

①

STATISTICS



DESCRIPTIVE

INFERENTIAL

1# DESCRIPTIVE

Central Tendency }
or
Average.

Mean (Arithmetic Mean)

median
mode

Normal
(Average)
Mean



Most common
Number

1 1 2 3 4 → 1
1 1 2 3 4 4 → ?

1 1 [2] 3 4

(after sorting)

1 1 [2 3] 4 4

$$\frac{1+1+2+3+4}{5}$$

Arithmetic mean
of these number ⇒ 2.5

ex:-

3 3 [3 3] 3 100

$$\rightarrow \text{mean: } \frac{115}{6} = 19 \frac{1}{6} \leftarrow \text{central Tendency}$$

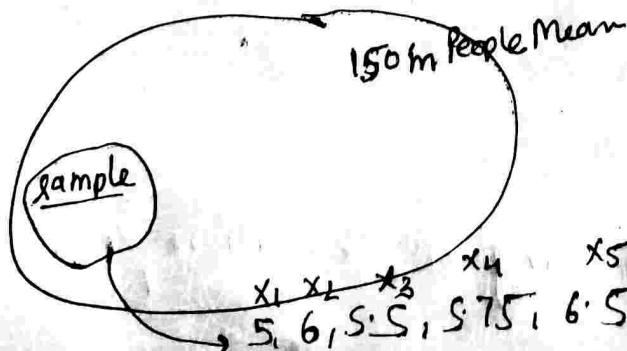
$$\rightarrow \text{median: } \frac{3+3}{2} = 3$$

→ Mode : 3.

2# Sample Population :-

μ = Population mean

\bar{x} = sample mean



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

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

Measures of Dispersion

ex: 2 2 3 3,

$$\mu = \frac{2+2+3+3}{4} = 2.5$$

→ close to 2.5

$$\mu = \frac{0+0+5+5}{4} = \underline{\underline{2.5}}$$

→ away from 2.5

more dispersion

How to measure this?

⇒ Variance: (σ^2)

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

ex.: 2 2 3 3 →

$$\sum_{i=1}^4 (x_i - \mu)^2 = 1$$

i	x_i	μ	$x_i - \mu$	$(x_i - \mu)^2$
1	2	2.5	-0.5	0.25
2	2	2.5	-0.5	0.25
3	3	2.5	0.5	0.25
4	3	2.5	0.5	0.25

$$\sigma^2 = \frac{1}{4} = 0.25 \quad (\text{Average squared distance from mean})$$

distance from mean

ex: 0, 0, 5, 5

$$\sigma^2 = \frac{0+0+25+25}{4} = 12.5$$

Note: Everything we are doing in Population.

$$\Rightarrow \sigma^2 = \text{variance} = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}; \quad \mu = \frac{\sum_{i=1}^N x_i}{N}$$

\Rightarrow Very hard to calculate for population. /

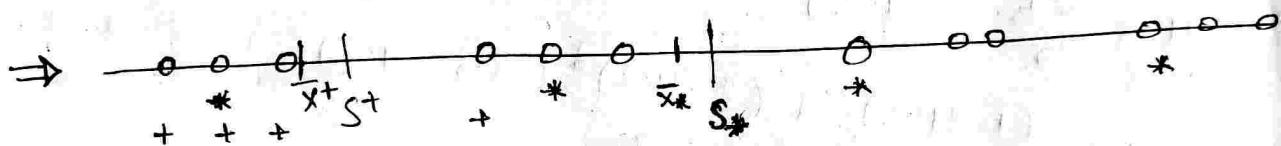
\Rightarrow Calculate parameters in Sample & make Inferences.

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

Sample mean ↑

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

Sample Variance ↑



$\Rightarrow *$ → sample 1-items;

\bar{x}_* → sample 1 mean ; S_*^2 → sample 1 variance

$\Rightarrow +$ → sample 2-items & (Skewed Sample)

\bar{x}_t → sample 2 mean ; S_t^2 → sample 2 variance

$\Rightarrow S_t$ → underestimating the population variance.

