

- ① Type of statistics  
 Descriptive → Tells where the problem is  
 & for eg. how many century Kohli scored  
 Predictive → look at data & tell me what might  
 happen. finding relationship among  
 object & predict  
 \* for eg. how many century Kohli will score in WC  
 Prescriptive → action design to do something  
 \* for eg. how do we get Kohli to score more  
 century in WC
- Descriptive → Describe the data without doing  
 any prediction on it  
 telling you the way it is

\* Standard deviation s.d  $\sqrt{\frac{1}{n-1} (x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}$

measure of spread  $\bar{x} = \frac{1}{n} (x_1 + \dots + x_n)$  it's  $n-1$  bcoz  
 now std dev → centered diff. is already  
 considered the element  
 in std dev → squared avg.

\* square is done as we want to see both  
 positive & negative deviation. Else it will  
 neutralize it.

Mean absolute deviation (mad)  $\frac{1}{n-1} (|x_1 - \bar{x}| + \dots + |x_n - \bar{x}|)$

\* Distribution  $F(x) = \text{Prob}(X \leq x)$   
 $f(x) = F'(x)$



- \* distribution function is Integral of density fn  
 \* density function is derivative of distribution fn

univariant , bivariant

↓

sense of correlation  
in two variables

(2)

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

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \rightarrow \text{covariance}(x, y) \quad \sigma_{xy}$$

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

↓ variance x

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

$$= \text{variance}(x)$$

$$= [\text{std dev}(x)]^2$$

correlation  $\rightarrow \text{cov}(x, y)$

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \quad -1 \leq \rho \leq 1$$

↓ std dev x      ↓ std dev y

covariance  $\rightarrow$  nature of relationship between  
two Random variable

but not the value of how strong the  
relationship ( becoz value is depend  
on unit of measurement)

correlation  $\rightarrow$  it cancel out the unit of  
measurement due to deviation of std.  
dev. It is a measure of relationship  
between two random variable.

③ variance → measure the dispersion of set of data point around their mean  
std dev →  $\sqrt{\text{variance}}$  → measure for spread of Data

relative std dev → coefficient of variation

$$\frac{\text{std. dev}}{\text{mean}}$$

used for comparing two data set

panda → df.corr() → will give correlation matrix

- be used  
any regression model,  
simply descriptive  
✓ it predicts  
✗ it's prescriptive

Descriptive Stat

univariate

- \* location → mean, median,  $Q_3$ ,  $Q_1$
- \* variation → std dev., variance, Range  
Inter Quartile Range

Quartile

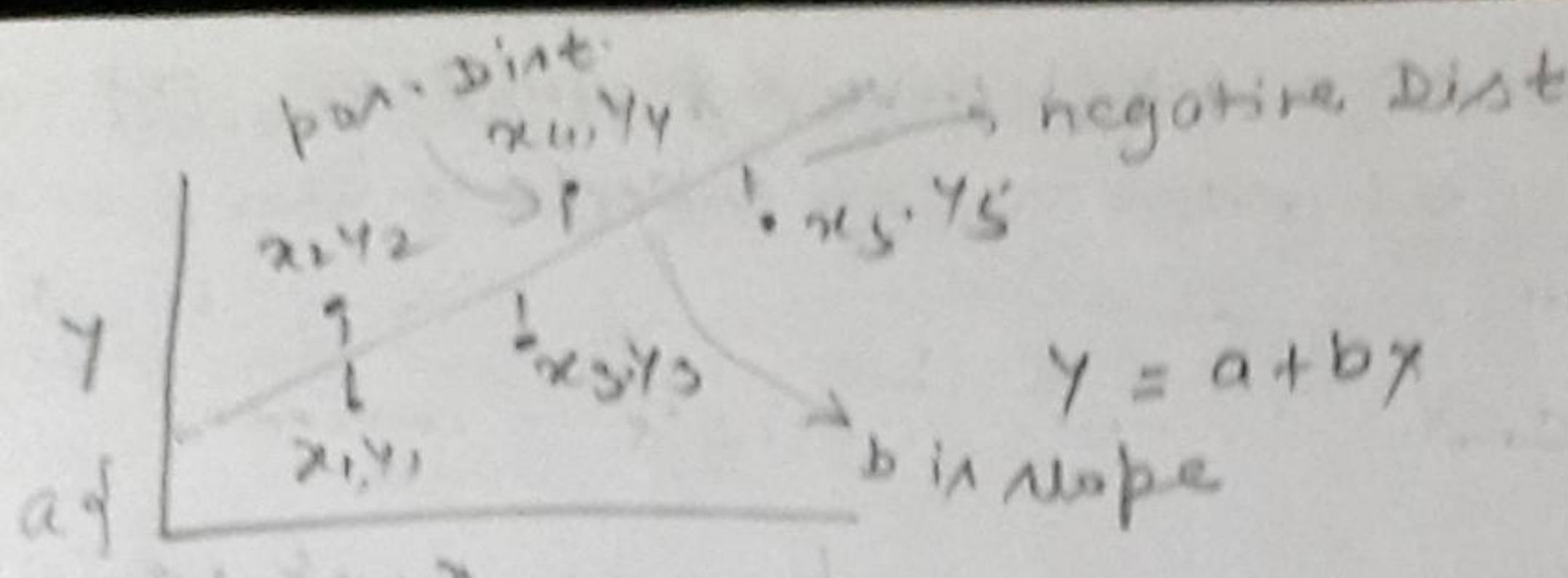
bivariate → covariance, correlation

multivariate → linear Regression

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

it can be used for prediction & forecasting as well. but in Descriptive, it's just to describe Relationship

(4)



$$(y_1 - (a + bx_1))^2 + (y_2 - (a + bx_2))^2 + \dots$$

Tells how far the line

is from the data  $\rightarrow$  the smaller the closer

find the value of  $a, b$  so that ~~above exp.~~ is  
smallest - i.e. choose  $a, b$  to minimize  
above exp.  $\rightarrow$  thin in what linear

Regression  $\Rightarrow$

least square algo

$$\hat{b} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

coef-

$$y = \hat{a} + \hat{b}x$$

$$\hat{a} = \bar{y} - \hat{b}\bar{x}$$

Intercept -

(5) Probability  $\rightarrow$  likelihood of particular event  
 event  $\rightarrow$  outcome of experiment  
 experiment  $\rightarrow$  process performed to understand  
 and observe possible outcome

Set of all outcome of experiment is called  
 "Sample space"

3 component

$$P(D) = 10\%$$

Sample space

G G G	G A D	G D G	D G G
G D D	D G D	D D G	D D D

D  $\rightarrow$  Defective  
 G  $\rightarrow$  Good

4 chances of atleast 2 defective component

$$P(G D D \text{ or } D G D \text{ or } D D G) = P(G D D) + P(D G D) + P(D D G) \xrightarrow{\substack{(\text{OR}) \\ (\text{OR})}} \text{Disjoint}$$


if one happen  
other two can't happen

$$3 (.9 \times .1 \times .1) = \underbrace{P(G) \times P(D) \times P(D)}_{\text{independent event}} \rightarrow P(G) \text{ and } P(D) \text{ and } P(D)$$

when Event are independent  $\rightarrow$  Probability can be multiply

when events are disjoint  $\rightarrow$  probability can be summed

mutually exclusive event  $\rightarrow$  these are not independent but if one occur, other event can't

independent event  $\rightarrow$  one event can't impact other event and both can occur simultaneously.

## Rules for computing Probability

④

1) addition Rule  $\rightarrow$  mutually exclusive events

$$P(A \cup B) = P(A) + P(B)$$

A present occurrence  
 & B or vice versa  
 i.e., A & B can not occur at same time

$P(A \text{ or } B)$

$$P(A \cup B) = P(A) + P(B) - P(A \text{ and } B)$$

this part is 0 in  
 case of mutually  
 exclusive events

addition Rule for  
not mutually exclusive  
events

e.g.

\* Probability of getting King or Queen card

$$P(K) = \frac{4}{52} = \frac{1}{13} \quad P(Q) = \frac{1}{13}$$

$$P(K \text{ or } Q) = \frac{1}{13} + \frac{1}{13} = \frac{2}{13}$$

placed

\* Probability of King or Diamond card

$$P(K) = \frac{4}{52} \quad P(D) = \frac{13}{52} \quad P(D \text{ & } K) = \frac{1}{52}$$

$$P(K \text{ or } D) = \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

(only one card  
 can be King of  
 Diamond)

2) multiplication Rule  $\rightarrow$  independent event

$$P(A \cap B) = P(A) \times P(B)$$

Probability of Intersection

if occurrence of A is in no way influenced by occurrence of B or vice versa

\* if events are not independent

$$P(A \cap B) = P(A) \times P(B/A) \rightarrow \text{if A has occurred}$$

&  $P(B) \times P(A/B) \rightarrow \text{if B has occurred}$

conditional  
prob

$$\textcircled{1} \quad P(B/A) = P(\text{'B given A})$$

find Probability of event A occurring & then find whether event B will occur along with event A

I can also be written as  $\frac{P(A \cap B)}{P(A)}$  that will be a subset of  $P(A)$

e.g. Drawing 2 card from back of card. what is probability of Drawing both cards as spade

let A = getting spade in first draw

let B = getting spade in 2nd draw if first card is not Replaced

(there are 13 spades & 51 cards remaining)

$$P(A) = \frac{13}{52} \rightarrow P(B/A) = \frac{12}{51} \xrightarrow{\substack{\text{dependent} \\ \text{Event as} \\ \text{first card is} \\ \text{not Replaced}}}$$

$$\text{if } \xrightarrow{\substack{\text{independent event}}} \text{first card is replaced } P(B) = \frac{13}{52}$$

it  $P(B/A) \rightarrow$  becoz both card has to be spade  
that possible if first drawn card is already a spade. if yes. what is the probability of 2nd card to be an spade

★ marginal Probability  $\rightarrow$

★ conditional probability

$$\hookrightarrow P(B/A) = \frac{P(A \cap B)}{P(A)}$$

★ joint probability

Family	income less than 10 lakh	income $\geq 10 L$	total
Buyer of car	38	42	80
Non buyer	82	38	120
Total	120	80	200

Q1 → what is prob of randomly selected family buy a car  $\rightarrow \frac{80}{200} \rightarrow$  (marginal prob.)

