

Statistic

→ Statistics is the discipline that concerns the collection, organization, analysis, interpretation and presentation of data.

Descriptive Statn

① Analyzing data, Summarizing

data, Organizing the data,

Avoid whole info. of data,

in the form of

"Numer and Graph"

② Bar plot, Histogram, pie chart

PCDF, PDF .

"Visualisation"

③ Measure of central tendency

"Mean, median, mode"

④ Measure of variance.

⑤ All Distribution function

Inferential statn

Election-

2 million

sample

Exit poll

Party 1 →

Party 2 →

from entire population taking

some sample, we are try
inferring some info. and

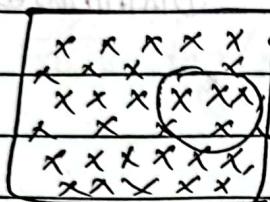
come to the conclusion

which party win (chance)

→ Z-test, T-test, chi-square test.

→ Confidential interval.

- Population → the whole group of entire observation or data
- Sample → A Collection of data from part of the Population.



Population

Sample

$$\text{mean } \mu = \frac{\sum x_i}{N}$$

$$\text{mean } \bar{x} = \frac{\sum x_i}{n}$$

Random variable

Discrete Continuous.

How we choose the sample from population?

$x \otimes x x$	\rightarrow	$xx \otimes x \otimes$	\rightarrow	$xx \otimes x$
$xx \times x \otimes$		$x \otimes x \otimes x$		$x \otimes x x$
$x \otimes x x x$		$\otimes x \otimes xx$		$x \otimes \otimes \otimes$
$\otimes x \otimes x x$	Bent	$xx \otimes x \otimes$		$\otimes \otimes \otimes \otimes$
$x \otimes x x \otimes$		$\otimes x x \otimes x$		$x x x \otimes$

① Random Sample

(Pick Randomly from the list)

② Systematic Sample

(Such as every 3rd).

follow a pattern

③ Stratified Sample

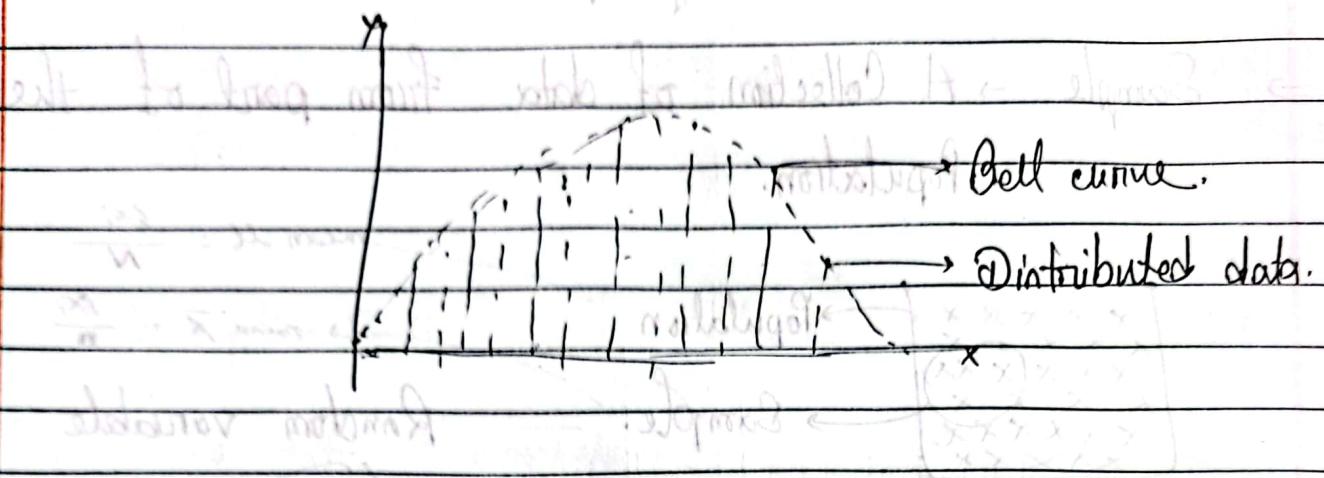
(Randomly, but in relation to group size)