Unbiased Sample Variance \hat{s}^2

$$\Rightarrow \boxed{S^2 = S_{n-1}^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

* Standard deviation $\sigma = \sqrt{\sigma^2}$; $S = \sqrt{S^2}$

* Why standard deviation is required:-

Suppose x is in meter; $x \rightarrow$ i/p data.

Variance will be 2.5 (suppose) meter²

So, here variance is in square meter.

Now, to get the value in appropriate unit we use standard deviation.

examples 1 2 3 8 7

$$\mu = \frac{21}{5} = 4.2$$

$$\sigma^2 = \frac{(1-4.2)^2 + (2-4.2)^2 + (3-4.2)^2 + (8-4.2)^2 + (7-4.2)^2}{5}$$

$$\sigma^2 = \frac{38.8}{5} = \underline{\underline{7.76}} \quad \sigma = \underline{\underline{2.79}}$$

$$S^2 = \underline{\underline{9.7}} \quad ; \quad S = \underline{\underline{3.11}}$$

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

$$\rightarrow \frac{\sum_{i=1}^N x_i^2 - 2\mu \sum_{i=1}^N x_i + \mu^2 \sum_{i=1}^N}{N}$$

$$\rightarrow \frac{\sum_{i=1}^N x_i^2 - 2\mu \sum_{i=1}^N x_i + \mu^2 N}{N}$$

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

$$\Rightarrow \frac{\sum_{i=1}^N x_i^2}{N} - 2\mu^2 + \mu^2 = \frac{\sum_{i=1}^N x_i^2}{N} - \mu^2$$

faster way of calculating variance

$$\sigma^2 = \frac{\sum_{i=1}^N x_i^2}{N} - \mu^2 \quad \xrightarrow{\text{raw score method}}$$

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

We don't have to calculate mean ahead of time.

Normal Distributed : (Gaussian Distribution)

$$p(x) = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

standard Z-score

$$p(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{(x-\mu)^2}{2\sigma^2} \right)$$

how much standard deviation we are away from mean,

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} = \frac{1}{\sqrt{2\pi\sigma^2} e^{-\frac{1}{2}z^2}}$$

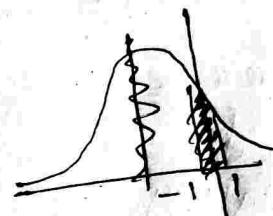
\downarrow \downarrow \uparrow

$$\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \left(e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}\right)^{-\frac{1}{2}}$$

④ To know the probability from probability density graph, we need to take small area (around that point) under the curve.

so,

$$\int_{-1}^1 \frac{1}{10\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x+5}{100})^2} dx$$



But; not an easy integral to evaluate. so can be done numerically.

A function is defined called cumulative distribution function, that is useful tool to figure out this area

~~CDF~~
$$CDF(x) = \int_{-\infty}^x p(x) dx$$

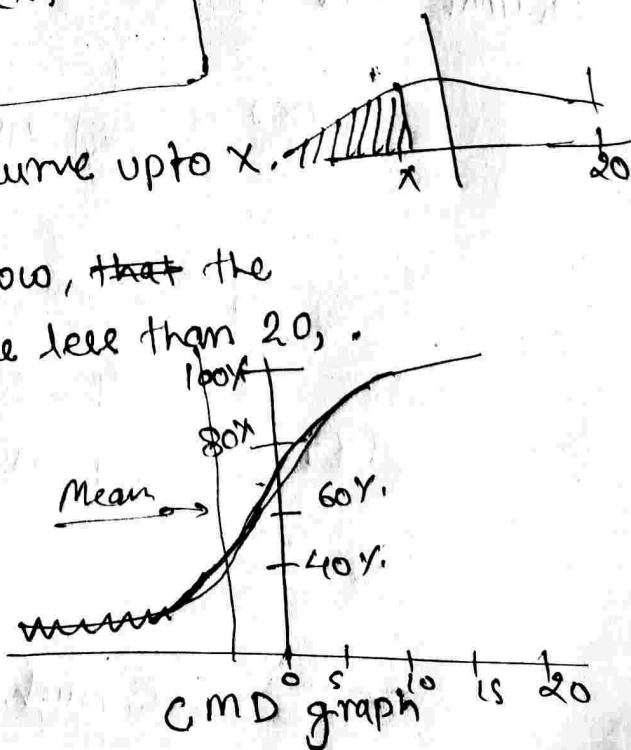
it tells you area under the curve upto x .