Q2 → Prob of randomly selected family is both car buyer & income  $\geq 10 L$

$$\frac{42}{200} \text{ (joint prob)}$$

Q3 → randomly selected family belongs to income  $\geq 10$ , find prob of this family buying a car

$$\frac{42}{80} \text{ (as per table)}$$

$$P(\text{car}/\geq 10L) = \frac{P(\text{car} \cap \geq 10L)}{P(\geq 10L)}$$

$$= \frac{42/200}{80/200} = \frac{42}{80}$$

↓  
conditional  
Prob.

conditional prob =	$\frac{\text{joint Prob}}{\text{marginal Prob.}}$
--------------------	---

① Bayes' theorem

$P(\text{spam} \neq \text{word})$

$$\hookrightarrow P(\text{word}) \approx P(\text{spam}/\text{word})$$

$$P(\text{spam}) \approx P(\text{word}/\text{spam})$$

$$P(\text{spam}/\text{word}) = \frac{P(\text{word}/\text{spam}) \approx P(\text{spam})}{P(\text{word})}$$

e.g. of Baye's theorem

$$P(\text{HIV}/+) = \frac{P(\text{HIV and } +)}{P(+)}$$

$$= \frac{P(+/\text{HIV}) * P(\text{HIV})}{P(\text{HIV and } +) + P(\text{not HIV and } +)}$$

$$= \frac{P(+/\text{HIV}) * P(\text{HIV})}{P(+/\text{HIV}) * P(\text{HIV}) + P(+/\text{not HIV}) * P(\text{not HIV})}$$

$$P(\text{HIV}) = 1\%$$

$$P(\text{NHIV}) = 99\%$$

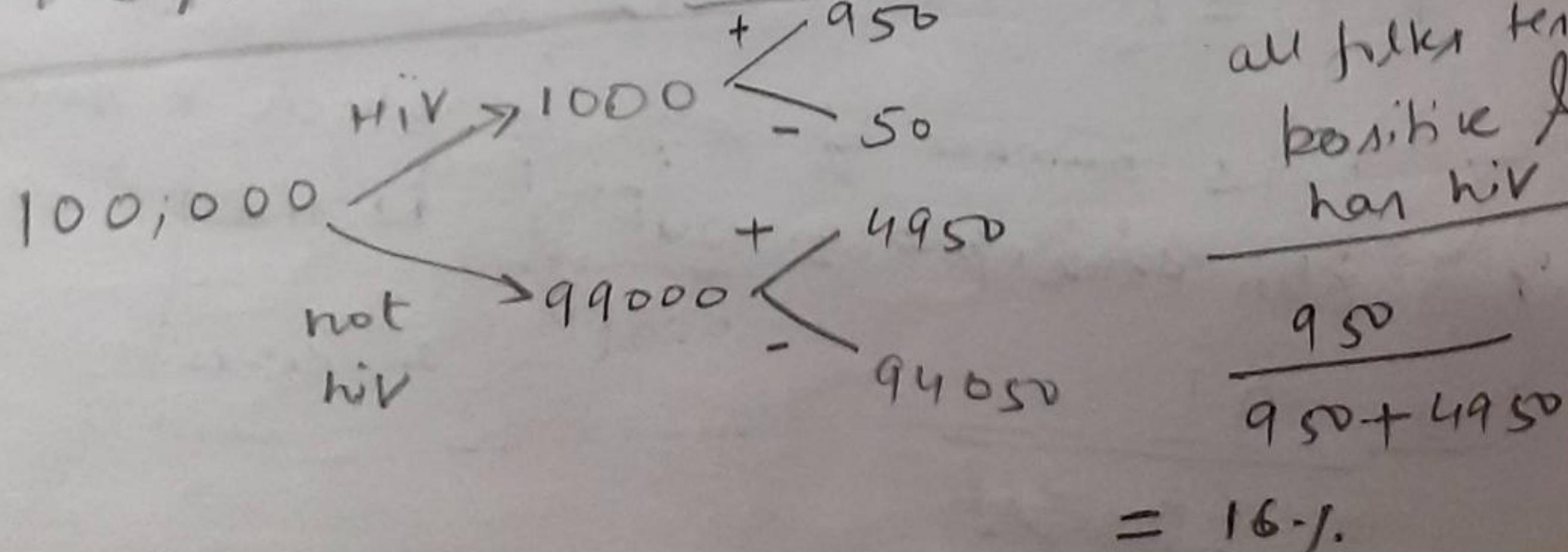
$$P(+/\text{HIV}) = 95\%$$

$$P(-/\text{not HIV}) = 95\%$$

$$\hookrightarrow P(+/\text{not HIV}) = 5\%$$

$$= \frac{.95 * .01}{(.95 * .01) + (.05 * .99)}$$

$$= 0.16$$



19

$$P(\text{spam}) = 30\%$$

$$P(\text{ns} \text{spam}) = 70\%$$

$$P(\text{congratulation} / \text{spam}) = 75\%$$

$$P(\text{congratulation} / \text{ns} \text{spam}) = 35\%$$

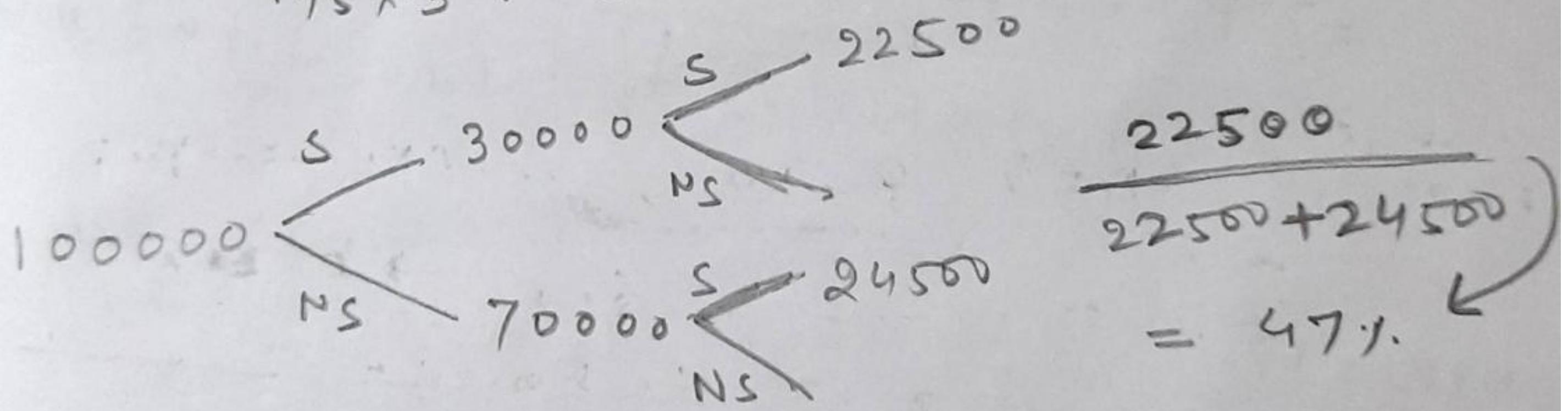
$$P(\text{spam} / \text{congratulation}) = ?$$

$$\rightarrow \frac{P(\text{cong} / \text{spam}) \times P(\text{spam})}{P(\text{cong})}$$

$$\rightarrow \frac{P(\text{cong} / \text{spam}) \times P(\text{spam})}{P(\text{cong} / \text{spam}) + P(\text{cong} / \text{ns} \text{spam})}$$

$$P(\text{cong} / \text{spam}) \propto P(\text{spam}) + P(\text{cong} / \text{ns} \text{spam}) \propto P(\text{ns} \text{spam})$$

$$= \frac{.75 \times .3}{.75 \times .3 + .35 \times .7} = 0.47$$



### Binomial Distribution

Probability of getting  $x$  success out of  $N$  trials

$$P(x) = \binom{n}{x} P^x (1-P)^{n-x}$$

$\binom{n}{x}$  is number of ways in which  $x$  success can take place out of  $N$  trials

$$\rightarrow \frac{n!}{x! (n-x)!}$$

(1) if  $p$  is prob of success of a single trial then what is the prob of getting  $x$  successes out of  $n$  trials

$n \rightarrow$  trials

$p \rightarrow$  success prob of each trial (based on past data)

$$P(x \text{ success}) = \binom{n}{x} p^x (1-p)^{n-x}$$

- E.g. → Bank issue CC to customer  
 → bank has found 60% of all accounts paying on time  
 → if 7 account is selected at random  
 → construct the Binomial prob Distribution of accounts paying on time

for ex.  $P(2 \text{ ontime}) = \binom{7}{2} (.6)^2 (.4)^5$

if 2 account among 7 pay on time  $\frac{7!}{2!5!} = 21$

through  
panda

$$n=7 \quad p=0.6 \quad k = n \cdot p. \text{arrange}(0,8)$$

$k = 0, 1, 2, 3, 4, 5, 6, 7 \rightarrow$  because there can be 8 combination for 7 random account → 0 means none of 7 paid

on time → 7 means all paid on time  
import scipy.stats as stats

$$\text{binomial} = \text{stats.binom.pmf}(k, n, p)$$

pmf → prob mass function (P)

mean & std dev. of Binomial Dist

(12)

\* mean  $\mu$  of Bin-dist.

$$\mu = E(x) = np$$

\* std. dev.  $\sigma$  is given by

$$\sigma = \sqrt{np(1-p)}$$

so for previous e.g., mean =  $7 \times 0.6 = 4.2$

$$\text{std} = \sqrt{4.2(1-0.6)} = 1.3$$

### Poisson Distribution

\* another Discrete Distribution which also plays a critical role in quality control in context of reducing the number of defect per std. unit

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$P(x)$  = Prob of  $x$  success given an idea of  $\lambda$

$\lambda$  = avg. number of success

$$e = 2.71828$$

$x$  = success per unit

$\lambda$  is parameter of Poisson Distri

mean of Poisson Distri =  $\lambda$

std. dev. of Poisson Distri =  $\sqrt{\lambda}$

$\lambda = np \rightarrow \frac{p}{n}$  (Probability)  
n (sample given)