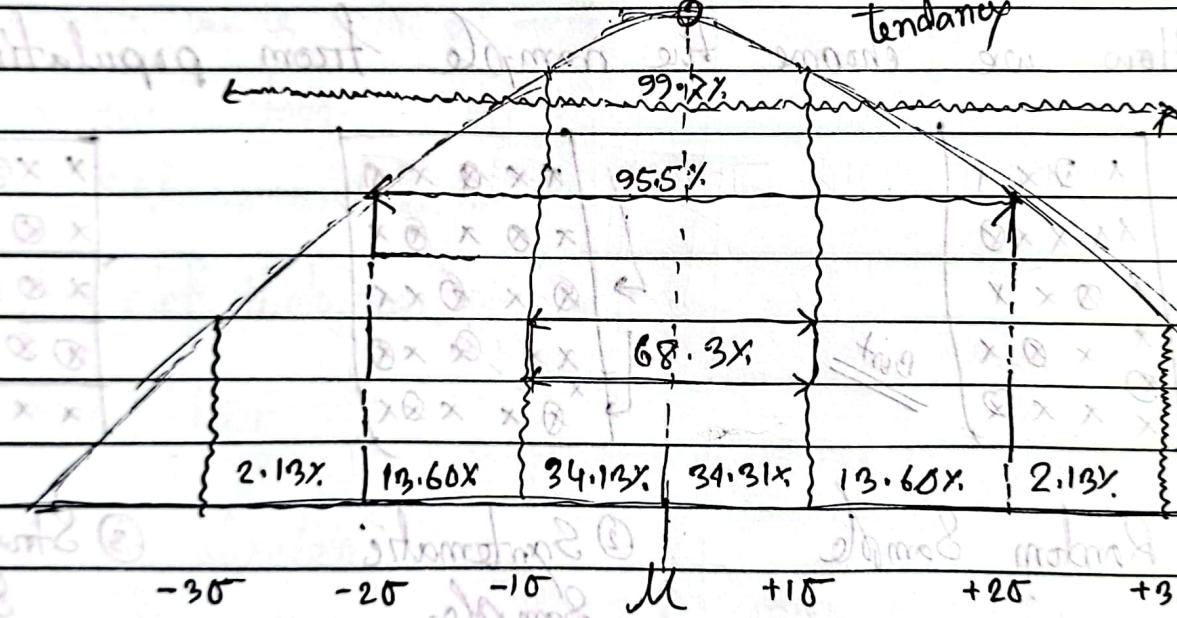
④ Cluster Sample

choose whole group randomly.

* Gaussian Distribution or Normal Distribution.



Measure of central tendency



$\sigma \rightarrow$ Standard deviation

$\theta \rightarrow$ Skew

Empirical Formula.

Standard Normal Distribution

$$\mu = 0$$

$$x_i - \mu$$

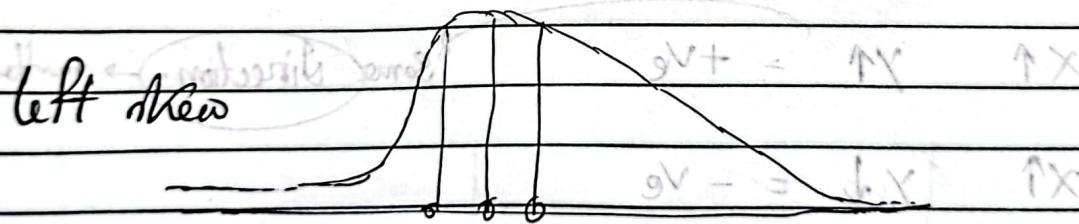
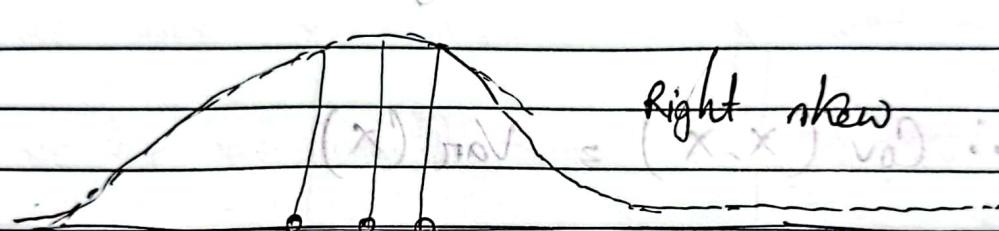
$$\sigma = 1$$

$$\sigma$$

* Log Normal distribution.

we have Random variable x .

→ Normal Distribution belongs to log normal distribution
if $\log(x)$ is Normally distributed.



log Normal Distribution.

feature 1

follow log ND

↓
Standard Normal
Distribution

feature 2

follow GND or ND

↓
Standard Normal Distribution

↓ Scale the feature ← them.

→ Two Random variables x, y .

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N}$$

$$\text{Var}(x) = \frac{\sum (x_i - \bar{x})^2}{N} \quad \sigma = \sqrt{\text{Var}}$$

$$= \frac{\sum (x_i - \bar{x})(x_i - \bar{x})}{N}$$

$N-1$ for sample

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

কোথায় তার স্থূল বিচ্ছিন্ন অবস্থার থাকে বা এখন \rightarrow variance.

$$x \uparrow \quad y \uparrow = +ve \quad \text{Same direction.} \rightarrow \text{matter.}$$

$$x \uparrow \quad y \downarrow = -ve$$

Pearson Correlation Co-efficient

$$① \text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N}$$

$$② \text{Pearson CC} = P_{(x,y)} = \frac{\text{Cov}(x, y)}{\sigma_x, \sigma_y}$$

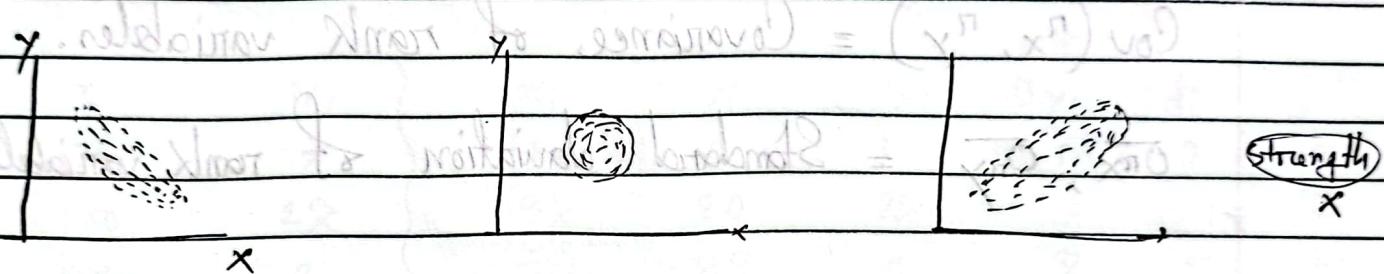
Covariance \Rightarrow Direction.