Helps you when you want to know, ~~that~~ the probability to get the value less than 20.

to calculate area b/w -1 to 1.

$\Rightarrow \int_{-\infty}^1 - \int_{-\infty}^{-1} \dots \Rightarrow$ desired Area.

~~equation~~ $\therefore \underline{\underline{CDF(1) - CDF(-1)}}$

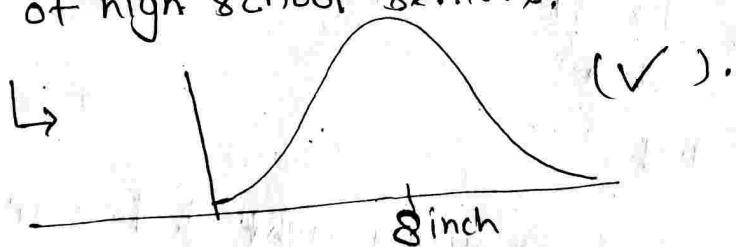


(4)

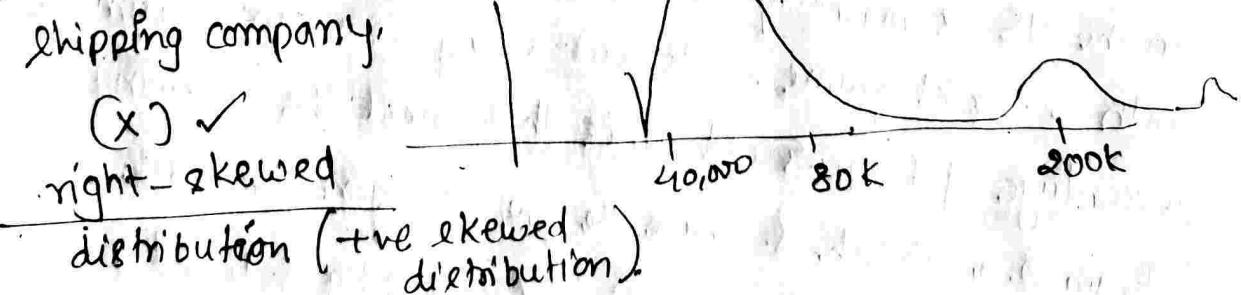
<u>Table!</u>	x	$(x-\mu)$	$(x-\mu)/\sigma$	$CDF(x)$	$P(x)$
		distance from mean	z-score How many standard deviation away from mean	cumulative distribution function	

example— Which of the following data sets is most likely to be normally distributed?

- (a) The hand span (measured from the tip of the thumb to the tip of the 5th finger) of a random sample of high school seniors.



- (b) The annual salaries of all employees of a large shipping company.



(mean is toward right, Hence it is called right-skewed)

example 2:

The grades on a statistic's mid-term for a high-school
are normally distributed with $\mu = 81$, $\sigma = 6.3$.

