) continuous → within a Range of Values we can have any values

2 Gaussian dist / Normal Distribution

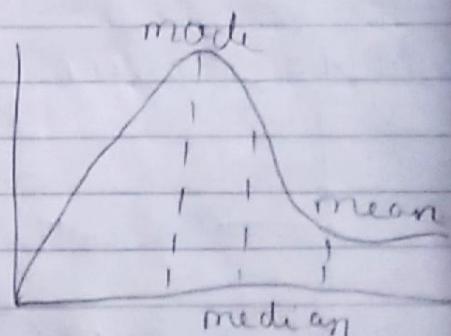
- They have approximate a wide variety of random variables
- Distributions of sample means with large enough sample sizes could be approximated to normal.
- All computable statistics are elegant
- Decision based on normal distribution insights have a good track record.

Skewness



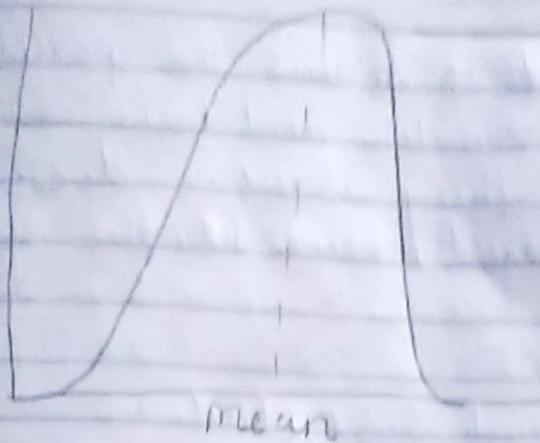
Negative direction

$$\text{mean} < \text{median} < \text{mode}$$



mode < median < mean

Positive direction



Symmetrical Normal Dist with no skewness

Mean - Median - mode

denoted as :-

Normal $\rightarrow N$

\curvearrowleft \rightarrow Distribution.

$\mu \rightarrow$ mean

$\sigma^2 \rightarrow$ variance.

$$N \curvearrowleft (\mu, \sigma^2)$$

\rightarrow A lower mean would result in the same shape of the distribution, but on the left side of the plane.

\rightarrow A higher mean would move to right.

Distribution :- A function that shows the possible values for a variable and how often they occur. Thing about a die numbered from 1 to 6.

The Standard normal Distribution

$$\rightarrow z = \frac{x - \mu}{\sigma} \quad z \sim N(0, 1)$$

\downarrow \downarrow
 mean s.d.

z-score

- \rightarrow Adding and subtracting values from all the data points does not change the standard deviation.

Covariance

The two variables are collected and which are correlated and the main statistic to measure this correlation is called Covariance.

It can be :-

- > 0, The two variable moves together.
- = 0, The two variables are independent.
- < 0, " " moves in opposite direction.

Sample Formula

$$s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

Population Formula

$$\sigma_{xy} = \frac{\sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)}{n}$$

Correlation

Co-relation adjusts covariance so that the relationship between the two variables become easy and intuitive to interpret.

$$\text{Correlation} = \frac{\text{cov}(x, y)}{\text{st.dev}(x) \times \text{st.dev}(y)}$$

$$\frac{s_{xy}}{s_x s_y}$$

$$-1 \leq \text{correlation coefficient} \leq 1$$

$$\text{Cov Coeff} = 1$$

The entire variability of one variable is explained by the other.

Causality = important to understand the direction of causal relationships.

Standard Deviation Formula

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

Sample Standard Deviation

Population Standard Deviation.

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

Variance of the Fall Data

$$\rightarrow 6, 8, 10, 12, 14, 16, 18, 20, 22, 24.$$

x_i	$d_i = x_i - 14$	Deviation from mean	$(x_i - \bar{x})^2$
6	-4	-9	81
8	-3	-7	49
10	-2	-5	25
12	-1	-3	9
14	0	-1	1
16	1	1	1
18	2	3	9
20	3	5	25
22	4	7	49
24	5	9	81
	5		330

14 = is assumed mean.

Formulas

$$\text{Skewness} = \frac{1}{n} \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}}^3$$

Variance and std. deviation

$$\text{Sample Variance formula: } s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

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

$$\text{Sample Standard deviation formula: } s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$\text{Population SD: } \sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{n}}$$

Central Limit Theorem

Sampling dist

$$N\left(\frac{\mu}{\sqrt{n}}, \frac{\sigma^2}{n}\right), n > 30 \rightarrow \text{sample size}$$

→ CLT allows us to perform tests, solve problems and make inferences using the normal distribution, even when the population is not normally distributed.

