

$X \sim \text{Log Normal Distribution } (\mu, \sigma)$



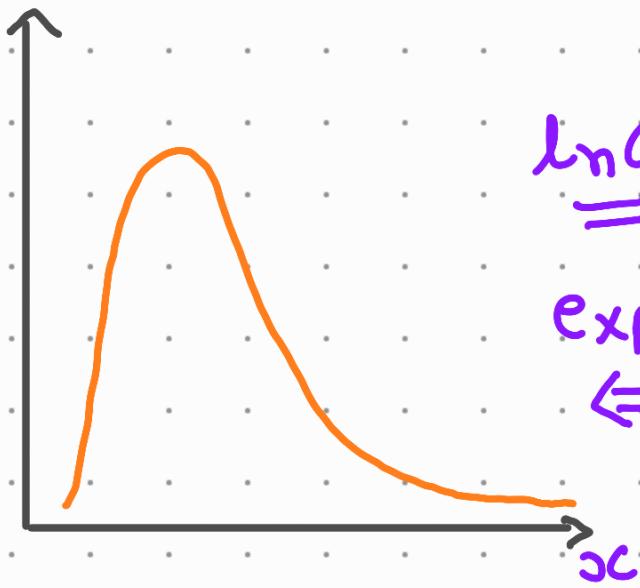
$Y \sim \ln(X) = \text{Normal Distribution}$



[\log_e]

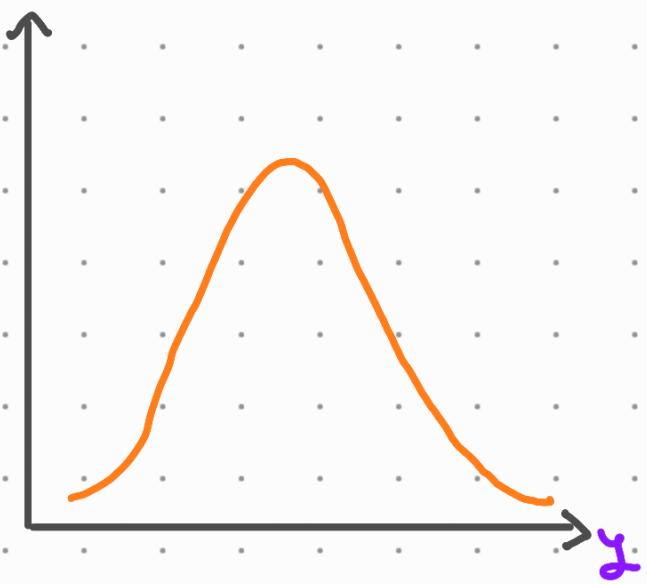


$X \approx \exp(Y) = \text{Log Normal Distribution}$



Log Normal Distribution

$$\begin{array}{c} \ln(x) \\ \Rightarrow \\ \exp(y) \\ \Leftarrow \end{array}$$

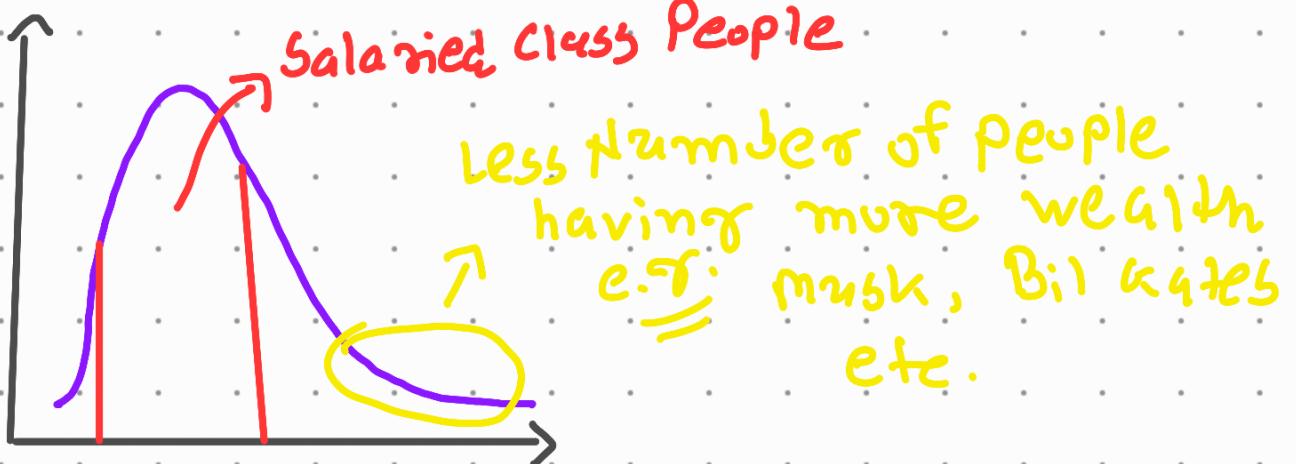


Normal Distribution



QQ Plot

e.g. ① The wealth distribution of the world

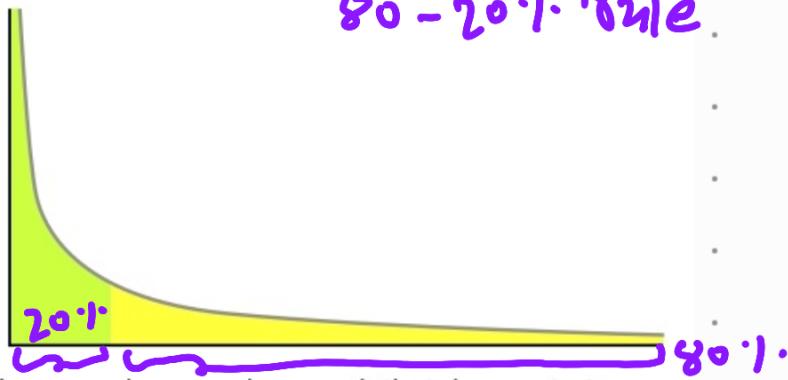


- ② Discussion Forum → length of the comments
- ③ Length of Chess Game
- ④ Dwell time on Online articles (joke, news)

⑧ Power Law Distribution

In statistics, a power law is a functional relationship between two quantities, where a relative change in one quantity results in a relative change in the other quantity proportional to a power of the change, independent of the initial size of those quantities: one quantity varies as a power of another.

80-20 rule



An example power-law graph that demonstrates ranking of popularity. To the right is the long tail, and to the left are the few that dominate (also known as the 80-20 rule).

e.g.

① → 20% of Team is responsible for winning 80% of matches.