Calculate Z-score? for each of the following exam

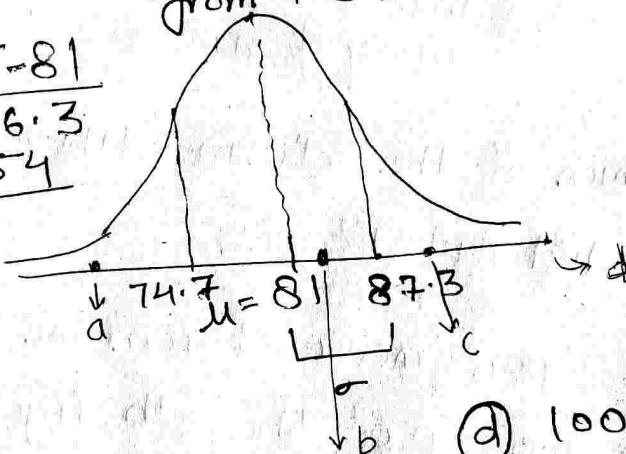
grades.

Sol: Z-score ← How many standard deviation away.
from the mean

(a) 65 $\Rightarrow \frac{65-81}{6.3} = -2.54$

(b) 83 $\Rightarrow \frac{83-81}{6.3} = 0.32$

(c) 93 $\Rightarrow \frac{93-81}{6.3} = \frac{x-\mu}{\sigma}$
 $\Rightarrow \underline{\underline{1.9}}$



(d) 100

$$\Rightarrow \frac{100-81}{6.3} \quad \begin{matrix} \text{distance from mean} \\ \text{magnitude (scale)} \end{matrix}$$

$$\Rightarrow 3.02$$

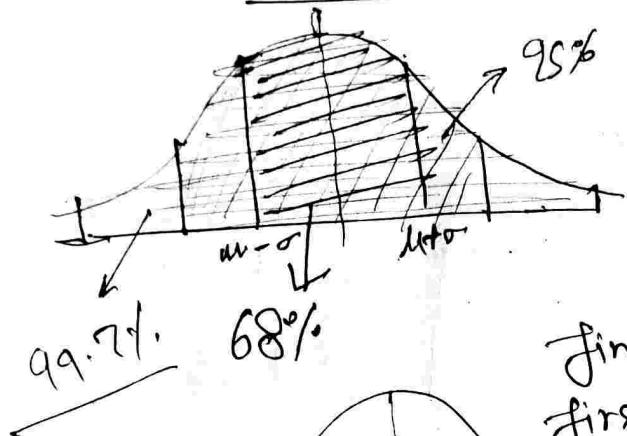
example: (3) Assume that the mean weight of 1-year-old girls
in the US is normally distributed with a mean of about 9.5 kg/m
with $\sigma \approx 1.1$ grams. Without using calculator, estimate the
percentage of 1-year old girls that meet the following conditions.

Draw & sketch & shade the proper region for each problem

(a) less than 8.4 kg (b) between 7.3 kg & 11.7 kg

(c) more than 12.8 kg.

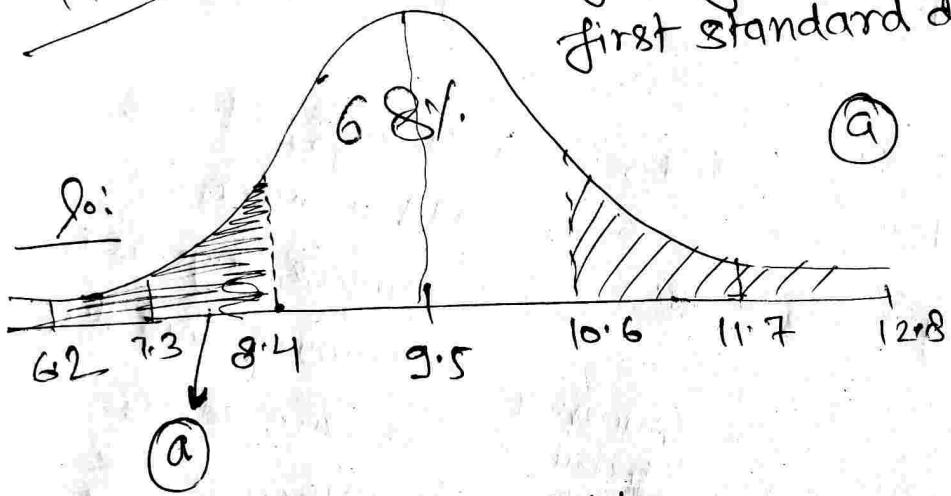
\Rightarrow Sol^e Empirical Rule:



(S)

$$\xrightarrow{\text{68-95-99.7}} \frac{99.7\%}{\text{b/w third deviation}}$$

68-95-99.7
 ↓
 68% chance
 b/w second deviation
 finding result b/w first standard deviation



(a) total area under curve

$$100 - 68\% = 32\%$$

left & right tail

(b) b/w 2 standard deviation is 95%

$$\Rightarrow \frac{32}{2} = 16\% \checkmark$$

(c) b/w 3 standard deviation = 99.7%

$$100 - 99.7\% = \frac{0.3\%}{2} = 0.15\% \checkmark$$

Note: for standard Normal distribution

$$\text{Mean} = 0$$

$$\text{standard deviation} = 1$$

* empirical Rule

$$68-95-99.7$$

Note! Z-score can apply to non-normal distribution.

Central Limit Theorem

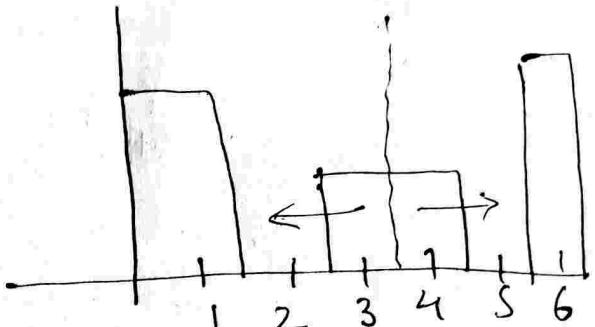
Sample size of $n = 4$

- $S_1 = [1, 1, 3, 6]$

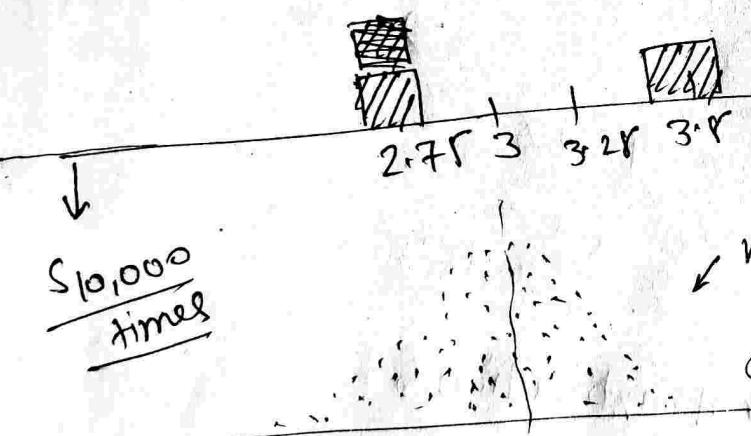
$$\bar{x}_1 = 2.75$$

- $S_2 = [3, 4, 3, 1]; \bar{x}_2 = 2.75$

- $S_3 = [1, 1, 6, 6]; \bar{x}_3 = 3.5$

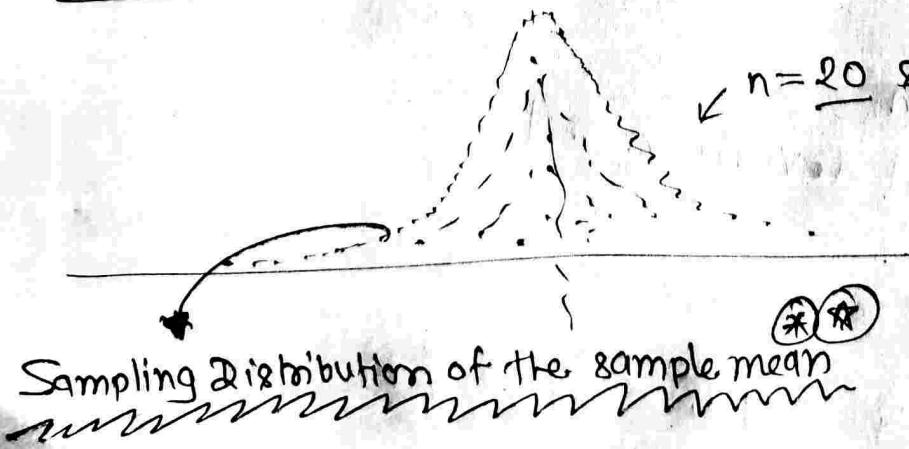