+ve and -ve value, but how much may
 $b \in -\infty \dots +\infty$. (actually we don't know)

Correlation \Rightarrow "Direction" and "Strength"

+ve and -ve correlation

-1 to 1



$$\text{Cov}(x, y) < 0$$

$$\text{Cov}(x, y) \approx 0$$

$$\text{Cov}(x, y) > 0$$

Perfect +ve
correlation

High +ve
correlation

Low +ve
correlation

+1 +0.9 0.5

-0.5
Low -ve
corr.

-0.9
High -ve
corr.

-1
Perfect -ve
corr.

0
No Conn.

"Spearman Correlation":

wikipedia.

Random variable x, y

$$r_s = \rho(r_x, r_y) = \frac{\text{Cov}(r_x, r_y)}{\sigma_{r_x} \cdot \sigma_{r_y}}$$

$r = \text{rank}$.

ρ = Pearson correlation co-efficient.

$\text{Cov}(r_x, r_y)$ = Covariance of rank variables.

$\sigma_{r_x}, \sigma_{r_y}$ = Standard deviation of rank variables.

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

$d_i = r(x_i) - r(y_i)$ in the difference between the two ranks of each observation.

Step 1: Sort the data by the first column (x_i). Create a new column r_x and assign it the ranks
1, 2, 3, 4 - - -

Step 2: Next, sort the data by the second column y_i . Create a fourth column r_y .

and Similarly assign it the rank values. 1, 2, 3, ..., n

③ Create a fifth column d_i to hold the difference between to rank Column (r_x_i, r_y_i).

④ Create one final Column d_i^2 to hold the value of column d_i squared.

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

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

$$= \frac{-29}{165} = 0.17$$

→ Mean, Median, Mode → Math is fun.

outliers ~~remove~~ remove or use Median.

কোনো কি প্রয়োজন নেই Krish NAIK.

→ Central Limit Theorem.

$$\text{Population mean } \mu = \frac{\sum x_i}{n}$$

$$\text{Sample mean } \bar{x} = \frac{\sum x_i}{n}$$

→ Actually আপনি এ সামগ্রেজলো ফিল্ড তথ্য দিবে তাহলে Sample জন।

→ Random variables,

Numerical

Discrete

Continuous.

Nominal

Ordinal.

Categorical

Central Limit Theorem.

$$X \approx \text{GD}(\mu, \sigma^2)$$

$$S_1 - - - - - x_{20} = \bar{x}_1$$

$$S_2 - - - - - x_{10} = \bar{x}_2$$

$$S_3 - - - - - x_5 = \bar{x}_3$$

$$S_4 - - - - - x_1 = \bar{x}_4$$

$$S_5 - - - - - x_5 = \bar{x}_5$$

$$\bar{X} \approx \text{GD}\left(\mu, \frac{\sigma^2}{n}\right)$$

Chebyshev's Inequality.

if $X \sim \text{GD}(\mu, \sigma)$.

X = Random variable
follows GD

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

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

$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99\%$$

if $Y \not\sim \text{GD}$

Y = Random variable
not following GD

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

$K = 1, 2, 3 \dots$

Let's, $K=2$,

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

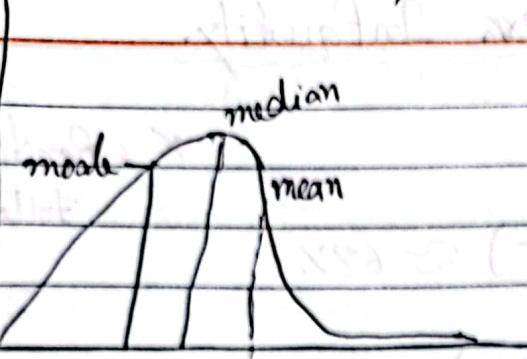
= 75%.

Skewness.

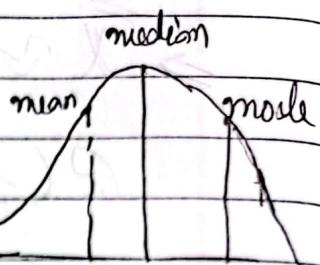
In Probability theory and Statistics, skewness is a measure of the asymmetry of the probability distribution of real-valued random variable about its mean. The skewness value can be +ve, -ve, 0.

$\Rightarrow \text{mean} > \text{median} > \text{mode}$

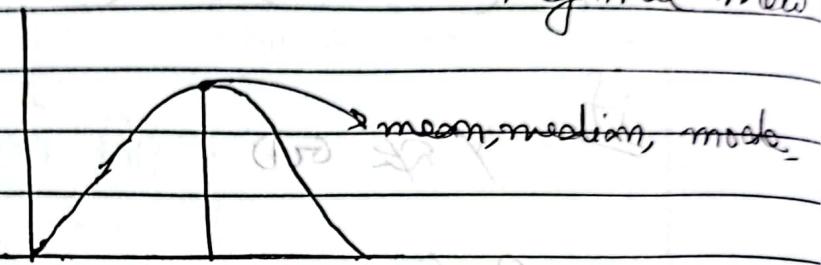
$\Rightarrow \text{mode} > \text{median} > \text{mean}$