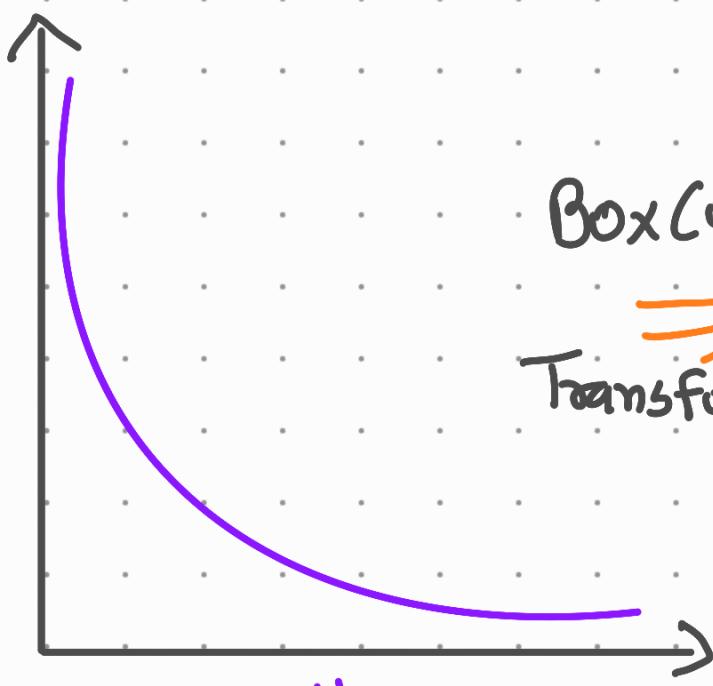
② → 80% of wealth are distributed with 20% of the total population.

③ 80% of the total oil is with 20% of the nation

④ Frequencies of words in most languages.

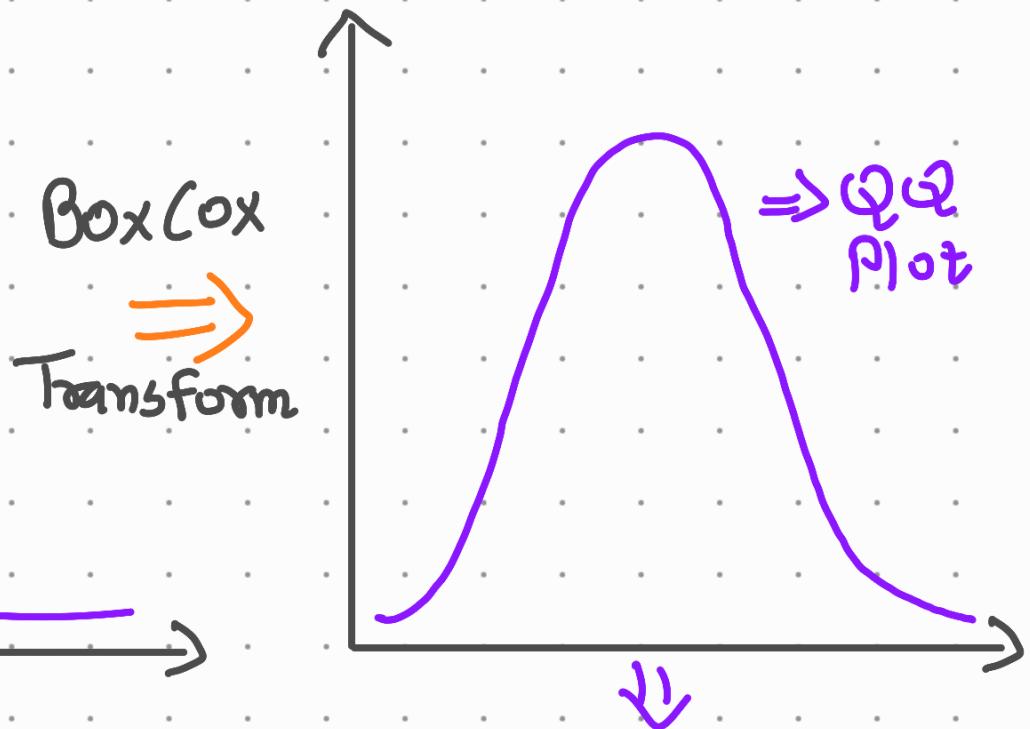
(5) 20% of the major defects fixes the 80% of the upcoming defects in a software product.

→ Distribution Graph



↓
Power Law
Distribution

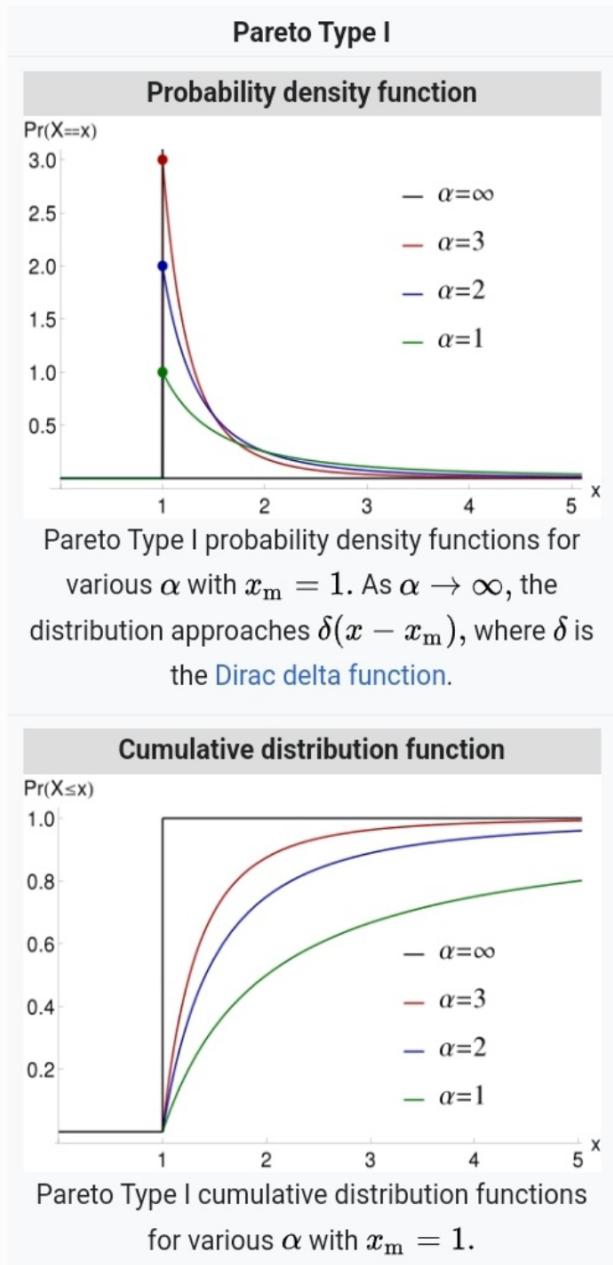
↓
Pareto
Distribution



↓
Normal Distribution

9) Pareto Distribution

The Pareto principle or "80-20 rule" stating that 80% of outcomes are due to 20% of causes was named in honour of Pareto, but the concepts are distinct, and only Pareto distributions with shape value (α) of $\log 45 \approx 1.16$ precisely reflect it. Empirical observation has shown that this 80-20 distribution fits a wide range of cases, including natural phenomena and human activities.