It says; if we add a bunch of actions samples together, assuming that they frequency distribution all have same distribution, or we take mean of all of those actions, and if we were to plot the frequency of those means, we will get Normal distribution



$n = 4$

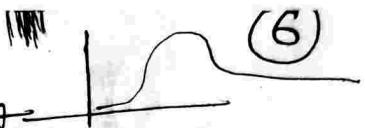
$n = 20$ sample size



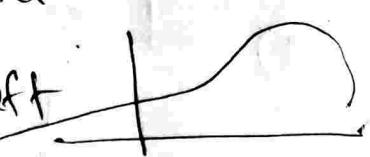
Sampling distribution of the sample mean

sample size $\rightarrow \infty$
↓ 10,000 times
Normal distribution

positive skewed :- long tail toward right

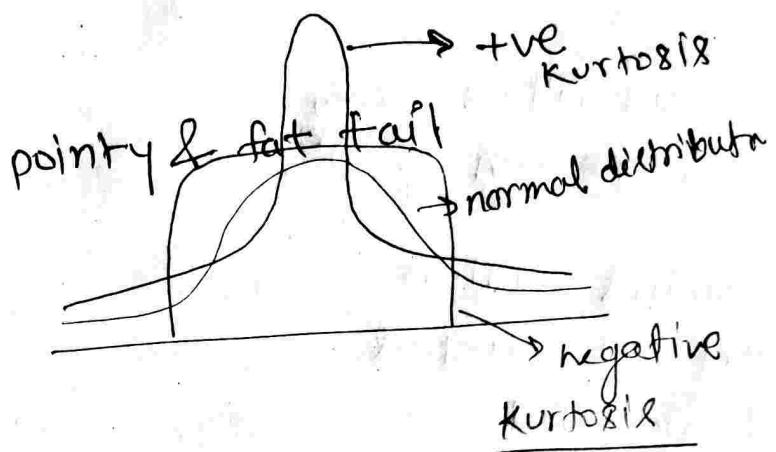


Negative skewed :- long tail toward left



Kurtosis

- +ve kurtosis is



Observations:

When sample size was

$$N=5$$

10,000 times

$$\text{skew} = 0.12$$

$$\text{kurtosis} = \begin{matrix} 0.04 \\ -0.26 \end{matrix}$$

$$N=25 \quad \checkmark \quad \begin{matrix} \text{Normal dist} \\ \text{tighter fit} \\ \text{to mean} \end{matrix}$$