Right skew or
Positive skew



Left skew or
Negative skew

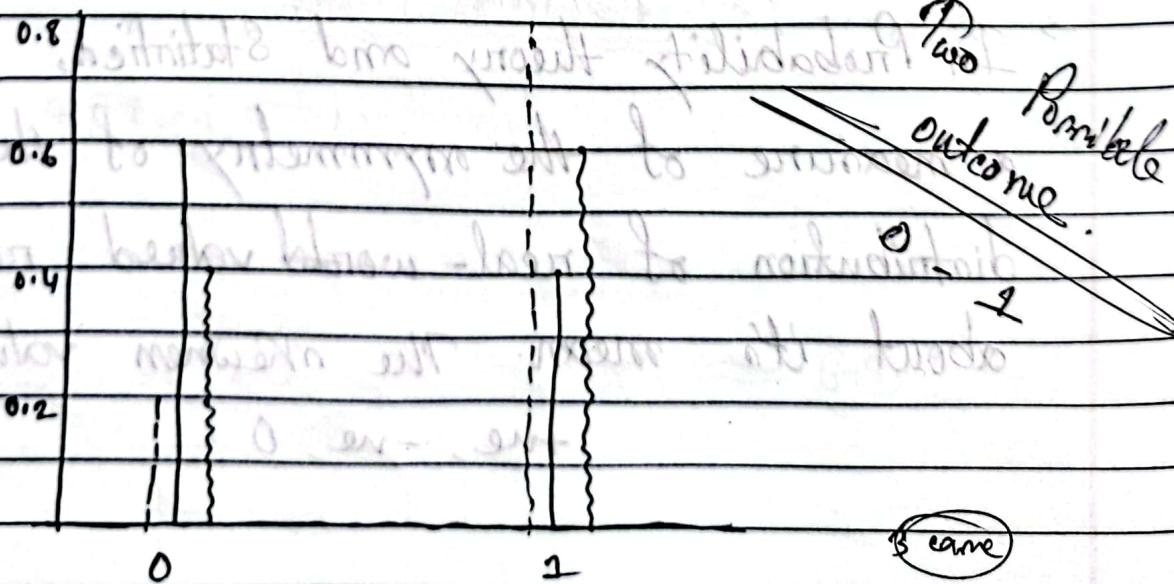


Symmetrical
Distribution.

$\Rightarrow \text{mean} \approx \text{median} \approx \text{mode}$.

Bernoulli Distribution.

1 Probability mass function.



$$\rightarrow P(n=0) = 0.2 \text{ and } P(n \geq 1) = 0.8$$

$$\rightarrow P(n=0) = 0.6 \quad " \quad P(n=1) = 0.4$$

$$\rightarrow P(n=0) = 0.4 \quad " \quad P(n=1) = 0.6$$

$$PMF = \begin{cases} p & \text{if } n=1 \\ q = (1-p) & \text{if } n=0 \end{cases}$$

$$p^k (1-p)^{1-k} = P(X=k)$$

$$\therefore P(X=x) = p^n (1-p)^{1-x} \rightarrow \text{Discrete value.}$$

outcome continuous value \rightarrow pdf \rightarrow Probability Density function.

$$x = 0 \text{ or } 1.$$

$$P(\text{Success}) = p$$

$$P(\text{fail}) = 1-p = q$$

$$P(X=0) = p^0 (1-p)^{1-0} = 1-p = q$$

$$P(X=1) = p^1 (1-p)^0 = p$$

"Interview Question"

Mean, variance, std.

$$\textcircled{1} \text{ Mean. } E(x) = \sum_{n=0}^1 n \cdot P(n) \quad n=0, 1$$

$$P(n=0) = 0.4$$

$$P(n=1) = 0.6$$

$$= (0 \times 0.4) + (1 \times 0.6) \cdot$$

$$= 0.6$$

$$= p$$

$$\textcircled{1} \text{ Variance } = p(1-p) = pq \quad (0=10) \quad 9$$

$$\textcircled{2} \text{ Std } = \sqrt{pq} \quad (1=10) \quad 9$$

$$\textcircled{3} \text{ if } q > p \text{ median } = 0 \quad 9 \quad \} = 7.149$$

$$p = q \text{ median } = 0.5$$

$$p > q \text{ median } = 1 - (q-1) \quad 9 = (x=x) \quad 9$$

~~Mode?~~

$$9 = (6.666666666666667) \quad 9$$

Confidence Intervals

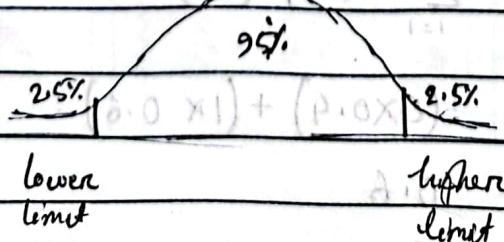
→ Point Estimate $\bar{x} \rightarrow \mu \rightarrow$ Population mean.

→ Range of values for an unknown parameter.