(13) Eg. of Poisson Distribution

- \* if on avg. 6 customer arrive at every 2 min in bank during busy hours (3 customer in avg)
- a) what is prob that exactly 4 customer arrive in a given min (8 customer every 2 min)
- b) what is prob that more than 3 customer arrive in a given min

Diff w.r.t. binomial Distribution  $\rightarrow$  sample size  
is not definite (there were 7 account) but  
here 6 is avg. number. It can be more or  
less than 6 also

from pandas      rate = 6  
 $n = np \cdot \text{arrange}(0, 20) \rightarrow$  can be set to  
any reason

Poisson = stats.poisson.pmf(n, rate)

op. do has been set from computational  
perspective  $\neq$  visualization perspective

a)  $P(x=8)$

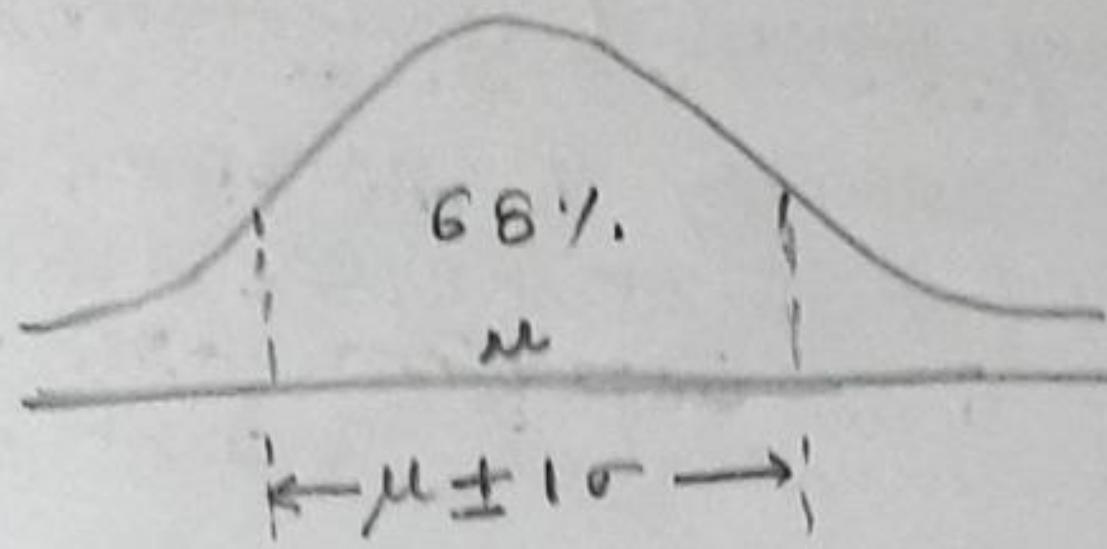
b)  $1 - P(x \geq 3) = 1 - P(x < 3)$

$$P(x \geq 3) = 1 - (P(0) + P(1) + P(2))$$

## Gaussian/normal Distribution

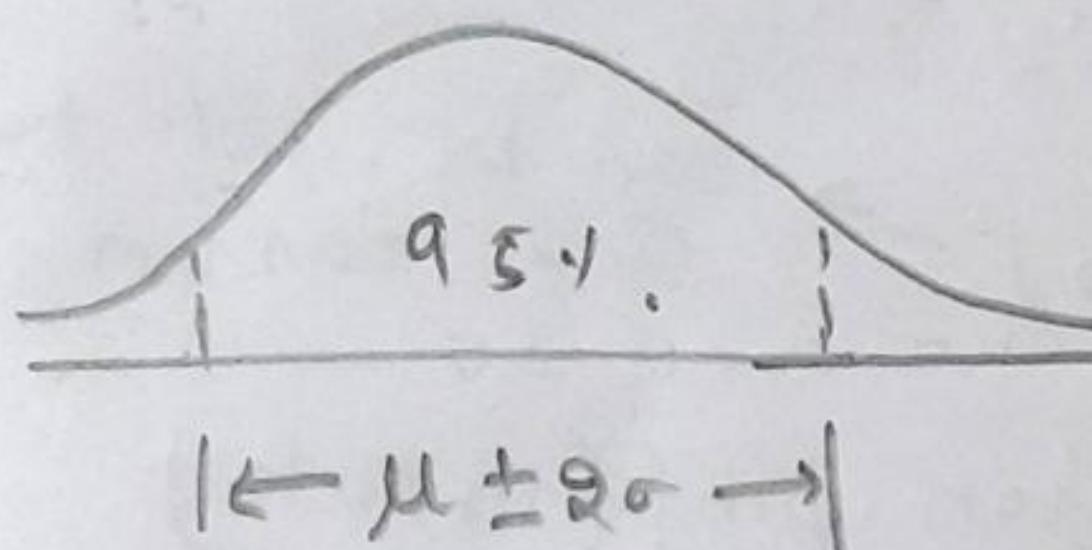
(14)

kg.

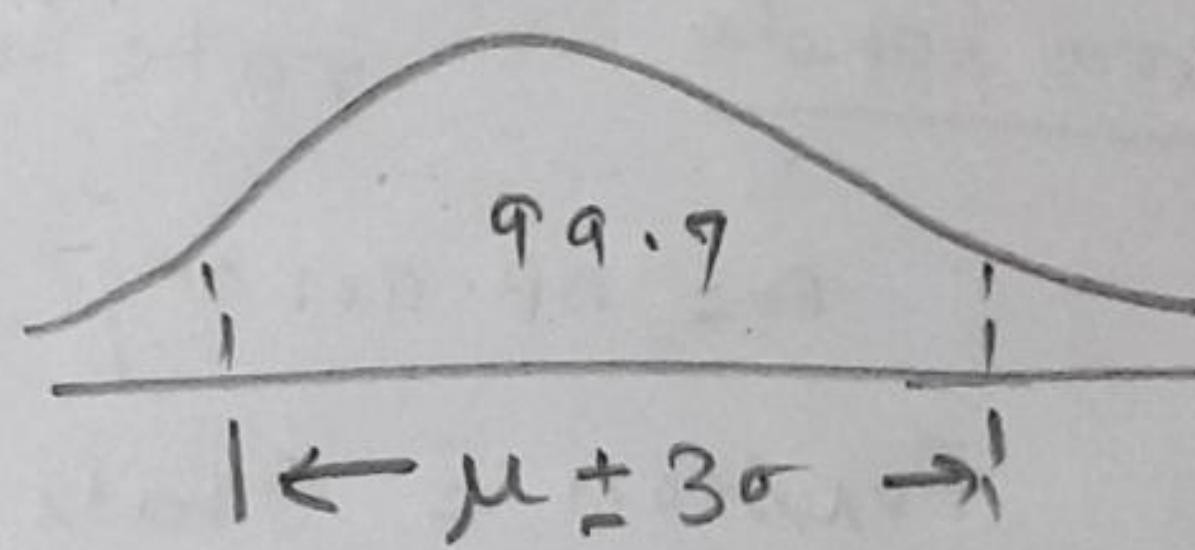


- derived from "central limit theorem"
- used when we don't have enough data but know mean & std dev.

- approximate 68% of the data in a bell shaped dist in within 1 std dev of the mean  
i.e.  $\mu \pm \sigma$  ( $\mu$  is mean  $\sigma$  is std.dev)
- This is empirical Rule approximates the variation of data in a bell shaped distribution



chances of being  
2 std dev



chance of being  
3 std dev

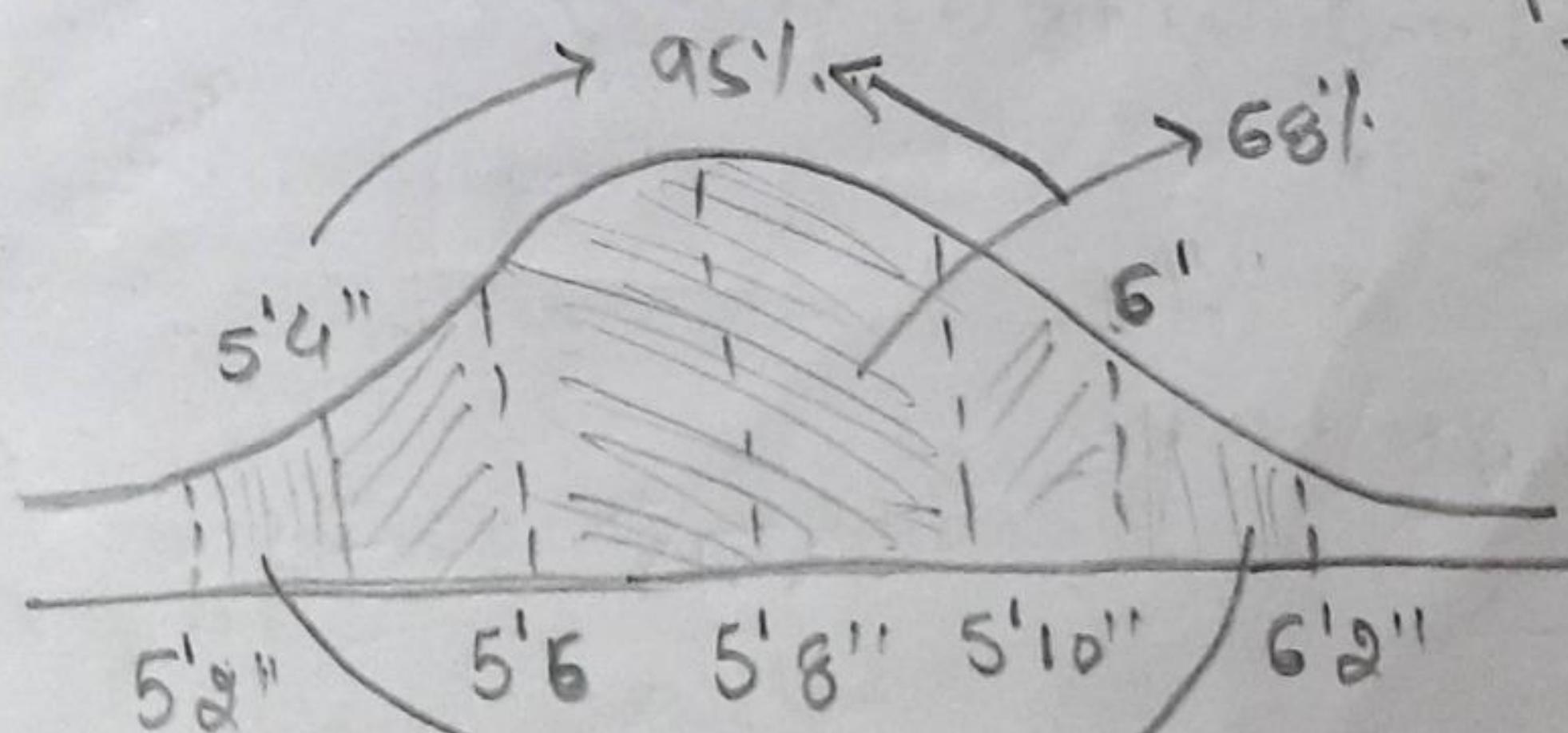
for ex:

- mean height = 5 ft 8 inch
- std dev = 2 inch

library

statn. norm. cd

cdf = cumu  
dist fun



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

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\left(\frac{z^2}{2}\right)}$$

$Z$  is number  
independent v

(15) E.g. mean of morning breakfast pack is 295 kg  
a std dev is 0.025 kg.

a) what is prob. that pack wt. is less than .28 kg.

b) " " " " more than .350 kg.

c) " " " " between .26 to .34 kg.

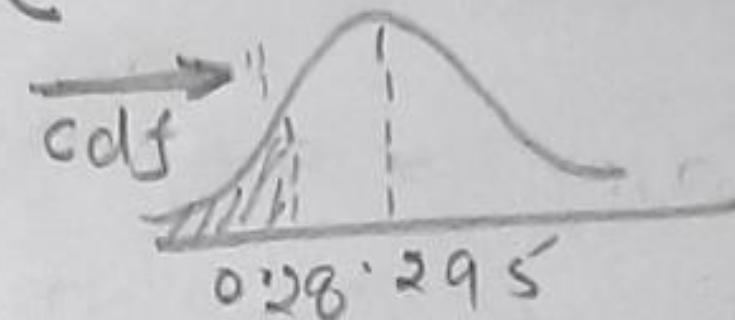
$$a) z = (.28 - .295 / 0.025) = -0.6$$

$$\text{statistic. norm. cdf}(z) = 27\%$$

OR

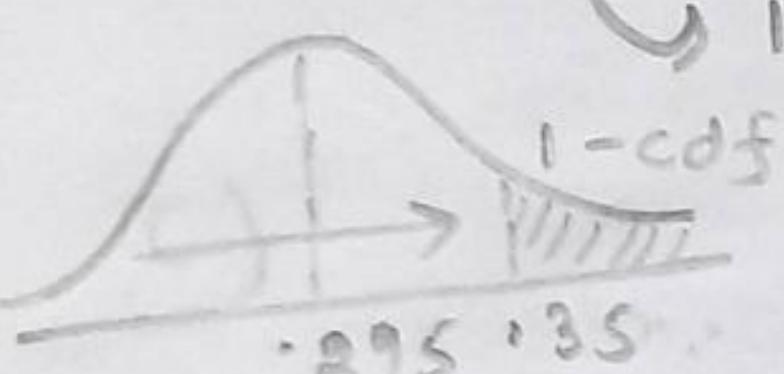
statistic. norm. cdf (.28, loc = .295, scale = 0.025)

$$\hookrightarrow = 27\%$$



$$b) 1 - \text{statistic. norm. cdf} (.350, .295, .025)$$

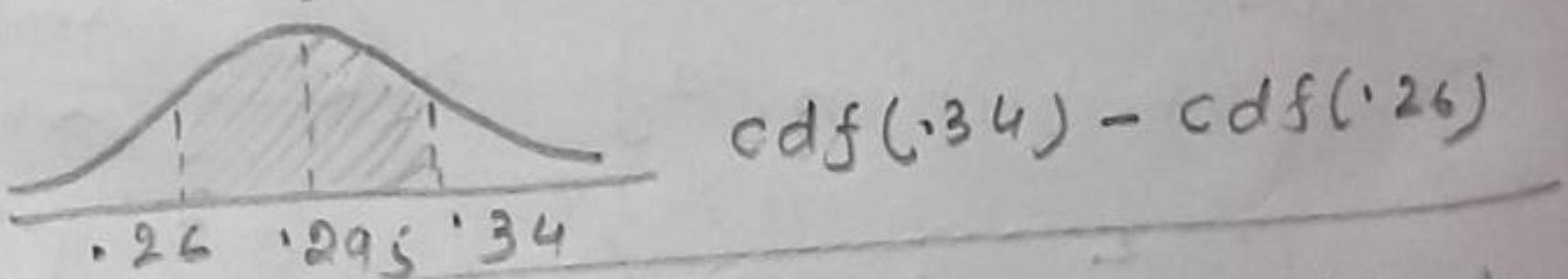
$\hookrightarrow 1.39\%$ . (whole prob is 1, so for more than we have to do  $1 - \text{cdf}$ )



$$c) \text{statistic. norm. cdf} (.340, .295, .025)$$

$$- \text{statistic. norm. cdf} (.260, .295, .025)$$

$$\hookrightarrow = 88\%$$



sampling Distribution  $\rightarrow \bar{x}$   
 $\rightarrow$  Total number

$N = 60 \rightarrow$  mean

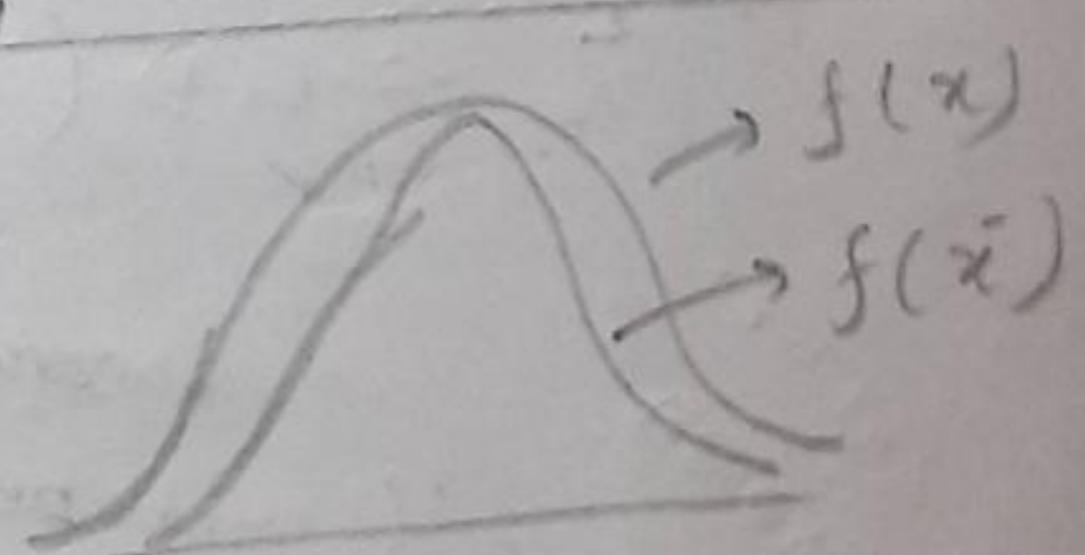
$\mu = 70 \rightarrow$  std dev.

$\sigma = 5 \rightarrow$  sample taken

$n = 4 \rightarrow$  sampling Distribution:  $\bar{x} \rightarrow$  take 4 Random

sample ( $60C_4$ )  $\rightarrow$  take mean of each sample

set  $\rightarrow$  plot distribution of each sample mean  
 $\rightarrow$  called sampling Distribution



$$\sigma(\bar{x}) = \frac{\sigma(x)}{\sqrt{\text{sample size}}}$$

sampling error  $|\bar{x} - \mu|$

(16)

- if  $x_1, x_2, x_3, \dots, x_n$  are  $n$  independent random samples drawn from Normal population with mean  $\mu$ , std dev  $\sigma$ , then sampling distribution  $\bar{x}$  follows a normal distribution with mean  $\mu$  & std dev  $\frac{\sigma}{\sqrt{n}}$ .  $\frac{\sigma}{\sqrt{n}}$  is known by term std. error.

central limit theorem

$$x \not\sim AD/ND \left( \mu, \sigma^2 \right)$$

Random variable  $\rightarrow$  may & may not in normal Distribution

min sample size  
theoretically  
to meet central  
limit theorem  
 $n \geq 30$

Sample,  $x_1, x_2, \dots, x_{30}$

$S_2$

:

$S_{100}$

$\bar{x}_1$  (mean of  $S_1$ )

$\bar{x}_2$

$\bar{x}_{100}$

if I take all above & try to plot  $\bar{x}$ , it will follow normal distribution irrespective what is distribution of  $x$  (as sample size becomes larger & variance becomes smaller)

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

when std dev is known

mean  
(will be approx same as mean of  $x$ )

$$\left( \frac{\text{std dev}}{\text{or}} \frac{1}{\text{std err}} \right) \frac{\sigma}{\sqrt{n}}$$

$\bar{x} - \mu \rightarrow$  sampling error

when proportion is known

$$P \pm Z_{\alpha/2} \cdot \sqrt{\frac{P(1-P)}{n}}$$

$$\bar{x} \pm Z_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right)$$

margin of error

$Z \rightarrow$  confidence coefficient

2.58  $\rightarrow$  99% confidence

1.96  $\rightarrow$  95% confidence

Refer calculate Z-Score video on youtube

### ⑦ Types of probability

- 1) Priori classical probability  $\rightarrow$  apply Rule to find prob.
  - 2) empirical prob  $\rightarrow$  data driven prob or Relative freq. Prob
  - 3) subjective prob
    - $\hookrightarrow$  start without data & ask experts
    - $\hookrightarrow$  learn & revise based on data availability over a period
- $P(A) = \frac{m}{n}$   $\rightarrow$  favorable outcome  
 $n$   $\rightarrow$  total number of outcome

for ex. in KBC

50-50 option  $\rightarrow$  classical prob

Audience poll  $\rightarrow$  Data Driven prob (empirical)

Phone a friend  $\rightarrow$  subjective prob

### Baye's theorem Exercise

drilling company

- $\rightarrow$  estimated 40% chance of striking oil for new well
- $\rightarrow$  detailed test has been scheduled for more info
- $\rightarrow$  Historically, 60% successful wells had detailed test
- $\rightarrow$  20% unsuccessful wells had detailed test
- $\rightarrow$  given that detailed test is scheduled, what is prob. of well will be successful.