Parameters $x_m > 0$ scale (real)
 $\alpha > 0$ shape (real)

Support $x \in [x_m, \infty)$

PDF
$$\frac{\alpha x_m^\alpha}{x^{\alpha+1}}$$

CDF
$$1 - \left(\frac{x_m}{x}\right)^\alpha$$

Quantile
$$x_m (1 - p)^{-\frac{1}{\alpha}}$$

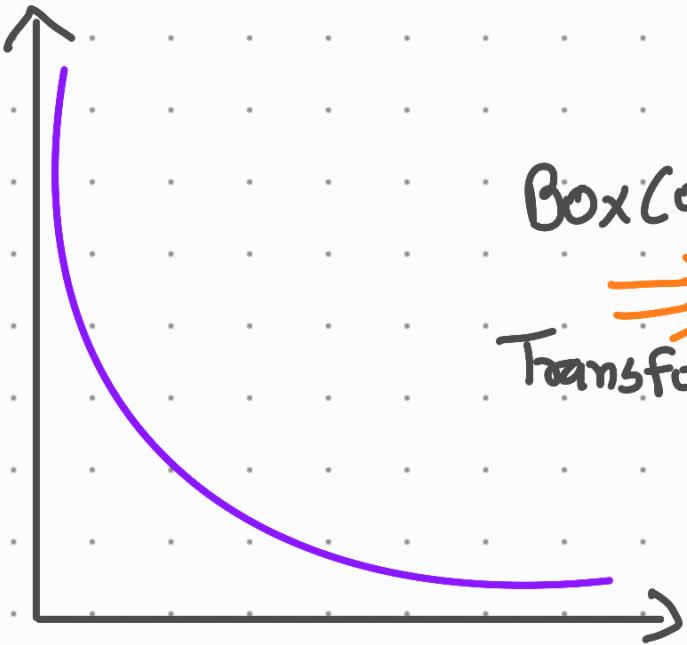
Mean
$$\begin{cases} \infty & \text{for } \alpha \leq 1 \\ \frac{\alpha x_m}{\alpha - 1} & \text{for } \alpha > 1 \end{cases}$$

Median
$$x_m \sqrt[{\alpha}]{2}$$

Mode x_m

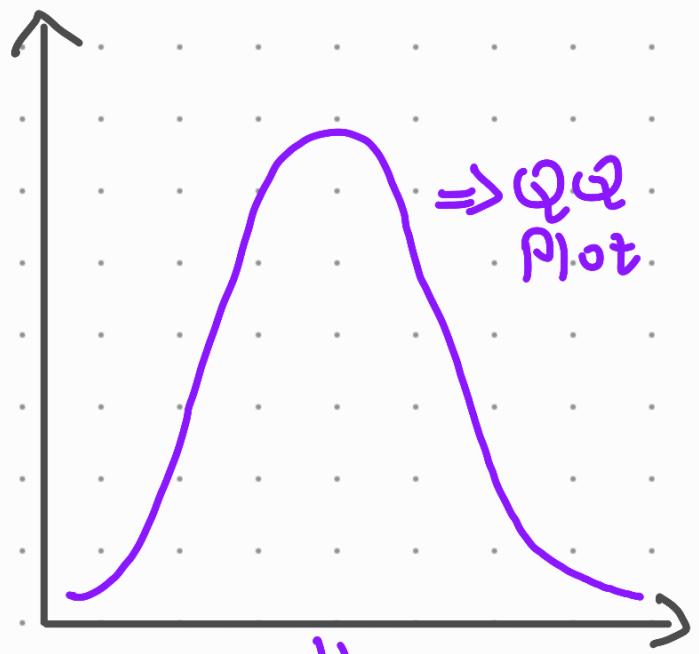
Variance
$$\begin{cases} \infty & \text{for } \alpha \leq 2 \\ \frac{x_m^2 \alpha}{(\alpha - 1)^2 (\alpha - 2)} & \text{for } \alpha > 2 \end{cases}$$

Skewness
$$\frac{2(1 + \alpha)}{\alpha - 3} \sqrt{\frac{\alpha - 2}{\alpha}} \quad \text{for } \alpha > 3$$



Box Cox
Transform

↙
Pareto
Distribution



⇒ Q-Q
Plot

↙ Normal Distribution

e.g.
//

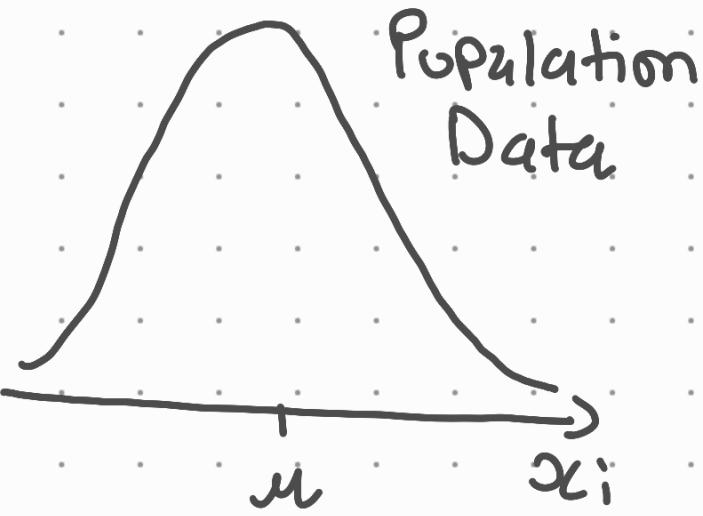
- ① 80% of the entire project is done by 20% of the team.
- ② 80% of defects can be solved by 20% of the main defects.

Central Limit Theorem

The central limit theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large number of samples taken from a population.

The central limit theorem says that the sampling distribution of the mean will always be normally distributed, as long as the sample size is large enough. Regardless of whether the population has a normal, Poisson, binomial, or any other distribution, the sampling distribution of the mean will be normal.

$$\textcircled{1} \quad \bar{x} \sim N(\mu, \sigma)$$



$n = \text{Sample Size} \Rightarrow f_{\bar{x}_n}$

$$S_1 = \{x_1, x_2, x_3, \dots, x_n\} = \bar{x}_1$$

$$S_2 = \{x_1, x_3, \dots, x_n\} = \bar{x}_2$$

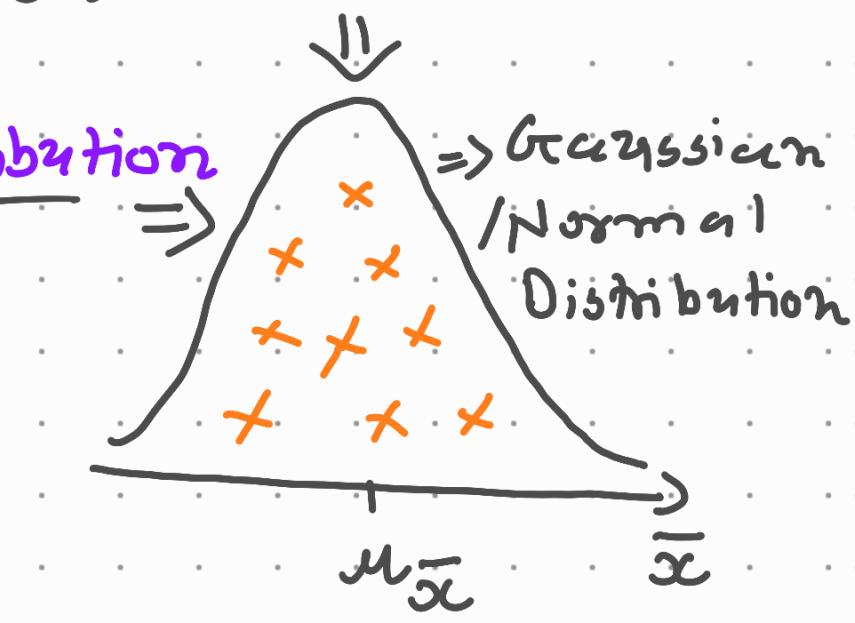
$$S_3 \quad \dots \quad - \quad -$$

$$S_4 \quad \dots \quad - \quad -$$

$$\vdots$$
$$S_m \quad \dots \quad - \quad -$$

$$\bar{x}_m$$

Sampling Distribution of Mean



② $x \notin N(\mu, \sigma)$

$n \geq 30 \Rightarrow$ Sample size

s_1

\bar{x}_1

s_2

\bar{x}_2

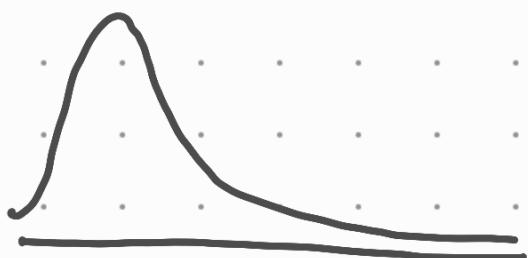
:

:

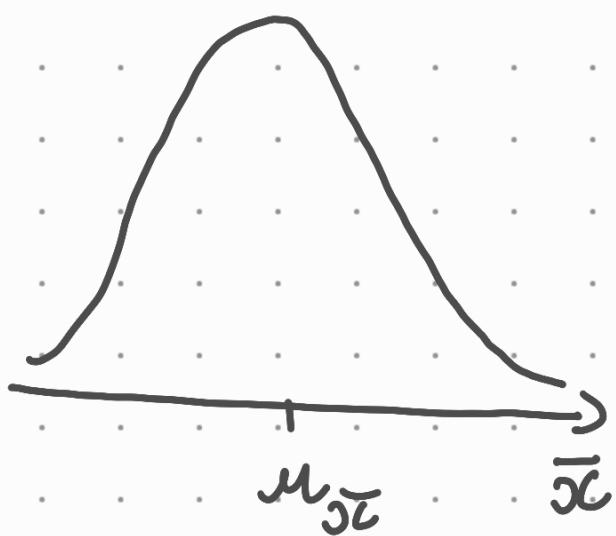
:

s_m

\bar{x}_m



↓ CLT



① Normal Distribution

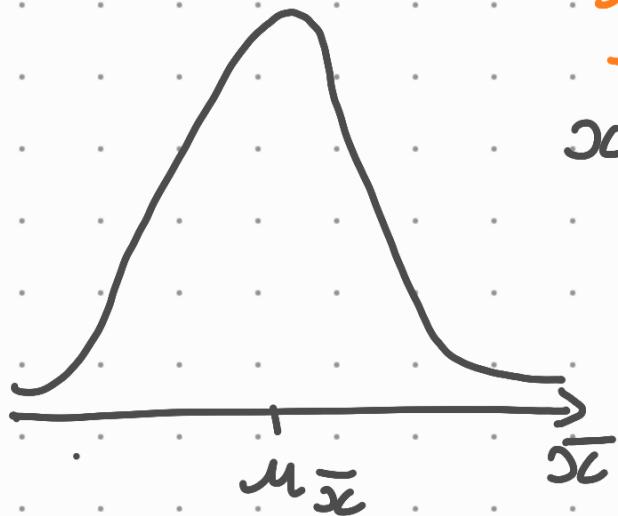
$\overline{\quad}$ $\overline{\quad}$



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



Sampling Distribution of Mean



σ = Population S.D.

μ = Population Mean

n = Sample Size

$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

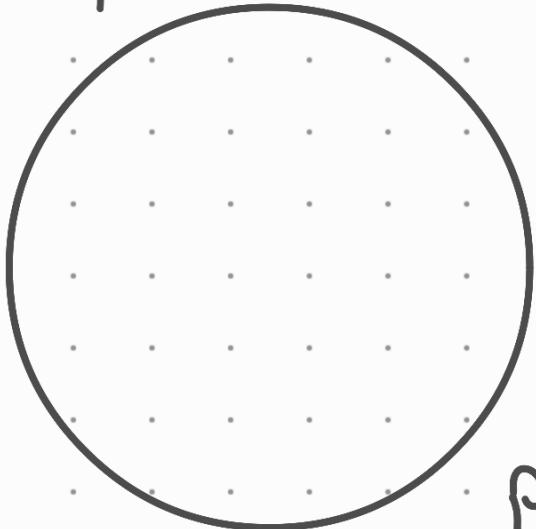
Estimates

It is specified observed numerical value used to estimate an unknown population parameter.

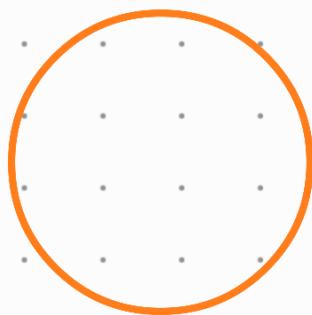
① **Point Estimate**: Single numerical value used to estimate an unknown population parameter.

e.g. Sample mean is a point estimate of population mean.

Population



Sample



Point Estimate

65μ ← \bar{x} 60

② Interval Estimate: Range of values is used to estimate unknown population parameters.