→ 95% Confidence level intended.

lower limit higher limit

→ Average size of sharks in the sea.



take sample. $\frac{C_U}{C_L}$

Suppose

Population std $\sigma = 100$

C.I = Point estimate \pm margin error.

sample $n = 30$

$$\bar{x} = 500$$

$$= \bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$= 500 \pm 2.05 \cdot \frac{100}{\sqrt{30}}$$

α Significance value

\rightarrow Z-table

$$= 5\% = 0.05$$

$$1 - 0.025 \leftarrow z_{0.025}$$

$$= 0.975$$

tail point.

$$= 500 \pm 1.96 \times \frac{100}{\sqrt{30}}$$

$$\text{Lower limit} = 500 - 1.96 \times \frac{100}{\sqrt{30}} = \approx 386$$

$$\text{Upper limit} = 500 + 1.96 \times \frac{100}{\sqrt{30}} = \approx 613$$

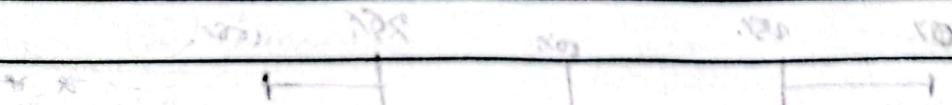
$\alpha \rightarrow 95\%$

$$\rightarrow 386 \leftrightarrow 613$$

when α is not given
apply t-test.

$$(SMP) z_{1-\alpha} = 1.9$$

$$(SMP) z_{\alpha/2} + 1.9$$



5 Number Summary and handle outliers using

IQR. → git hub → code bank.

z-test

5

① Minimum → 0%.

Percentile

② 25%.

1, 2, 3, 4, 5, 6, 7, 8, 9, 10

③ Median → 50%

Percentile = n% * (n+1)

④ 75%.

$$25\% = \frac{1}{4} * (13)$$

$$= \frac{1}{4} * \frac{13}{1}$$

⑤ Maximum → 100%.

$$50\% = \frac{13}{2} = 6.5 \rightarrow$$

handle Outliers → ① IQR. → Percentile

② z-test

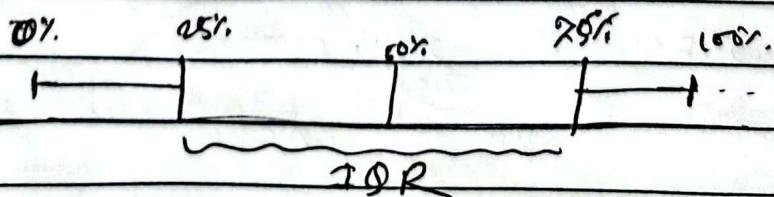
③ σ - standard deviation

$$IQR = Q_1 - Q_3$$

$$\text{Low range} = Q_1 - 1.5 \times IQR$$

$$\text{High } = Q_3 + 1.5 \times IQR$$

} out of this range outlier.



↓
outlier.

P-Value

- It is the probability for the "Null hypothesis" to be true.
- A p value used in "Null hypothesis" to help you support or reject the null hypothesis. The p value is the evidence against null hypothesis. The "smaller" the p-value, the stronger the evidence, that you should "reject" the null hypothesis.

Null hypothesis: Treats everything same, or equal.

H_0 = The coin is fair.

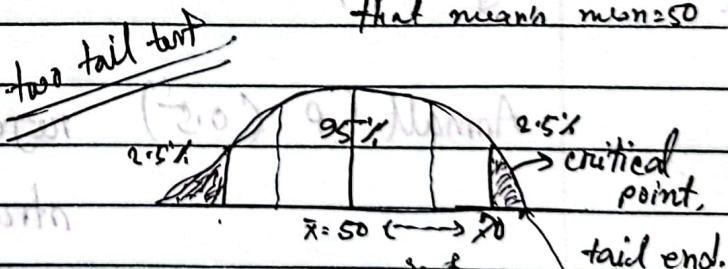
Run experiment 100 times

H_1 = The coin is not fair.

650 times (head).

$$P = 0.05$$

Significance value



(55) times

(30) times

Null hypothesis true

so we have fair coin.

away from the mean.

not fair

reject null hypothesis.

if $p > .10 \rightarrow$ no significant

if $p \leq .10 \rightarrow$ marginally significant.

if $p \leq 0.5 \rightarrow$ Significant.

if $p \leq 0.1 \rightarrow$ highly significant.

P value vs Alpha level.

→ Alpha levels are controlled by the researcher and related to confidence level. You set an alpha level by subtracting your confidence level from 100%. For example, if you want to be 98% confident in your research, the alpha level would be 2%. When you run the hypothesis test, the test will give you a value for p .

If small ($p < 0.5$) reject the null hypothesis. This is strong evidence that the null hypothesis is invalid.