\*  $S$  is success,  $\bar{S}$  is unsuccessful

$$P(D/S) = 60$$

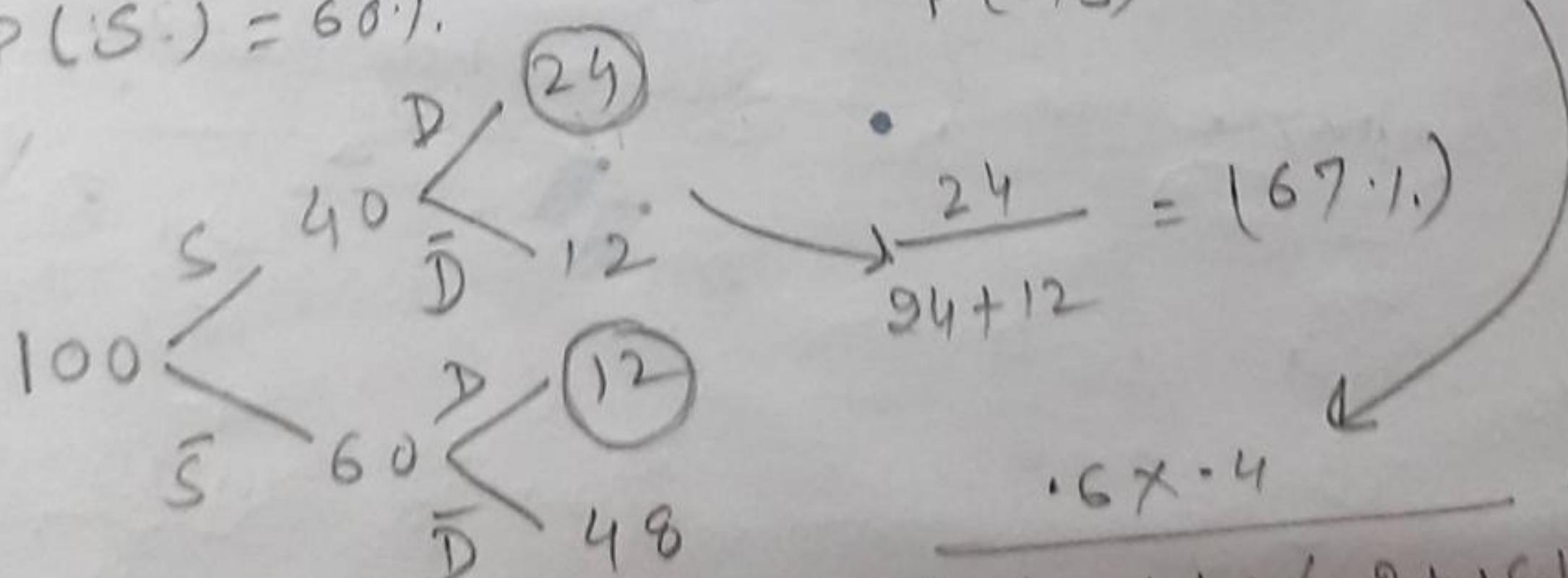
$$P(\bar{D}/\bar{S}) = 20\%$$

$$P(S) = 40\%$$

$$P(\bar{S}) = 60\%$$

$$P(S/D) = \frac{P(D/S) \times P(S)}{P(D)}$$

$$= \frac{P(D/S) \times P(S)}{P(D/S) \times P(S) + P(\bar{D}/\bar{S}) \times P(\bar{S})}$$



$$\frac{.6 \times .4}{(.6 \times .4) + (.2 + .6)} = \frac{.24}{.24 + .12} = 0.67\%$$

Prob. sample que

- Q1 Out of 37 men and 33 women
- 36 are completely refrain from alcoholic beverages
  - 9 women are non-smoker
  - 16 men smoke but do not drink
  - 13 men & 7 women drink but do not smoke

How many both drink & smoke : what is associated prob?

solution → (using marginality principle)

$$\begin{array}{ll} m \rightarrow \text{men} & m + \bar{m} = 70 \\ \bar{m} \rightarrow \text{women} & D + \bar{D} = 70 \rightarrow D = 34 \\ & \downarrow \\ & 36 \end{array}$$

$$D = DS + D\bar{S} \rightarrow DS = 34 - D\bar{S}$$

$$D\bar{S} = mD\bar{S} + \bar{m}D\bar{S} = 13 + 7 = 20$$

$$DS = 34 - 20 = 14$$

$$P(DS) = \frac{14}{70} = 20\%$$

- Q2 Employee of company were surveyed on question regarding college degree

→ of 600 → 400 has college degree

→ 100 are single

→ 60 were single college student

		Total		$\bar{m} \text{ or } D = \bar{m}\bar{D} + \bar{m}$
Sol. - I	340	D	$\bar{D}$ 160	
m	MD	$m\bar{D}$	$\bar{m}\bar{D}$	$m\bar{D} + \bar{m}\bar{D}$
single	400	200	40	600

$m\bar{D}$  is already counted

(19)

$$\bar{m} \text{ or } D = mD + \bar{m}D + \bar{\bar{m}}D$$

$$= 340 + 60 + 40 = 440$$

$$P(m \text{ or } D) = \frac{440}{600} = 77.3\%$$

approach - 2 (not mutually exclusive approach)

$$P(m \text{ or } D) = P(m) + P(D) - P(m \cap D)$$

$$= \frac{400}{600} + \frac{100}{600} - \frac{60}{600}$$

$$= \frac{440}{600} = 77.3\%$$

Q3. \* prob of house rate will increase next month is

- \* prob of interest rate going up in same period is .74
- \* prob of house rate or interest rate go up
- \* prob of house rate increase in next 6 months is 0.89

find prob. of both house rate & Interest Rate will Increase in next 6 month is ?

$$SI + \bar{S}I + S\bar{I} = .89$$

approach - 1

S	I	$\bar{I}$		
	SI	$\bar{S}\bar{I}$	0.25	
$\bar{S}$	$\bar{S}I$	$\bar{I}$	$\circled{.75}$	
.74	$\circled{.26}$			1

$$SI + (.25 - SI) + (.74 - S\bar{I}) = .89$$

$$SI = .99 - .89 = .1$$

approach 2

$$P(S \text{ or } I) = P(S) + P(I) - P(S \cap I)$$

$$P(S \cap I) = P(S) + P(I) - P(S \cup I)$$

$$= .25 + .74 - .89 = .1$$

Q4. company has 2 machine

- \* old machine produce 23% defective widget
- \* new machine produce 8% defective widget
- \* in addition new machine produce 3 times as widget as older machine

give the randomly chosen widget, what is prob it's produced by new machine?

Approach - 1

N	$\bar{N}$
75	25
.08	.23

(20)

$$D = DN + D\bar{N}$$

$$= .08 \times 75 + .23 \times 25$$

$$P(DN) = \frac{DN}{DN + D\bar{N}} = \frac{.08 \times 75}{.08 \times 75 + .23 \times 25} = 0.51$$

Approach - 2

D	$\bar{D}$	75
N	$.08 \times 75 = 6$	69
$\bar{N}$	$.23 \times 25 = 5.75$	19.25
		<u>25</u>
		<u>100</u>
	<u>11.75</u>	<u>88.25</u>

$$P(DN) = \frac{6}{11.75} = 0.51$$

Q5 → in a used car lot

- 1) 70% have A/c
- 2) 40% CD player
- 3) 20% have Both

	A/c	$\bar{A}/\bar{C}$	Total
CD	• 2	• 2	• 4
$\bar{C}D$	• 5	• 1	• 6
	• 7	• 3	1

1) how many car. have A/c but no CD → .5 i.e. 50

2) what is prob of car has CD given it has

$$A/c \text{ also } \rightarrow \frac{2}{7} \Rightarrow \frac{20}{70} = .2857$$

$$P(CD/Ac) = \frac{P(Ac/CD) \times P(CD)}{P(Ac)}$$

$$\text{by Baye's theorem} = \frac{\frac{2}{7} \times .4}{P(Ac/CD) \times P(CD) + P(Ac/\bar{CD}) \times P(\bar{C}D)}$$

$$= \frac{\frac{2}{7} \times .4}{.2 + \frac{5}{6} \times .5} = \frac{.2}{.7} = .2857$$

$P(CD) \& P(Ac) \rightarrow \text{marginal}$  /  $P(CD/Ac) \rightarrow \text{conditional prob}$   
 $P(CD) \& P(Ac) \rightarrow \text{joint}$  /  $P(Ac) \rightarrow \text{joint prob}$

② sample size problem

Q1 → marketing manager of food chain wishes to estimate avg. yearly amount that family spend on fast food. He wants an estimate to be within 100 with a confidence level of 99%. It is known from earlier pilot study that std dev. of family expenditure on fast food is 500. How many family must be chosen for this problem.