55 - 65

↓
CI

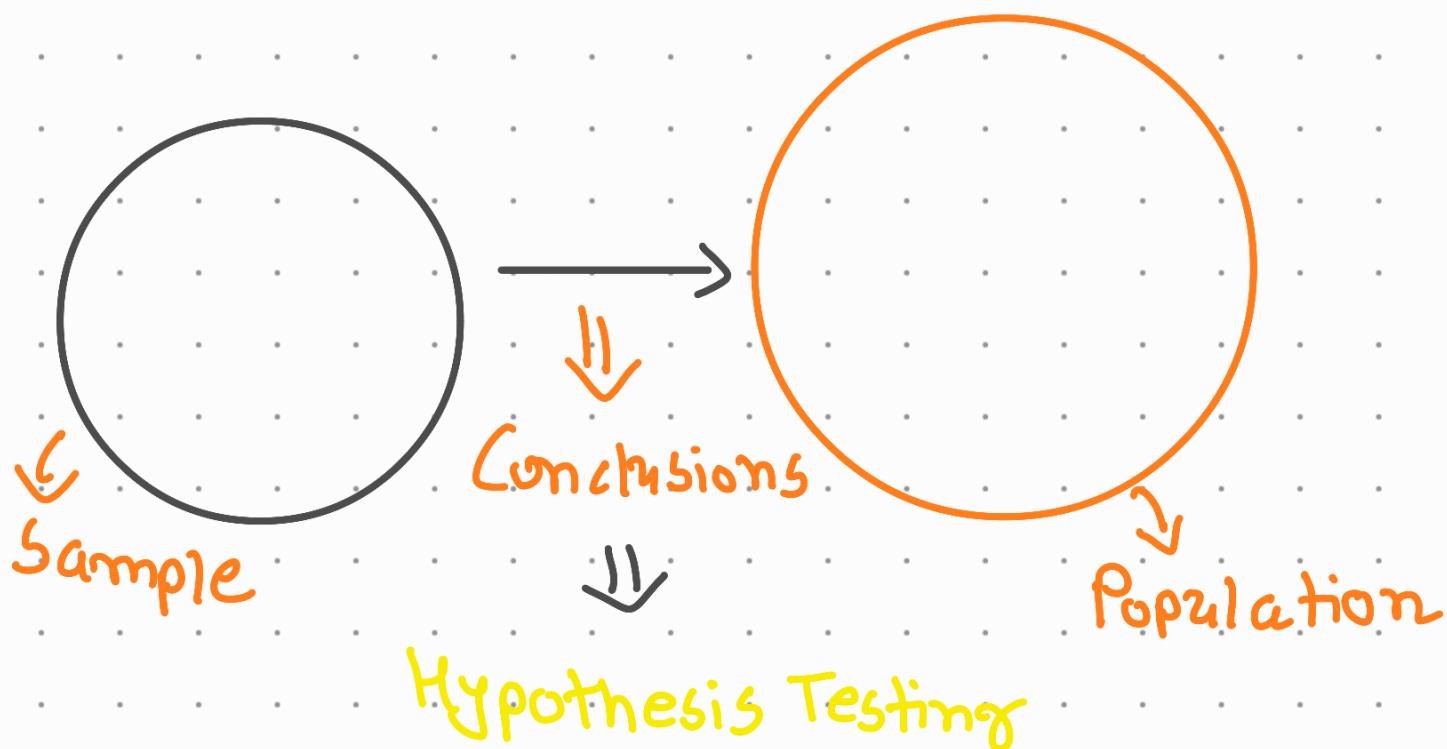


Point Estimate

Confidence Interval

A Hypothesis And Hypothesis Testing Mechanism

→ Inferential Stats : Conclusion or Inference



→ Hypothesis Testing mechanism

① Null Hypothesis (H_0) - The person is not guilty

→ The assumption you are beginning with

② Alternative Hypothesis (H_1) - The person is guilty

→ Opposite of Null Hypothesis

③ Experiments - Statistical Analysis

→ Collect Proof (e.g. DNA, Finger Print)

④ Accept the null Hypothesis or
Reject the null Hypothesis

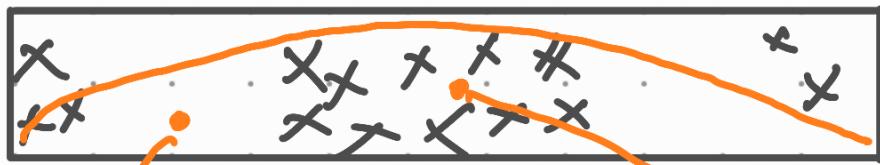
Eg: Colleges at District A stats its ^{Average} pass percentage of Students are 85%. A new college opened in the district and it was found that a sample of student 100 have a pass percentage of 90% with a standard deviation of 4%. Does this school have a different pass percentage.

Ans Null Hypothesis (H_0) = $\mu = 85\%$.

Alternate Hypothesis (H_1) = $\mu \neq 85\%$.

P-Value

A p-value, or probability value, is a number describing how likely it is that your data would have occurred under the null hypothesis of your statistical test



Space Bar $\rightarrow P\text{-value} = 0.8$

$P\text{-value} = 0.2 \Rightarrow$ Out of 100

touches on space bars, we touch around 20 times in this region.

Hypothesis Testing

e.g. Coin is Fair or Not {100 Tosses}

$$P(H) = 0.5 \quad P(T) = 0.5$$

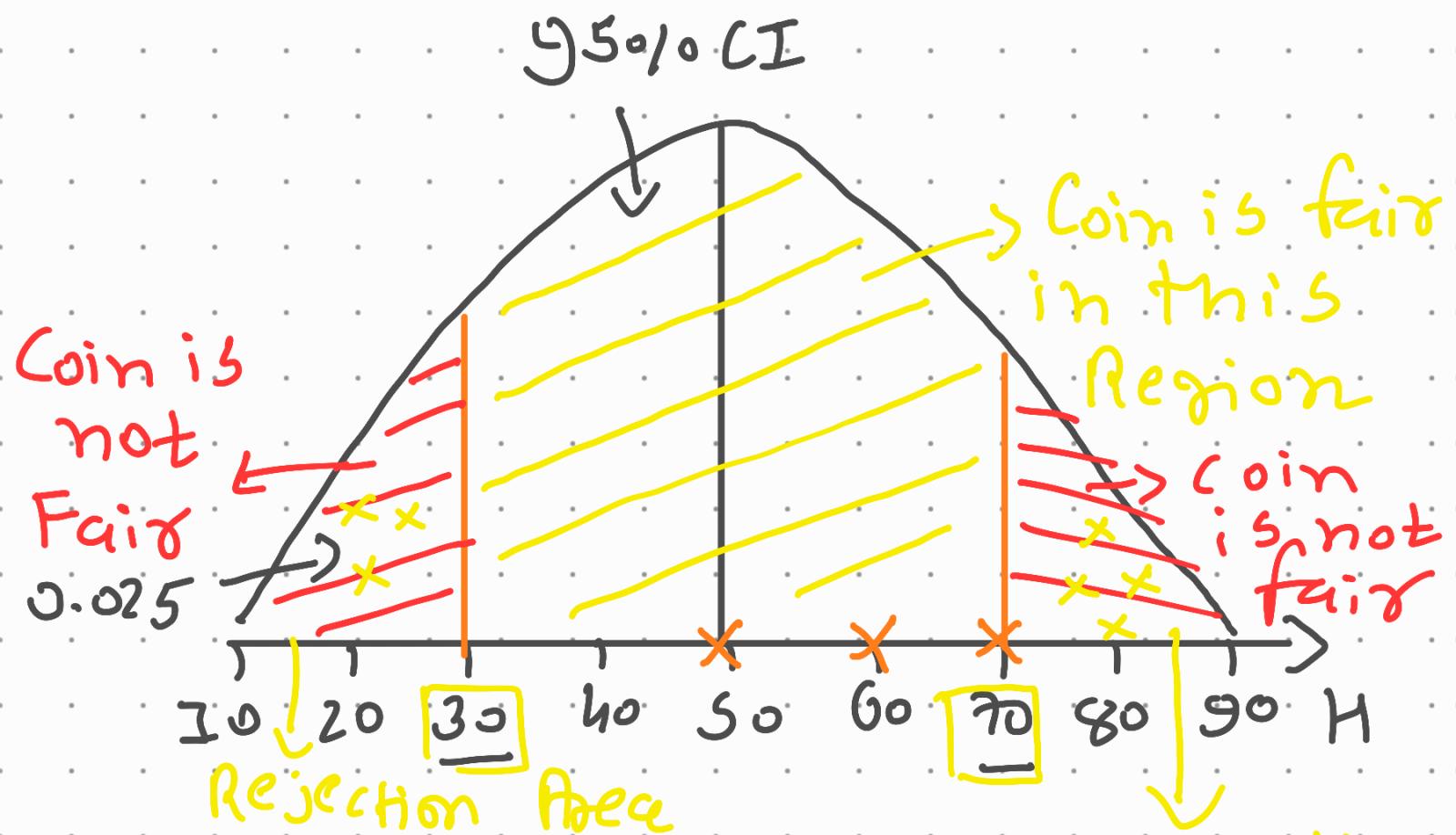
$$\text{if } P(H) = 0.6 \quad P(T) = 0.4$$

$$\text{if } P(H) = 0.7 \quad P(T) = 0.3$$

① Null Hypothesis: H_0
↳ The Coin is Fair

② Alternate Hypothesis: H_1
↳ The Coin is Not Fair

③ Experiment: 100 times



④ Significance Value Value Rejection Area

$$\rightarrow \alpha = 0.05$$

$$\rightarrow CI = 1 - 0.05 = 0.95$$

→ $P < \text{Significance Value}$

Reject the Null Hypothesis

→ Else

Fail to Reject the Null Hypothesis

Statistical Analysis

① Z-Test

$\bar{y} \Rightarrow \text{Average}$

② T-Test

\Downarrow
Z-Table \rightarrow Z-score & P-Value
T-Table

③ CHI SQUARE Test \Rightarrow Categorical Data

④ ANNOVA Test \Rightarrow Variance of Data

① Z-Test

Conditions To Use

i) Population Std

ii) $n \geq 30$

With a $\sigma = 3.9$

1) The average heights of all residents in a city is 168cm. A doctor believes the mean to be different. He measured the height of 36 individuals and found the average height to be 169.5cm.