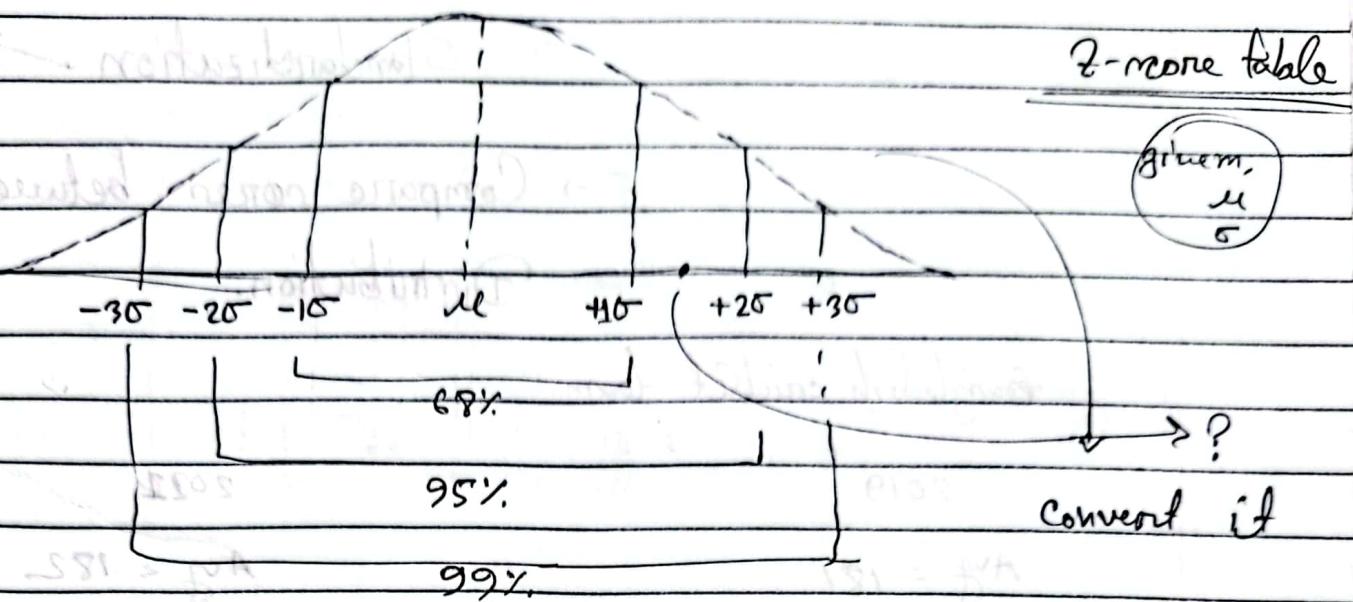
If large (> 0.05) means the alternative hypothesis is weak, so we do not reject the H_0 .

mittlere Varianz unter 9 gleichverteilten Werten

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

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

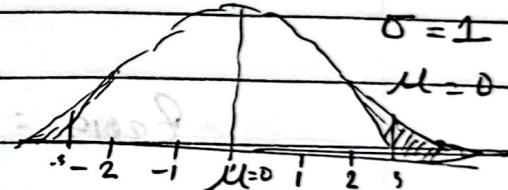
Z-Score



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

SND
Standard Normal Distribution.

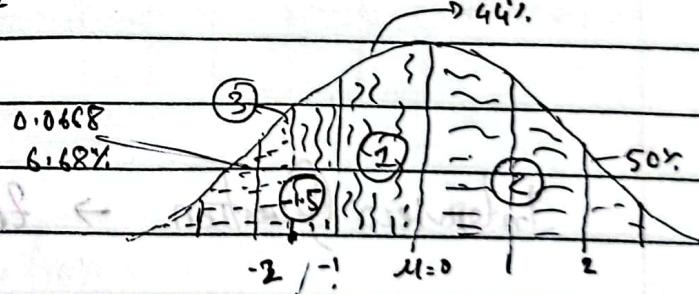
1, 2, 3, 4, 5



$\mu = 3$

$z = -2, -1, 0, 1, 2$

$\sigma = 1$



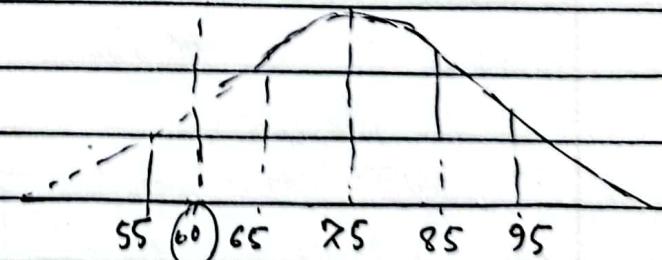
lets Problem $\mu = 25$

$\sigma = 10$

$\rightarrow P(X > 60)$

(a) problem for binomial dist.

fix value around 1,



$$100 = x + 50\% + 6.68$$

$$x = 100 - 56.68$$

$$\approx 44$$

Z -score - Application \rightarrow [Outlier remove]

\rightarrow Standardization $\xrightarrow{\text{same scale}}$

\rightarrow Compare scores between different distribution.

Bangladesh cricket team.

2019

$$\text{Avg} = 181$$

$$\sigma = 12$$

Total score = 182
final

NZB

2021

$$\text{Avg} = 182$$

$$\sigma = 5$$

Total score = 185
final

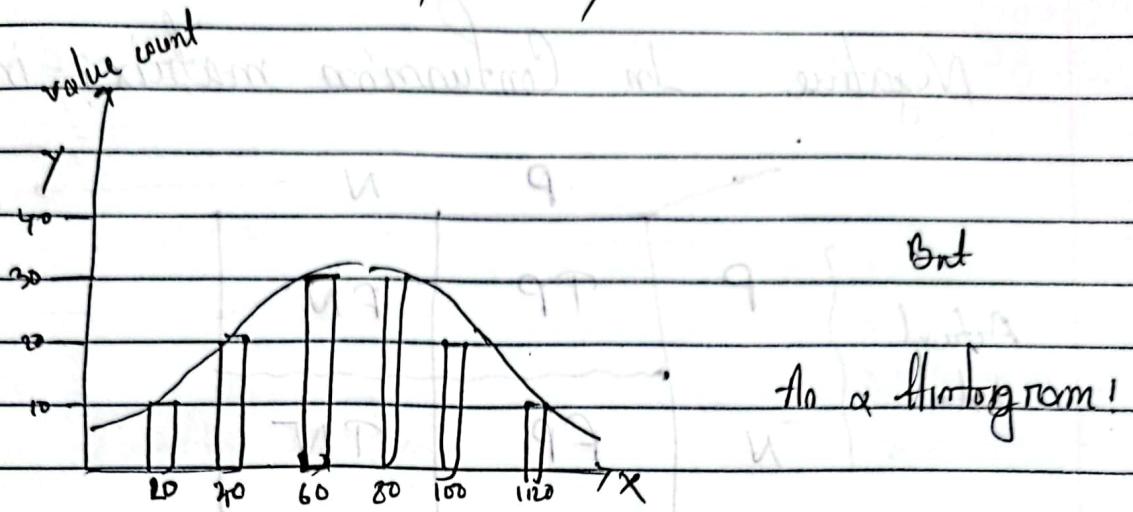
$$Z_{2019} = \frac{182 - 181}{12} = 0.5$$