for 99% confidence → value of  $Z = 2.58$

& std. dev. is known

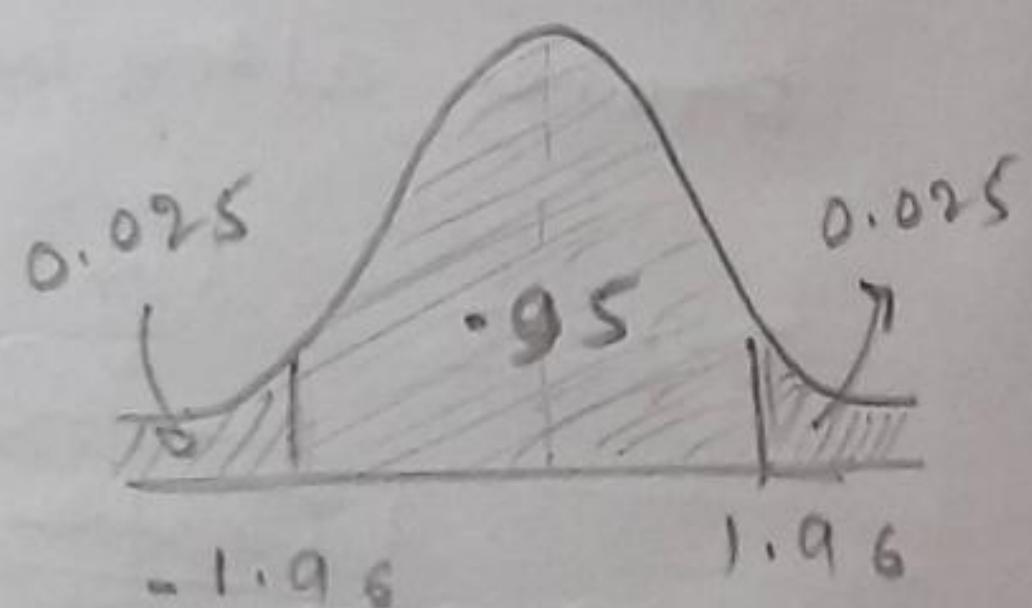
$$100 = 2.58 \times \frac{\sigma}{\sqrt{n}} \Rightarrow 100 = 2.58 \times \frac{500}{\sqrt{n}}$$

called margin of error calculate for  $n \rightarrow$  will be sample size number of above problem

Q2 → sports good manufacturer want to estimate the proportion of cricket player among high school student in India. The company wants the estimate to be within .03 with a confidence level of 95%. A pilot study done earlier reveals that out of 80 high school student, 36 play cricket. What should be sample size for study.

for 95% confidence →  $Z = 1.96$

$$1.96 \times \sqrt{\frac{P(1-P)}{n}} = 0.03$$



$$P = \frac{36}{80} = 0.45$$

$$1.96 \times \sqrt{\frac{0.45 \times 0.55}{n}} = 0.03$$

calculate  $n \rightarrow$  will be sample size for above solution

$$\downarrow \\ (1050)$$

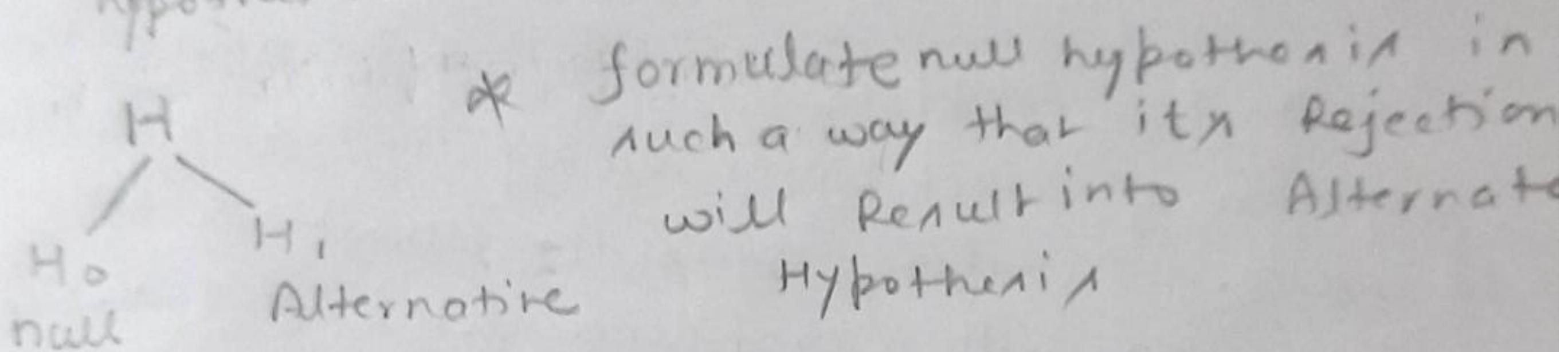
Q3 → wait if in Q2, P is not given as part of problem statement. How to calculate sample size in such case is when P is unknown take max value of P & 1-P i.e., .5 each

$$1.96 \times \sqrt{\frac{0.5 \times 0.5}{n}} = 1.02$$

calculate  $n \rightarrow$  will be my sample size  
 ↓  
 (1067)

### Hypothesis (statistical)

a statistical hypothesis is a statement about a population parameter. It may or may not be true. The manager has to ascertain the truth of hypothesis.



- \* A null hypothesis is status quo.
- \* Researcher / Decision maker generally want to prove alternate hypothesis.

### Type I and Type II error (Null Hypothesis)

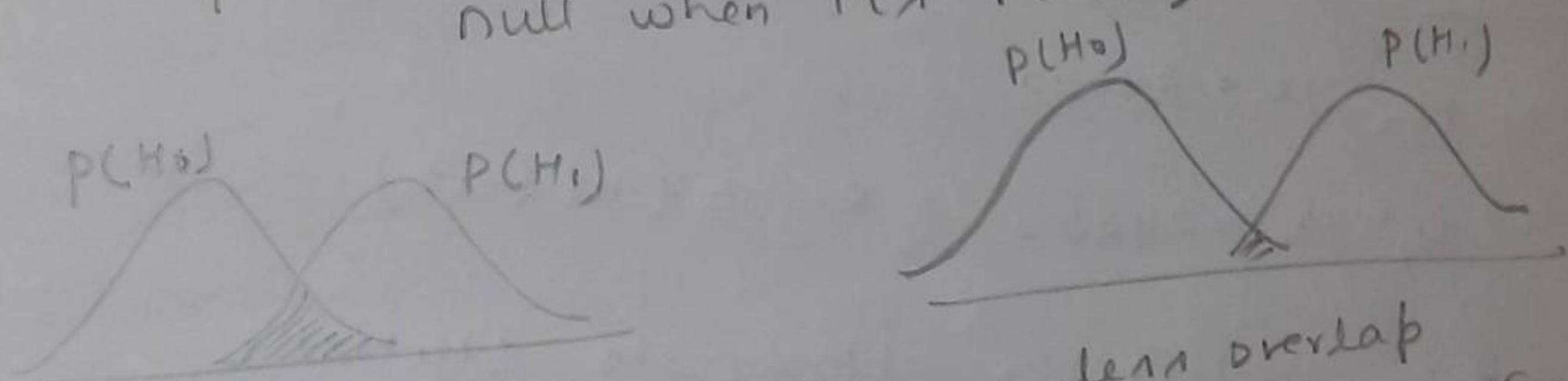
		True	False
Reject	True	Type I error ( $\alpha$ )	No Error
	False	No Error	Type II error ( $\beta$ )
Accept			

- \* Rejecting Null hypothesis when it's True type I error.
- \* Accepting the null hypothesis when it's False type 2 error.

(23)  $\alpha \rightarrow$  prob of committing the type I error  
also called the level of significance of  
the test (statistical level of significance)

$\beta \rightarrow$  prob of committing the type II error

$\checkmark 1 - \alpha \Rightarrow$  confidence level of the test (prob of accepting null when it's true)  
 $1 - \beta \Rightarrow$  power of the test (prob of rejecting null when it's false)



$\alpha \uparrow \beta \downarrow$        $\alpha \downarrow \beta \uparrow$       its tradeoff

( $1 - \beta$  can also be written as)

$\checkmark 1 - \beta \rightarrow$  prob of accepting alternative when it's true  
e.g. launching a product line in a new market

Survey of random sample of 400 households  
in that market showed a mean income per  
household of \$30000. Std. dev. based on a  
pilot study of 400 households is \$8000.

- \* Karen believes the product line will be  
adequately profitable only in market where  
mean household income is greater than \$39000
- \* Should Karen introduce product line into  
new market?

$$n = 400$$

$$\bar{x} = \$30000$$

$$\sigma = \$8000$$

(24)

- ✗ alternate hypothesis can be  $>$ ,  $<$ ,  $\neq$  but can never have equality sign
- ✗ Null hypothesis is always tested for threshold as should have equality sign

Step-1 formulate Null & alternate hypothesis

$$H_0: \mu = 29000$$

$$H_1: \mu > 29000 \text{ (right tail test)}$$

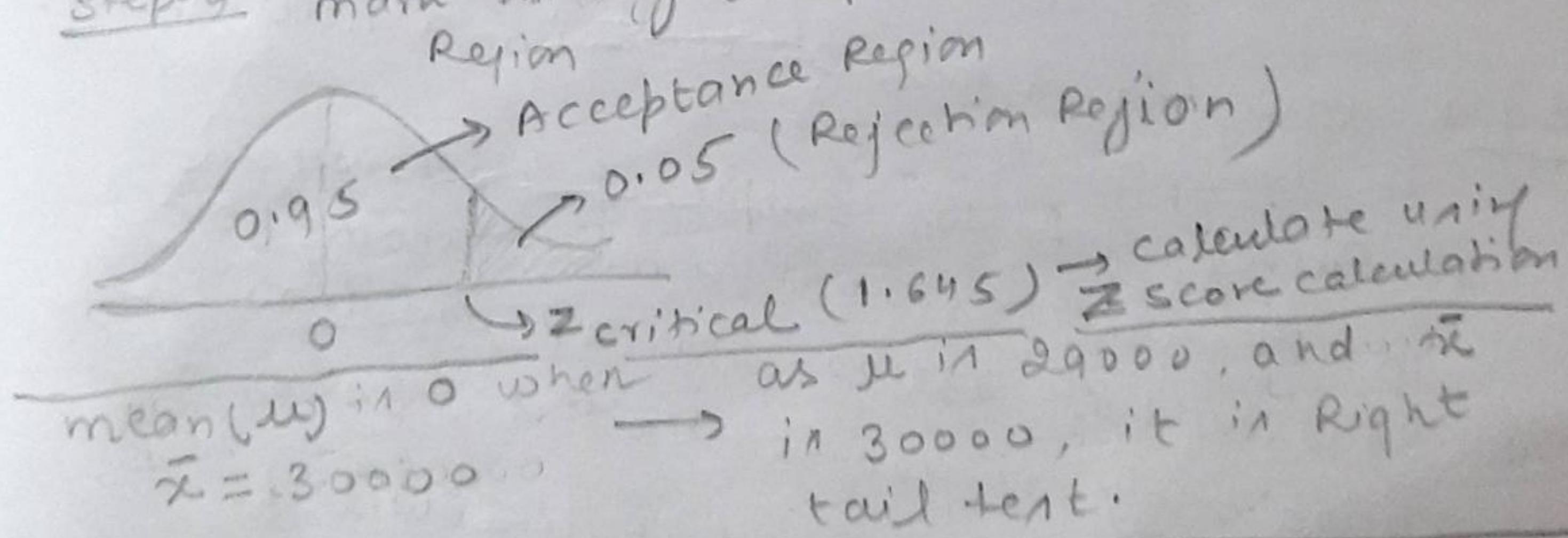
Step-2 Select the appropriate statistic

$$Z_{\text{stat}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \rightarrow \text{sampling error}$$