→ The mean of the samples we extract will be closer to normally distributed if we extract more samples.

→ The distribution of the sample mean is expected to have a mean equal to the mean of the original dataset.

→ The distribution of the sample mean is expected to have a variance equal to the variance of the original dataset, divided by the sample size.

Standard Error

→ The standard deviation of the dist formed by the sample means

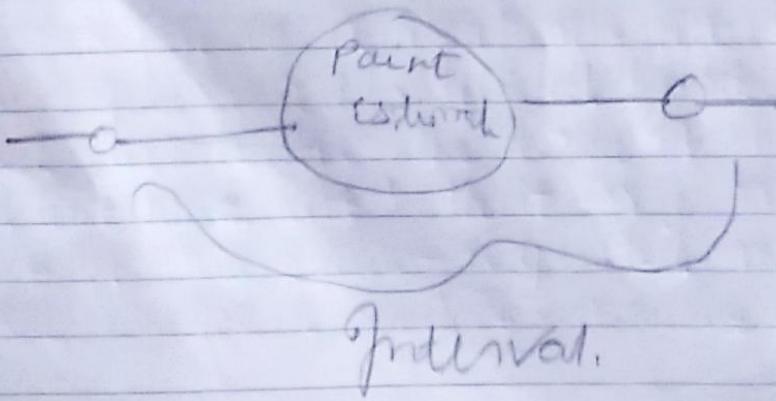
$$n\sigma\left(\frac{\mu}{\sqrt{n}}\right) \rightarrow \text{variance}$$

$$\text{S.E} = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$$

→ Standard Error decreases when sample size increases.

Estimator and Estimates

- ① Point Estimates
- ② Confidence Interval Estimates



Chebyshov's Inequality

$$\rightarrow 1 - \frac{1}{k^2}$$

\rightarrow when $y \neq c.s.$

Pearson Correlation Coefficient

$$\text{Covariance} = \text{cov}(u, v) = \frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(v_i - \bar{v})$$

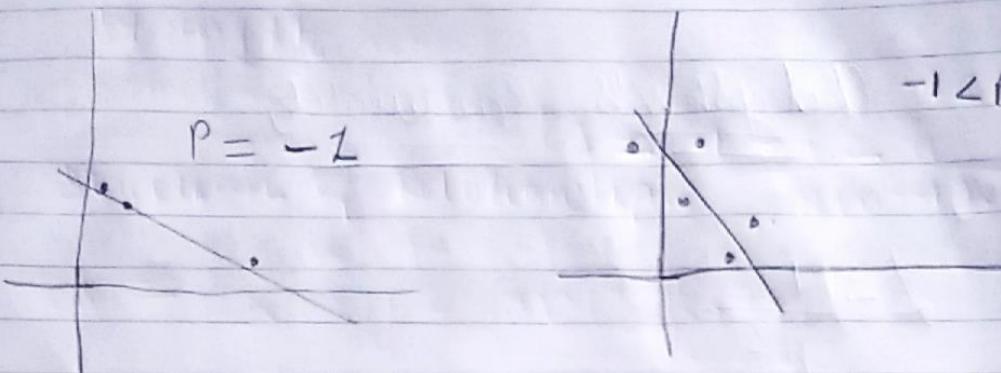
$$\text{Pearson } \rho = \rho(u, v) = \frac{\text{cov}(u, v)}{\sigma_u \sigma_v}$$

- \rightarrow It measures linear correlation between two variables u and v .
- \rightarrow It has value between +1 and -1.

+1 \rightarrow Total positive linear correlation

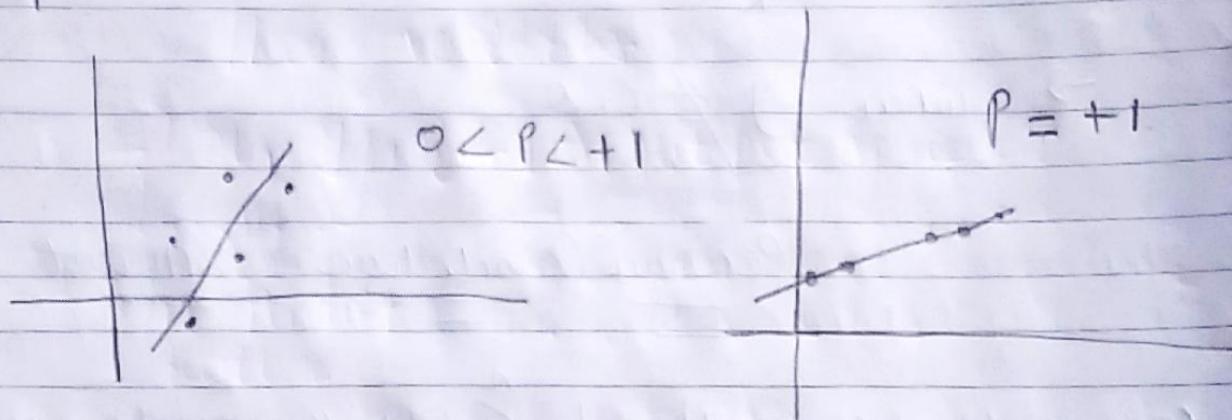
-1 \rightarrow negative linear correlation.