$$Z_{2021} = \frac{185 - 182}{5} = 0.6$$

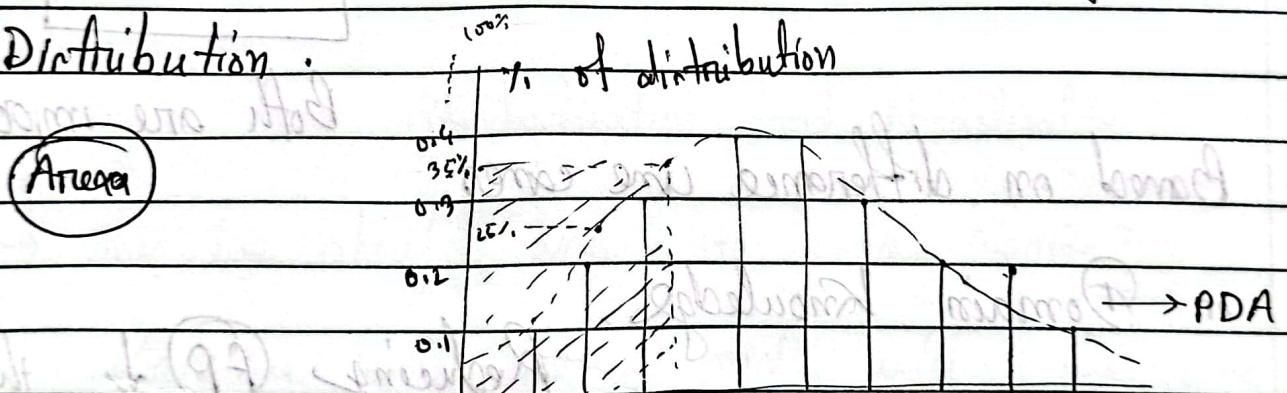
Better Performance

Interview Question \rightarrow Z-score - Application.

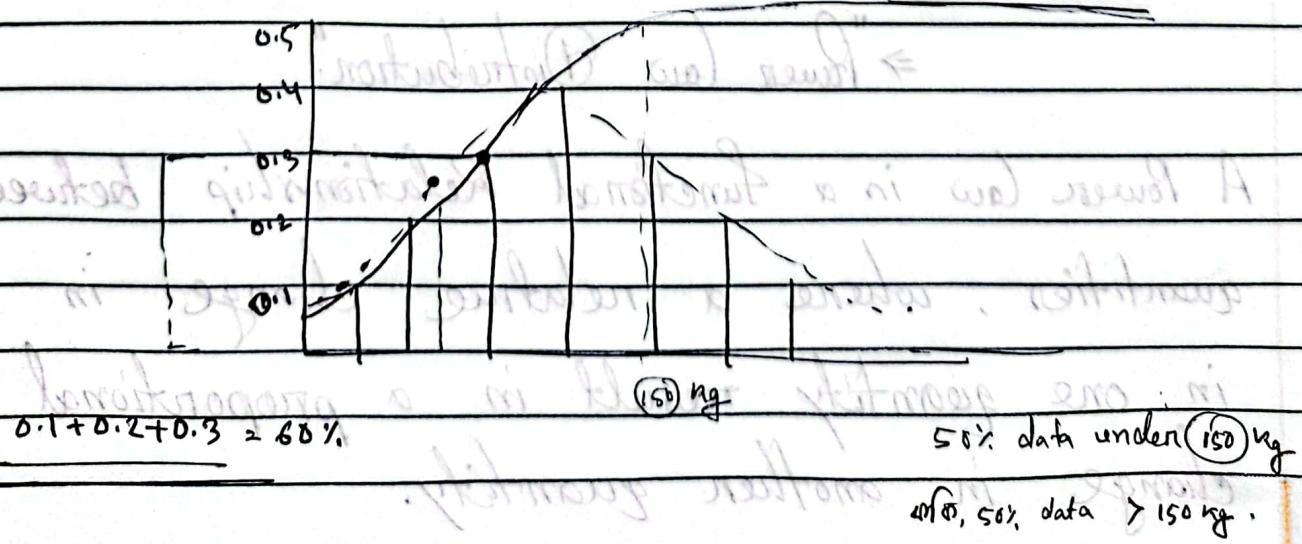
⇒ Probability Density Function.



But in a Probability Density Function. in Y curve axis
not a value of count, there's in percentage of
Distribution.