(a) State null and Alternate Hypothesis

(b) At a 95% confidence level, is there enough evidence to reject the null hypothesis.

$$\rightarrow \mu = 168 \text{ cm}$$

$$\sigma = 3.9$$

$$n = 36 \geq 30$$

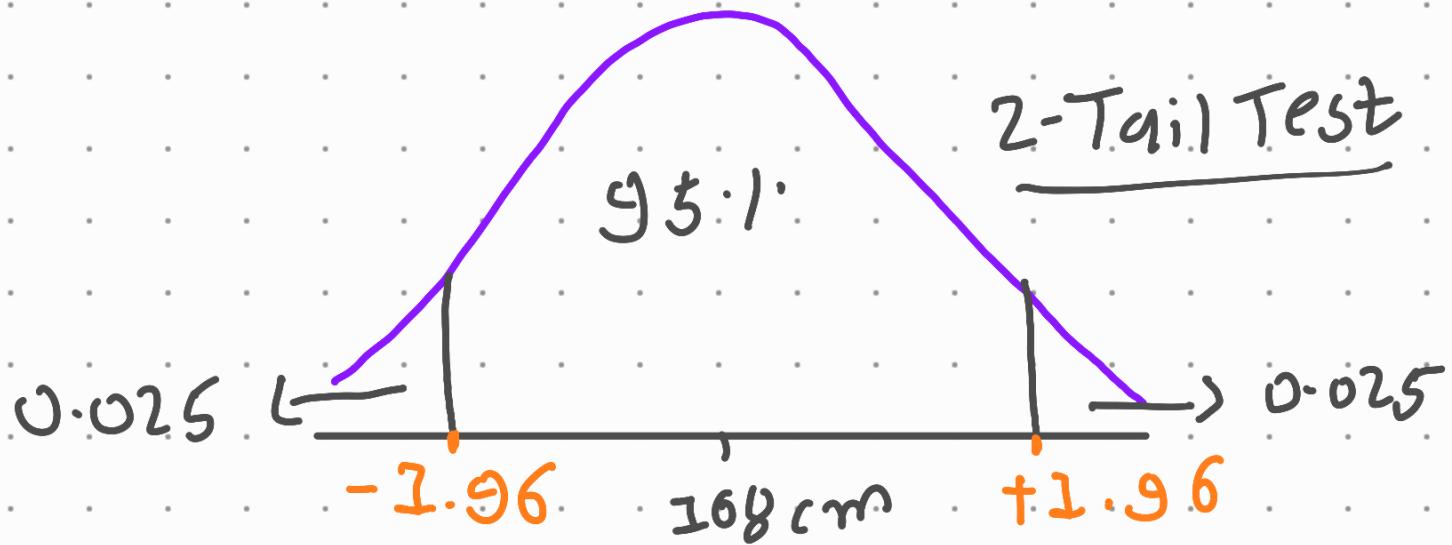
$$\bar{x} = 169.5 \text{ cm}$$

$$CI = 0.95 \quad \alpha = 1 - 0.95 = 0.05$$

(a) Null Hypothesis $H_0: \mu = 168 \text{ cms}$

(b) Alternate Hypothesis $H_1: \mu \neq 168 \text{ cms}$

(c) Based on C.I. Draw decision boundary



$$\begin{aligned} &= 1 - 0.95 \\ &= 0.05 \quad \left\{ \Rightarrow \text{Z-score} = 1.96 \right. \\ &\quad \text{(find from Z-table)} \end{aligned}$$

↑
Area

→ if Z is less than -1.96 or greater than $+1.96$, Reject the NULL HYPOTHESIS.

→ Z-Test

$$Z_d = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \Rightarrow \text{Sample Data}$$

$$= \frac{169.5 - 168}{3.9 / \sqrt{36}}$$

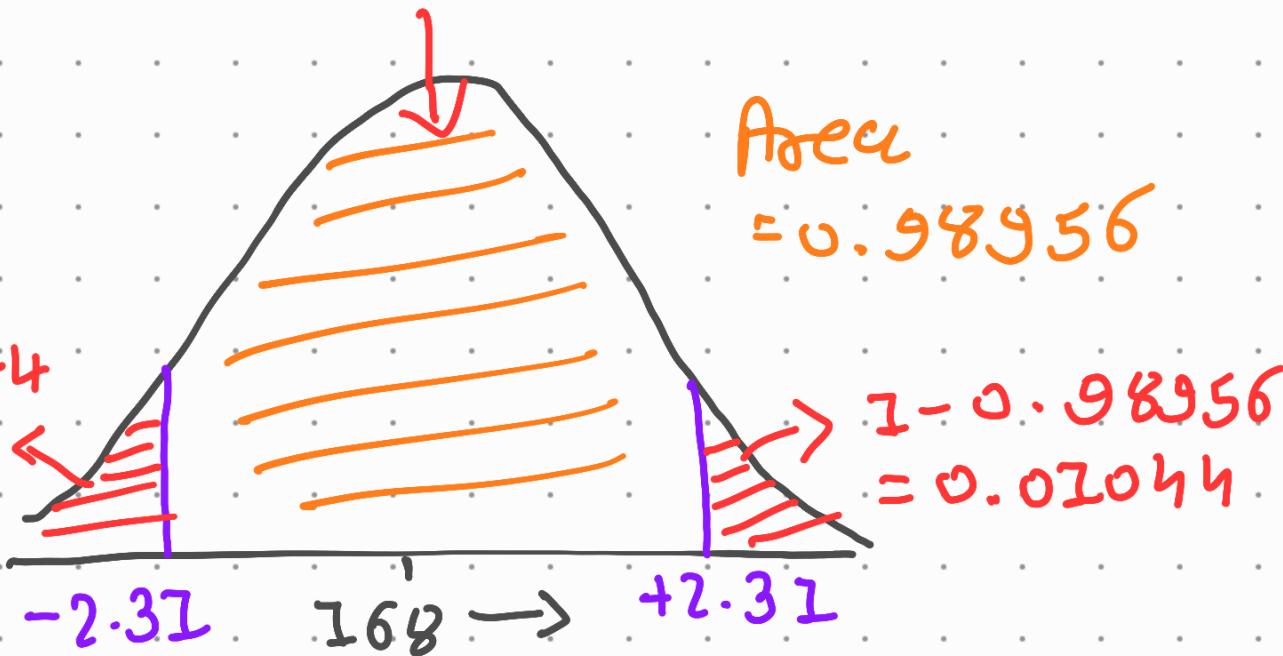
$$= \frac{1.5}{0.65}$$

$Z_d = 2.31 > 1.96$ so we reject

the Null Hypothesis

$$\underline{P < 0.05}$$

$$A_{dec} = 0.97922$$



$$\rightarrow P\text{-value} = 0.01044 + 0.01044 \\ = 0.02088$$

$$P < 0.05$$

\rightarrow So we reject the NULL Hypothesis.