0 \rightarrow no linear correlation.



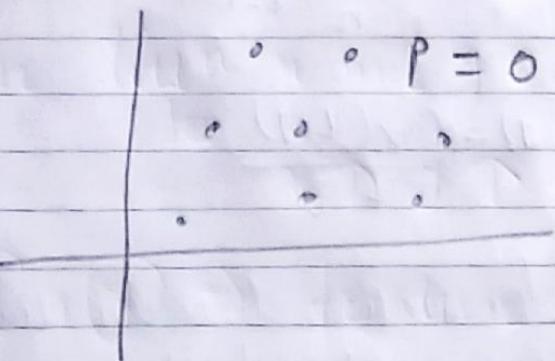
$$P = -1$$

$$-1 < P < 0$$



$$0 < P < +1$$

$$P = +1$$



$$P = 0$$

Defi :- Pearson Correlation Co-efficient is the covariance of the two variables divided by the product of their standard deviations. The form of definition involves a "Product moment", that is, the mean of the product of the mean-adjusted random variables. hence the modifier product-moment in the name.

For a Population :-

$$r_{xy} = \frac{\text{cov}(u, j)}{\sigma_u \sigma_j}$$

Cov → Covariance

σ_u → Standard deviation of u

σ_j → Standard deviation of j

For a sample :-

$$r_{xy} = \frac{\sum_{i=1}^n (u_i - \bar{u})(j_i - \bar{j})}{\sqrt{\sum_{i=1}^n (u_i - \bar{u})^2} \sqrt{\sum_{i=1}^n (j_i - \bar{j})^2}}$$

where :

n = Sample size

u_i, j_i are the individual sample points
marked with i .

$$\bar{u} = \frac{1}{n} \sum_{i=1}^n u_i \quad (\text{The sample mean});$$

$\text{for } \bar{j}$

1 Strength

2 Direction of Relationship

Three Variables.

$\{n, y\}, \{z\} \rightarrow$ output variable
↓

Independent
feature

Can but n and y is 2.
So

$n \uparrow y \uparrow$ so can is 2.

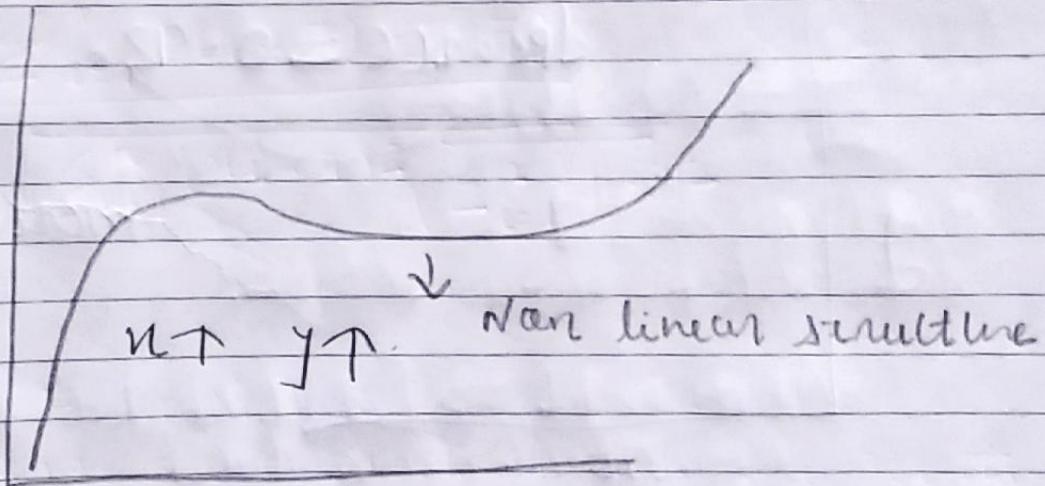
(n), (y) → Feature is same.

so drop in (2) of the
feature that can be
 n or y .