Step-3 compute  $Z_{\text{stat}}$

$$Z_{\text{stat}} = \frac{30000 - 29000}{8000/\sqrt{400}} = 2.5$$

Step-4 mark in diag. acceptance Region & Rejection



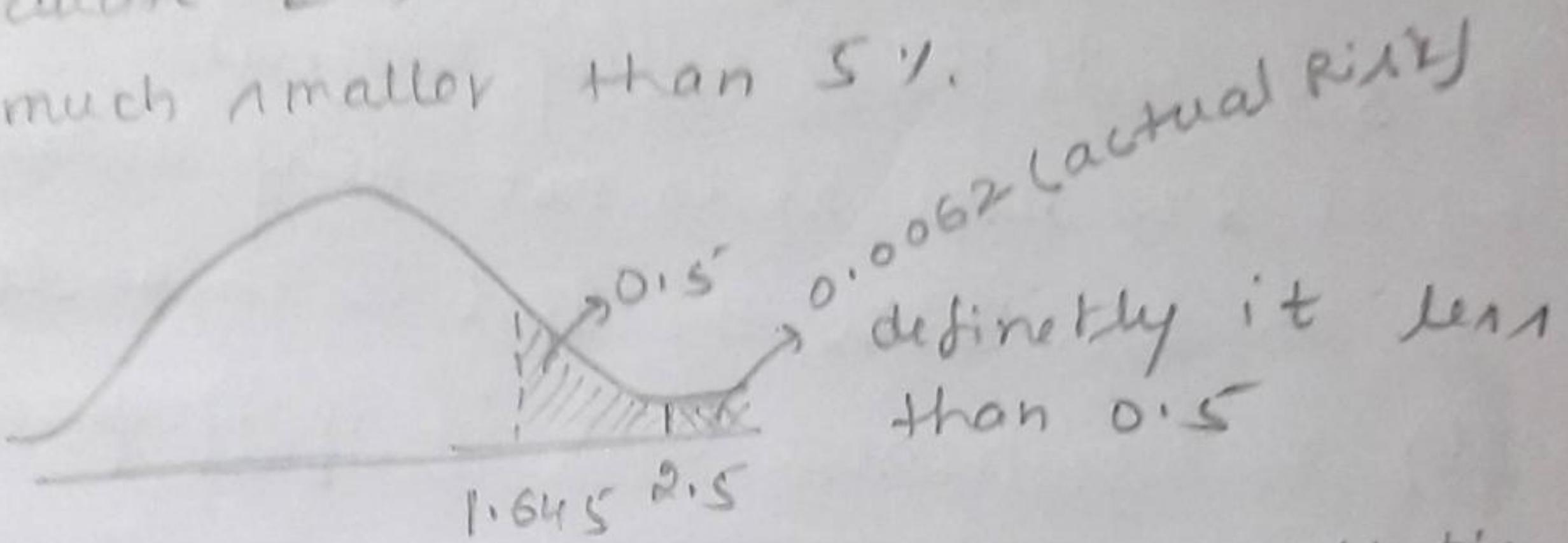
Theoretically when nothing is mentioned  $\rightarrow$  take  $\alpha$  as 5%. that means 95% confidence and level of significance is 5%. that means

$$P(Z) \geq 1.645 = 0.05 \text{ i.e. } 5\%.$$

(25) Step 5 → find out where computed Z falls,  
Z falls in Rejection Region as it  
is greater than 1.645

∴ hence Null hypothesis is Rejected  
if  $P(Z) \geq 1.645 = 0.5$  means, we are  
allowed to take risk of 5%.

\* calculate Z is 2.5 i.e. actual Risk  
is much smaller than 5%.



$P(0.5)$  → called P-value of statistic  
(Actual Risk)

Step 6 → As Null hypothesis is Rejected,  
Alternate hypothesis is Accepted, hence  
Keren can launch the product in new  
market.

as desired risk is 5%, Actual Risk is  
0.62%, means I am 99.38% confident  
on the decision made.

if  $P_{val} \leq \alpha$  (in above case  $0.62\% \leq 5\%$ )  
Null hypothesis is Rejected  
↳ (another way to approach above problem)

E9.2 A retailer is weighing strawberries to sell as 250 g. punnets. customer complained that what he bought was less than 250 g. retailer decided to check weight of 36 punnet. He finds avg wt. is 248.5 g. with std dev. of 4.8 g. Using significance test to judge whether he is selling under wt. punnet, which of following is true.

- a) at 5% level he is selling underwt.
- b) at 5% level he is not selling underwt
- c) at 5% level the test is inconclusive
- d) A significance test is inappropriate in this case.

Step 1

$$\bar{x} = 248.5 \quad \mu = 250$$

$$n = 36$$

$$H_0: \mu = 250$$

$$\sigma = 4.8$$

$$H_1: \mu < 250 \quad (\text{left tail test})$$

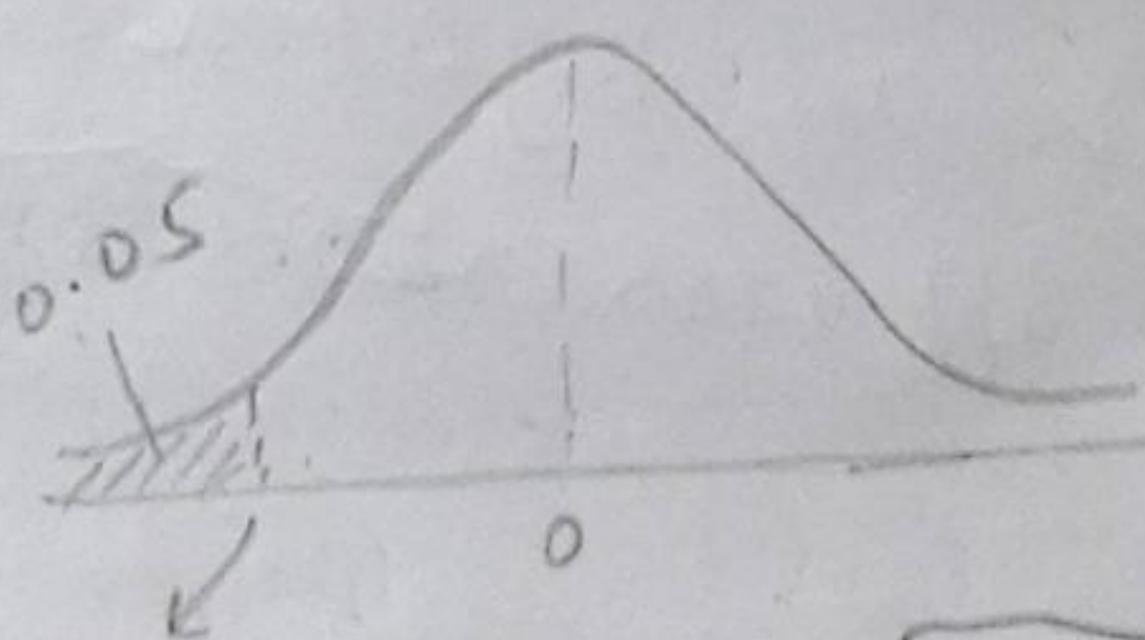
Step 2

$$z_{\text{stat}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

Step 3

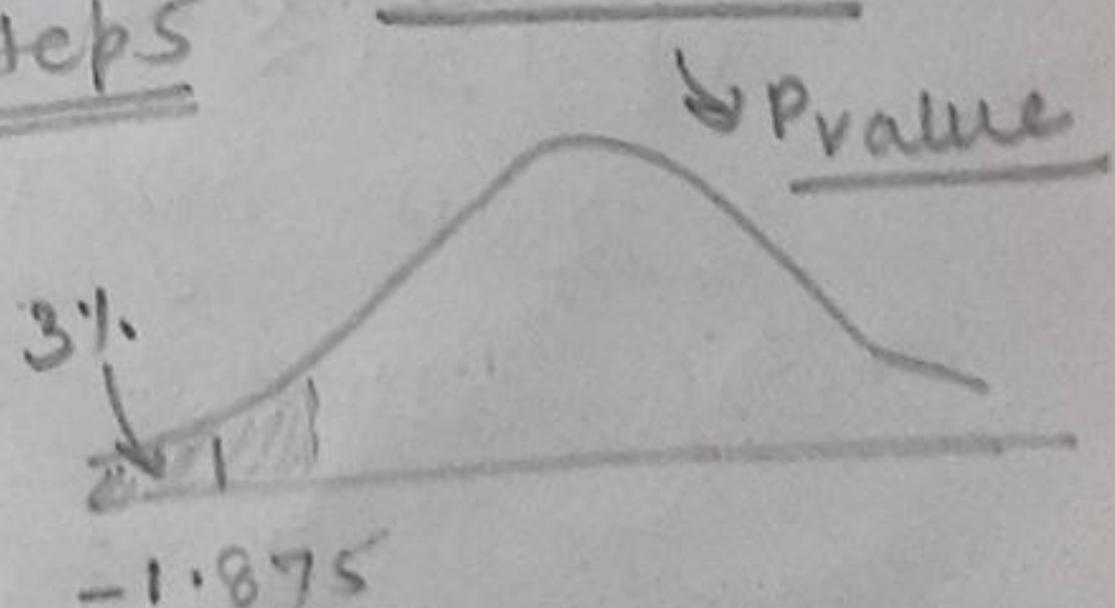
$$z_{\text{stat}} = \frac{248.5 - 250}{4.8/\sqrt{36}} = -1.875 \quad P(z_{\text{stat}} < -1.875) = .0304$$

Step 4



$$z_{\text{critic}} = -1.645$$

Step 5



p-value < cc

↓ falling in Rejection region

Hence

at 5% confidence level, he is selling underwt. (Null hypothesis rejected)

(27) State null & Alternative hypothesis for the following

a) Is the avg. waiting time for the customer of super market at checkout greater than 15 min  
 $H_0 \rightarrow \mu = 15$        $H_1 \rightarrow \mu > 15$        $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$

b) Is proportion of household owning color TV in Chennai less than 0.4 (proportion is always in %, and n is always a number)  
 $H_0 \rightarrow \pi = 0.4$        $H_1 \rightarrow \pi > 0.4$        $Z = (P - \pi) / \sqrt{\pi(1-\pi)/n}$

c) Is the avg. expenditure per household on eating out significantly higher in Bangalore than in Calcutta.