Final Conclusion

\rightarrow The average height $\neq 168$ cms.
The average height seems to be increasing based on the sample data.

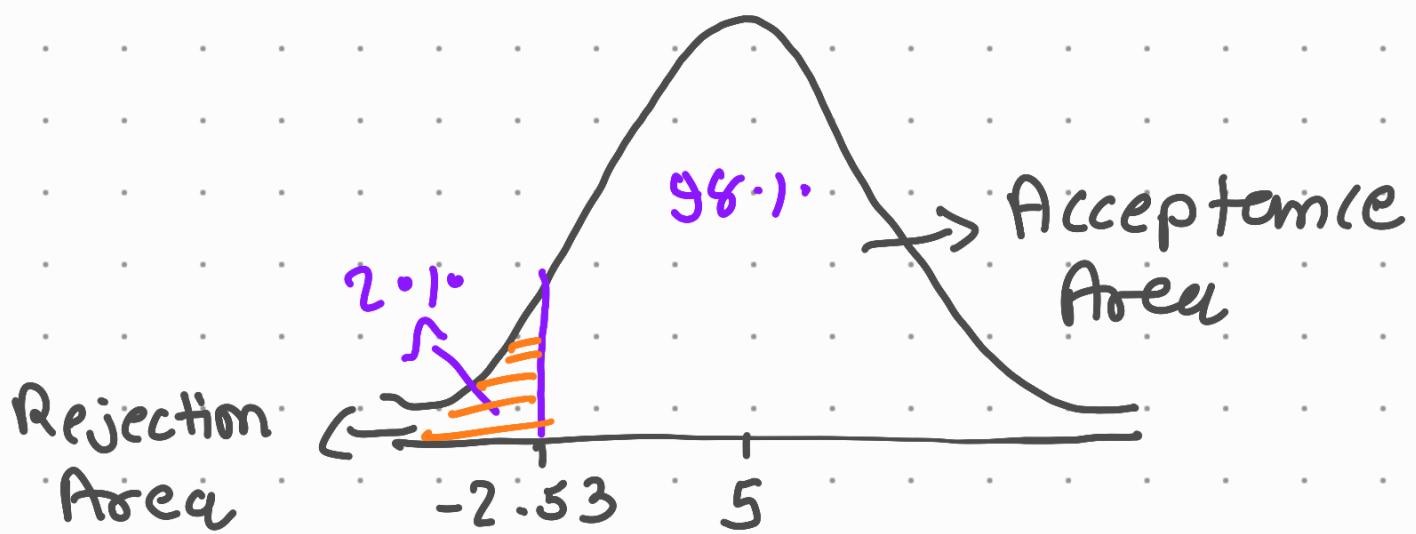
② A factory manufactures bulbs with an average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 years. He tests a sample of 40 bulbs and finds the average time to be 4.8 years.

- (a) State null and alternate hypothesis
- (b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised?

$$\rightarrow \sigma = 0.50 \quad \mu = 5 \\ n = 40 \quad \bar{x} = 4.8$$

- (a) Null Hypothesis $H_0: \mu = 5$
- (b) Alternate Hypothesis $H_1: \mu < 5$
{1-tail test}

(c) Decision Boundary



Z-Test

$$Z_d = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$
$$= \frac{4.8 - 5}{0.50 / \sqrt{40}}$$

$Z_d = -2.53 \Rightarrow$ Area under the Curve
 $= 0.00570$

P-Value = 0.00570

\rightarrow Compare P-Value with Significance Value

Significance Value $\alpha = 0.02$

$0.00570 < 0.02 \Rightarrow$ True

\rightarrow We reject the Null Hypothesis.

Student T Distribution

- In Z-stats when we perform any analysis using Z-score, we require σ (Population Standard Deviation) → is already known.
- How do we perform analysis when we don't know the population standard deviation?

↓
we use Student T Distribution

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \Rightarrow$$

$$t = \frac{\bar{x} - \mu}{S / \sqrt{n}}$$

→ t-test
or
t-stats

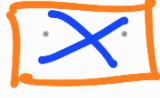
S = Sample Standard Deviation

→ We will use t-table.

Degree of Freedom

$$dof = n - 1$$

3 People \Rightarrow



\rightarrow First will have 3 choices where to seat

\rightarrow Second will have 2 choices where to seat $\rightarrow dof = 3-1=2$

\rightarrow But First will have to seat on the spot remaining $\rightarrow dof = 3-2=1$

T-Stats With T-test Hypothesis Testing

(One Sample T-test)

① In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? ~~one tailed~~

$$\rightarrow CI = 95\% \quad n = 30$$

$$\alpha = 0.05 \quad \bar{x} = 140$$

$$\mu = 100 \quad s = 20$$

(a) Null Hypothesis $H_0: \mu = 100$

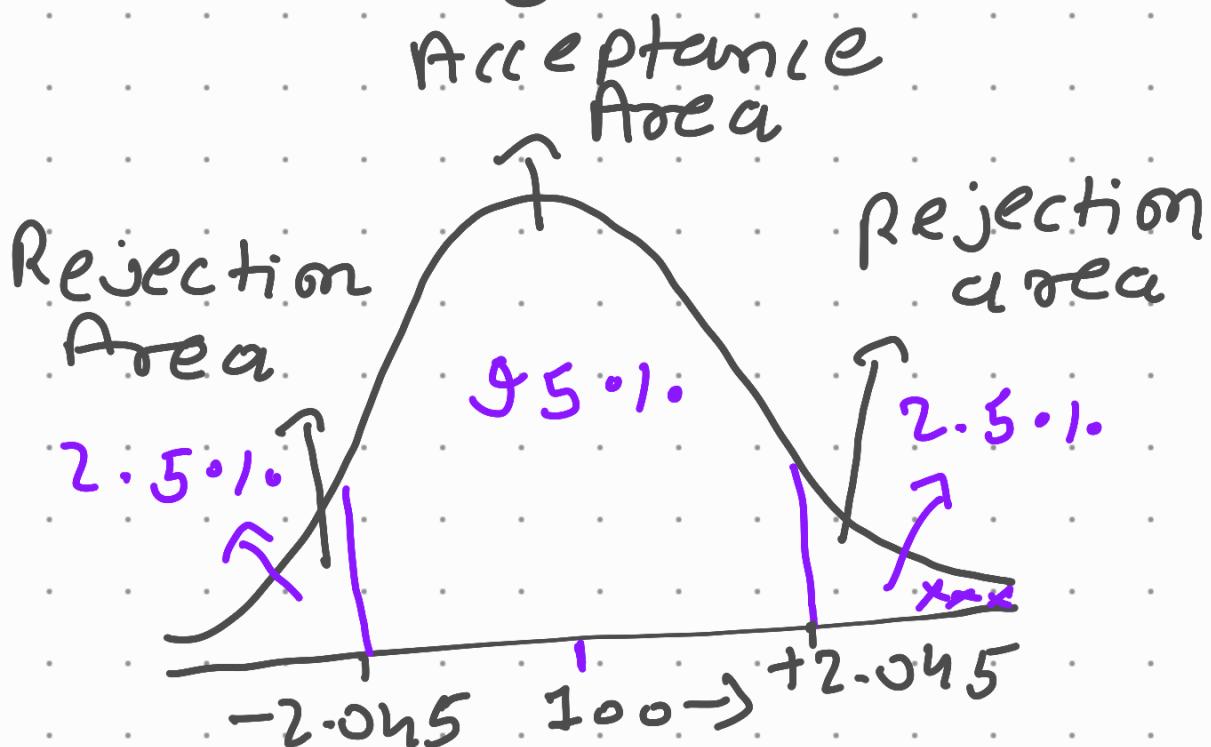
(b) Alternate Hypothesis $H_1: \mu \neq 100$