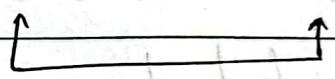
Cumulative Density Function.



Whether we should Reduce False Positive or
Negative In Confusion matrix? - interview Q?

		P	N	Predicted.	
Actual data	P	TP	FN		
	N	FP	TN		

$$\text{TP} \uparrow \quad \text{TN} \uparrow \quad \text{FN} \downarrow \quad \text{FP} \downarrow$$



Both are important.

Based on difference line comes

on Domain Knowledge

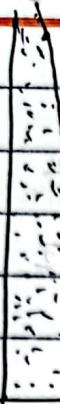
Reducing $\text{FP} \downarrow$ then

$\text{FN} \downarrow$

\Rightarrow "Power Law Distribution".

A Power law is a functional relationship between two quantities, where a relative change in one quantity result in a proportional relative change in another quantity.

"Strong & balanced business is"



Pareto Principle

80/20 rule.

80% result get
20% effort.

Interview Question Example.

Main focus \rightarrow 20% Effort.

get 80% result.

Standardization and Normalization.

\rightarrow why we should use std... and Norm.

\rightarrow when we use std... and when we use Norm.

Brief - toward main net.

most important topics

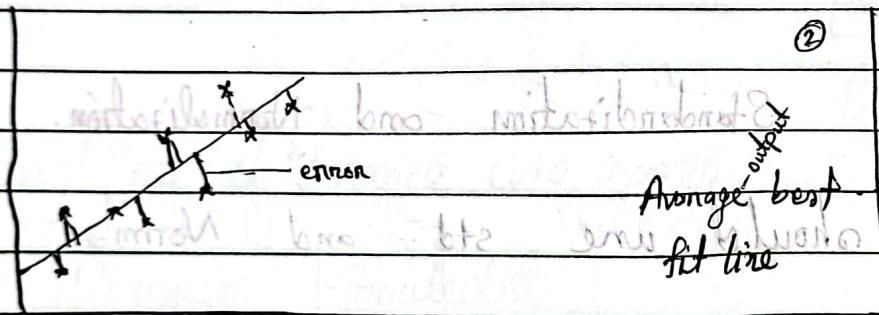
"R Square" and "Adjusted R Square"

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

① SS_{res} = Residual or Error (Sum of square)

② SS_{tot} = Sum of Average total

①



$$SS_{\text{res}} = \text{Sum} (y - \hat{y})^2$$

$$SS_{\text{tot}} = \text{Sum} (y_i - \bar{y}_{\text{avg}})^2$$

$R^2 < 0$, when my best fit line worse than than the average

$$SS_{\text{res}} > SS_{\text{tot}} = > 1 \text{ and get -ve.}$$

best fit line.

$R^2 \rightarrow$ check goodness of best fit line.

⇒ Follow Krish Naik's and
Analyticvidhya article

$$\text{Adjusted } R^2 = \left\{ 1 - \left[\frac{(1-R^2)(n-1)}{n-k-1} \right] \right\}$$

when,

① added feature not correlated

$R^2 \downarrow$

$k \uparrow$

to the target value.

$(n-k-1) \downarrow$

$\left[- \frac{1}{n-k-1} \right] \uparrow$

Adjusted R^2 decrease

② added feature correlated to the

target variable.

when

$R^2 \uparrow$

and $(n-k-1) \downarrow$

$(1-R^2)(n-1) \downarrow$

⇒ Adjusted R^2 Increase

$\text{Adj } R^2 < R^2$

Analyticvidhya article →

R^2 vs Adjusted R^2

Difference?

Hypothesis testing.

→ Evaluates 2 Mutual Exclusive statement, } → P values

↓
on

Population

Using Simple Data.

→ T test

→ Anova test

→ Z-test

→ F-test

Step:

example

① Make Initial Assumption (H_0 - Null hypothesis)

Criminal Case

② Collecting data [Info, Evidence]

H_0 } Innocent

H_1 } Guilty

③ Gather Evidence to Re check.

Reject or Not Reject Null

hypothesis.

↓ (1-n) (n-1) on Confusion matrix

	H_0	H_1	Truth
H_0	Do not Reject (except H_0)	OK	Type II error
H_1	Reject	Type I Error	OK

Type I and Type II Errors play major role.

Gender	Age	Weight	Height	
M	Teen	20	1.4	$P \leq 0.05$
F	Adult	65	1.2	
M	Adult	65	1.4	Significance
M	child	20	1	value.
F	Adult	25	1.3	
M	Teen	80	1.3	

Two categories (Gender, Age)

One categorical

H_0 --- There is no difference

H_1 There is difference

Proportion in Gender

Test → Chi-Square test

Test → One Sample Proportion

Test

$P \leq 0.05$ Reject H_0

Accept H_1

height mean test $\rightarrow 1.3$

based on previous sample this mean is difference or not?

One Continuous feature

Two Continuous feature

H_0 ---

H_1 ---

Test \rightarrow T-Test

Test \rightarrow Correlation.

Population mean vs
Sample mean

Test \rightarrow F-test

(Let's) One Numerical or Continuous feature :

One or two categorical feature

(many) categories.

Then we apply more than two

Anova test

→ if only two categories.

Suppose one numerical
one categorical
and in categorical variable

And we have two categories

we apply

T-test

→ Chi-Square test

① Chi-Square Test for Independence

② Chi-Square test for Goodness of fit.

Statlectures

YouTube