↳ it involves comparing two population  
 $H_0 \rightarrow \mu_1 = \mu_2$        $\mu_1 \rightarrow$  Bangalore eating  
 $H_1 \rightarrow \mu_1 > \mu_2$        $\mu_2 \rightarrow$  Calcutta eating  
 $H_0 \rightarrow \mu_1 - \mu_2 = 0$       ↳ T-test compare  
 $H_1 \rightarrow \mu_1 - \mu_2 > 0$

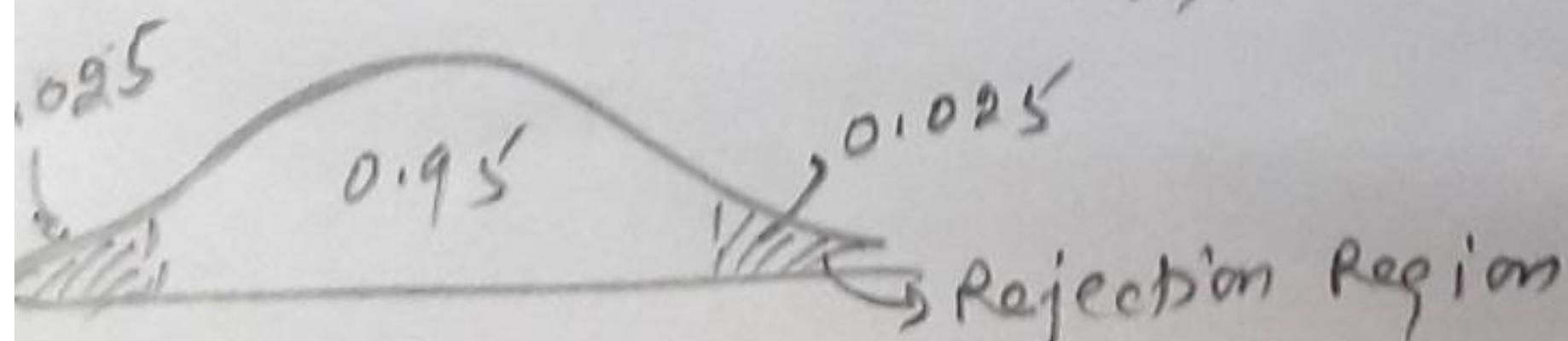
d) Two random sample survey conducted with 2 month gap to assess public opinion on outcome. The question that was posted was

"If the general election going to take place tomorrow, would you cast vote for or against the ruling party?"

↳ (this is an opinion poll problem)

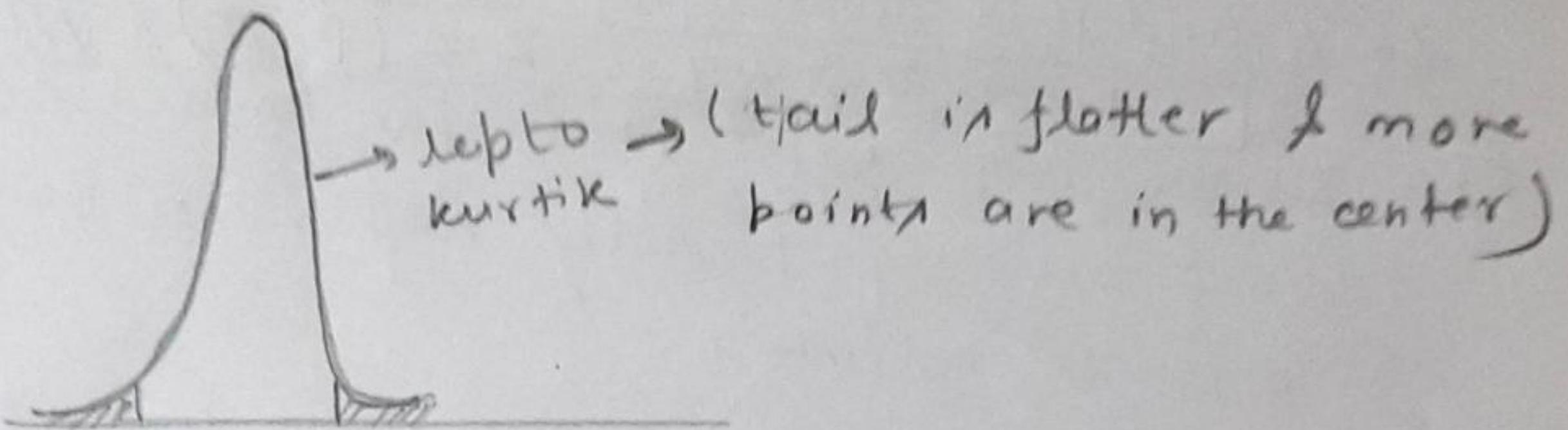
$H_0 \rightarrow \pi_1 = \pi_2$   
 $H_1 \rightarrow \pi_1 \neq \pi_2$  ↳ Two tail test

$\pi_1 \rightarrow$  Proportion of people in first poll vote for Ruling party



$\pi_2 \rightarrow$  Proportion of people in 2nd poll

- (26) use case (for eg. in case of scientific tests)  
 if  $n < 30$  &  $\sigma$  is unknown ( $\rightarrow$  Population std dev)
- \* To do T-test, assumption in distribution is normal
  - \* Use T-test  $\rightarrow t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$   
 $s \rightarrow$  sample std. dev.



### T-test application - one sample

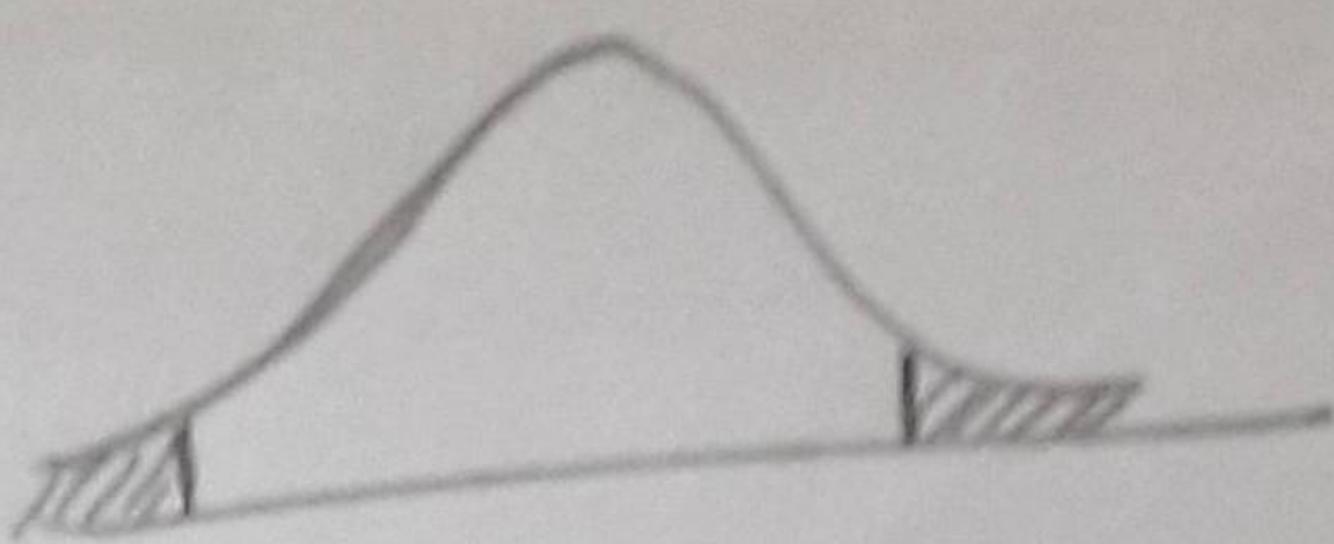
Experian Marketing Services reported that the typical American spend a mean of 144 min per day accessing Internet via mobile. In order to test validity of the name, you select a sample of 30 friends & family. The results for time spent per day accessing Internet via mobile device are stored in InternetMobileTime.  $\rightarrow$  Excel file name

- a. Is there evidence that the population mean time per day is diff from 144 min.  
 Use p-value approach and level of significance is 0.05.
- b. what assumption about population distribution is needed in order to conduct the t-test in (a)?

(29)

$$H_0 \rightarrow \mu = 144$$

$$H_1 \rightarrow \mu \neq 144 \quad \text{Two tail test}$$



$\bar{x} = \boxed{175.2667}$  mean (of all entries of given excel sample)

$s = \text{std dev (of all entries of given excel sample)}$

$$\boxed{139.8368} \rightarrow 1.224674$$

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

$P\text{value} = 2 \times P_{\text{value}}(t) \rightarrow (\text{due to two tail test})$   
left tail + Right tail.

↓  
calculated as 0.2305

$P_{\text{value}} > \alpha$  → Null hypothesis is accepted  
(83.1%) (5%)  
↓  
means it's not diff from  
144 minute