{2-Tail Test}

(c) Degree of Freedom

$$dof = n - 1 = 30 - 1 = 29$$

(d) Decision Rule



→ if t -test is less than -2.045 or greater than +2.045 we reject the Null Hypothesis.

(c) Calculate Test Statistics

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$
$$= \frac{140 - 100}{20/\sqrt{30}}$$

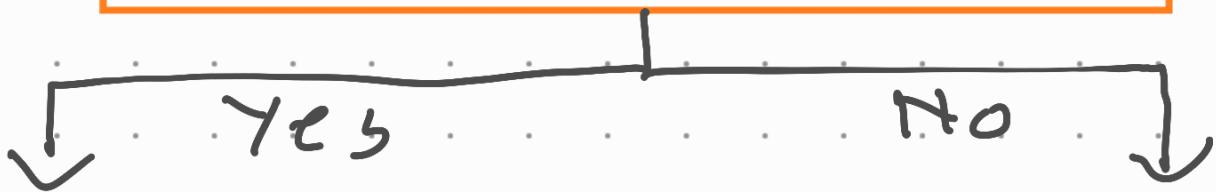
$$t = 20.96$$

→ Since $t = 20.96 > +2.045$, so Reject The Null Hypothesis

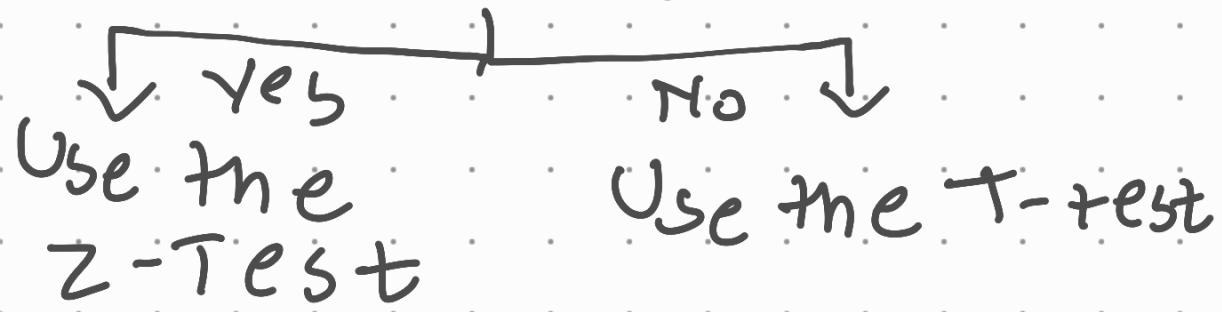
→ Conclusion: medication used has affected the intelligence. It has increased the IQ.

When To Use T-test Vs Z-Test

Do you know the
Population Standard
Deviation



Is the sample size
above 30? ($n > 30$)



Type-1 and Type-2 Errors

Reality : Null Hypothesis is True or
Null Hypothesis is False

Decision : Null Hypothesis is True or
Null Hypothesis is False

Outcome-1 :- We reject the Null Hypothesis when in reality it is False \rightarrow Good

Outcome-2 :- We reject the Null Hypothesis when in reality it is True \rightarrow Type-I Error

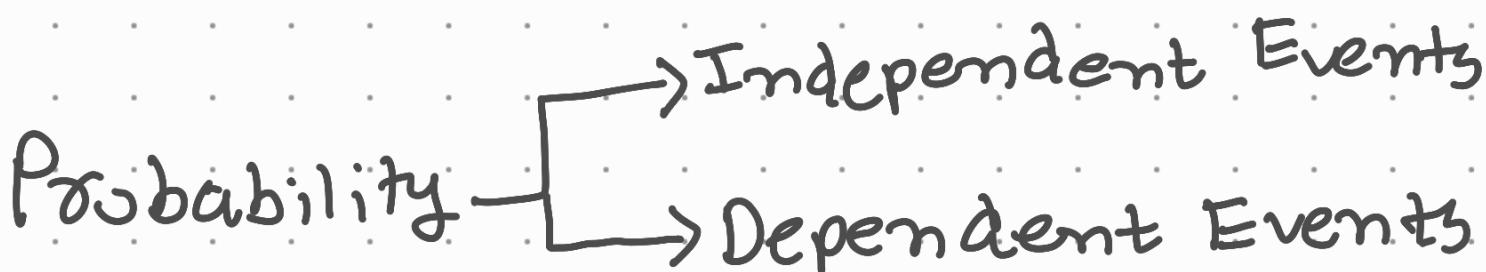
Outcome-3 :- We retain the Null Hypothesis, when in reality it is False \rightarrow Type-II Error

Outcome-4 :- We retain the Null Hypothesis, when in reality it is True \rightarrow Good

Bayes Statistics (Bayes Theorem)

Bayesian statistics is an approach to data analysis and parameter estimation based on Bayes' theorem.

Bayes' Theorem



① Independent Events

e.g. Rolling a dice

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

$$P_3(1) = 1/6$$

$$\vdots \quad \vdots$$

$$P_3(6) = 1/6$$

② Tossing a coin

$$P_S(H) = 0.5$$

$$P_S(T) = 0.5$$

② Dependent Events

e.g. ① 3 yellow, 2 red marbles in a bag.

$$P_S(\text{Get one red}) = 2/5 \text{ after this}$$

$$P_S(\text{Yellow}) = 3/4$$

$$\rightarrow P_S(\text{Red and Yellow})$$

$$= P_S(\text{Red}) * P_S\left(\frac{\text{Yellow}}{\text{Red}}\right)$$

↳ Conditional Probability

$$\rightarrow P_S(A \text{ and } B) = P_S(B \text{ and } A)$$

$$= P_S(A) * P_S(B/A)$$

$$= P_S(B) * P_S(A/B)$$

$$P_B(B|A) = \frac{P_B(B) \times P_A(A|B)}{P_B(A)}$$

→ Bayes' Theorem

$$P_A(A|B) = \frac{P_A(A) \times P_B(B|A)}{P_B(B)}$$

A, B = Events

$P_A(A|B)$ = Probability of A given B is True

$P_B(B|A)$ = Probability of B given A is True

$P_A(A), P_B(B)$ = Independent Probabilities of A and B