Pearman's Rank Correlation

→ The Spearman correlation coefficient is defined as the Pearson correlation coefficient between the rank variables.

→ Rank Variables :→ betⁿ a set of items such that for any two items) the first is Ranked higher , Rank lower than or Ranked equal to the second.



Most Important :- When the data are roughly elliptically distributed and there are no prominent outliers, the Spearman Correlation and Pearson Correlation give similar values.

Formula

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

$$R_s = \text{sgn}_x \text{sgn}_y = \frac{\text{cov}(\text{rg}_x, \text{rg}_y)}{\sigma_{\text{rg}_x} \sigma_{\text{rg}_y}}$$

ρ = Pearson correlation coefficient
 $\text{cov}(\text{rg}_x, \text{rg}_y) \rightarrow$ covariance of the Rank Variables.
 σ_{rg_x} and σ_{rg_y} are the standard deviation of the Rank Variables.

$$\text{SpearCC} = 0.92$$

J

Increasing

n

$$\text{SpearCC} = -0.91$$

J

Decreasing

n

Graph

x_i	y_i
106	4
100	27
86	2
101	50
99	28
103	29
97	20
113	12
112	6
110	17

Rank(y_i) This is created by the highest number in y_i taken.

⑦ \rightarrow 7th Rank

① \rightarrow highest value.

(i) Arrange them in ascending order

x_i	y_i	Rank(x_i)	Rank(y_i)	d_i	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	4	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

$$g = 1 - \frac{6\epsilon d^2}{n(n^2-1)}$$

tee Give

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

which evaluates to $g = -29/165$

$$= -0.175$$

with a p-value = 0.627

9

Finding an outlier in a Dataset

→ What is an outlier?

→ An outlier is a data point in a dataset that is distant from all other observations. A datapoint that lies outside the overall distribution of the dataset.

Techniques:

① Z-score.

② IQR → Inter quartile Range.

→ What is the reason for an outlier to exists in a Dataset?

③ Variability in the data.

④ An experimental measurement error.

→ What are the impacts of having outliers in a dataset?

⑤ It causes various problems during our analysis.

⑥ It may cause a significant impact on the mean and the standard deviation.

Various ways to find.

- ① Scatter plot
- ② Boxplot
- ③ Using Z-score
- ④ Using IQR.

Selecting outlier using Z-score :-

$$\text{Formula} = (\text{observation} - \text{mean}) / \text{standard deviation}$$

$$Z = (x - \mu) / \sigma$$

Outliers = []

def detect_outliers(data):

threshold = 3 [3 std if data is falling not an outlier]

mean = np.mean(data)

outlier, if

std = np.std(data).

faling away from
3 std its outlier

for i in Data:

Z-score = (i - mean) / std

if np.abs(Z-score) > threshold:
outliers.append(i)

return outliers.

outlier-pt = detect_outliers(dataset)
outlier-pt

IQR

75% - 25% Values in dataset

Steps

- ① Arrange the data 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 $q_3 + 1.5$

sorted(data set)

quantile2, quantile3 = np.percentile([dataset, [25, 75]])
print(quantile1, quantile3)

→ Find the IQR.

iqr-value = quantile3 - quantile2.
print(iqr-value).

Find the lower bound value and the higher bound value.

$$\text{lower-bound-val} = \text{quantile 1} - (1.5 * \text{IQR-value})$$

$$\text{upper-bound-val} = \text{quantile 3} + (1.5 * \text{IQR-value})$$

Print (lower-bound-val, upper-bound-val).



Values away from these bounds are considered as outliers.

Standardization vs Normalization

- ① Normalization - helps you to scale down your feature between 0 to 1.
- ② Standardization - helps you to scale down your feature based on standard normal distribution ($\mu_{\text{mean}}=0$, $\sigma_{\text{std}}=1$).

± Normalization - (min max Normal)

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

$df = pd.read_csv(\text{"IRIS Read data}, header=None, usecols=[0, 1, 2]).$

$df.columns = ["Class", "PetalLength", "PetalWidth"]$

$df.head()$.

From Akbari - preparing import minmaxScalar
scaling = minmaxScalar()
scaling = fit (dt [L'Almat', 'math'])

② Standardization (Z-score normalization)

All the features will be transformed in such a way that it will have the properties of a standard normal distribution with mean(μ) = 0 and standard deviation (σ) = 1.

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

From Akbari - preparing import standard scalar

scaling = standardScalar()

scaling = fit_transform (dt [L'Almat', 'math'])