$$\text{skew} = 0.03$$

$$\text{kurtosis} = -0.04$$

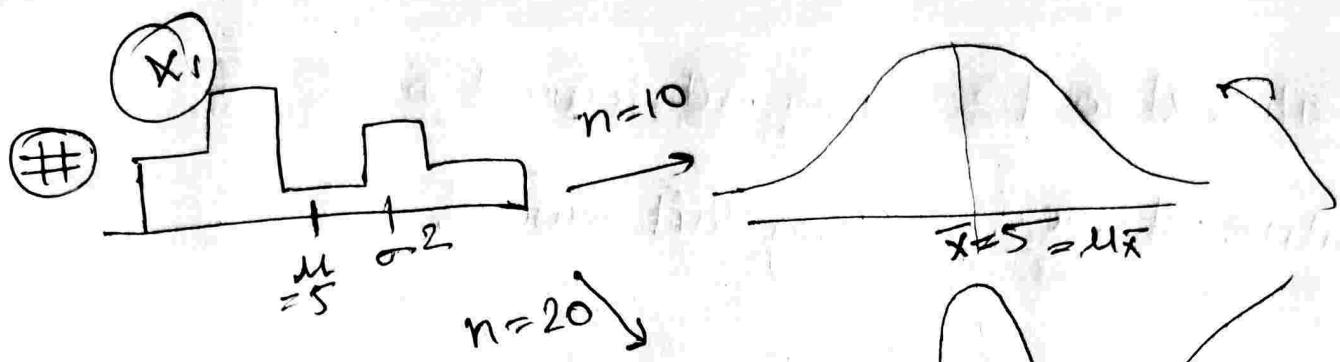
Conclusion: Central limit theorem says as you increase 'n' (sample size) → you will lead to normal distribution

$n \rightarrow \infty \rightarrow$ normal distribution

(sample size 4)

So; Sampling distribution } will → normal distribution.
of Sample mean

$n \uparrow \rightarrow$ standard deviation squeezed.

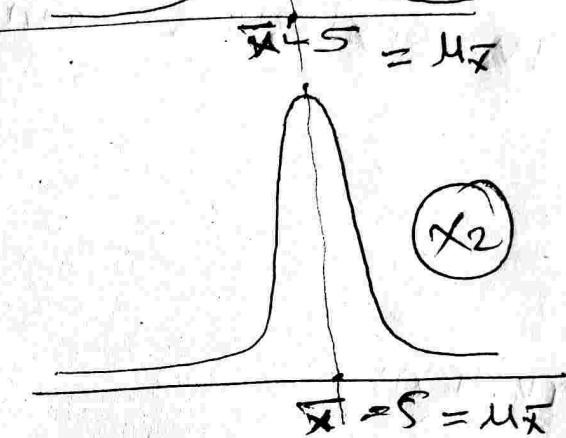


→ mean will be same,
irrespective of ' n '.

- standard deviation starts squeezing as ' n ' increases,

~~$n=10$~~

- $\mu_{\bar{x}} \rightarrow$ mean of sample mean



(#) Relation b/w probability distributed graph (X_1) & normal function's standard deviation → 's standard deviation

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} = S^2 \quad (\text{Variance})$$

$\sigma_{\bar{x}}$ → standard deviation of sampling distribution of sample mean
 (standard deviation of the mean)

or

Standard Error of the mean

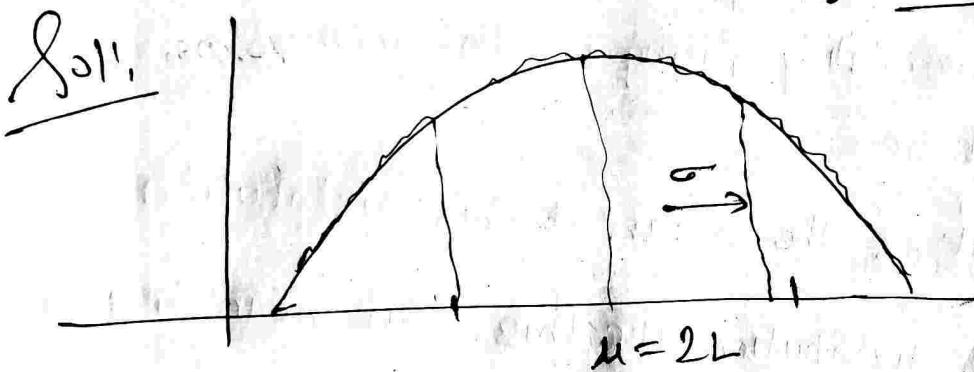
$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

standard deviation

example:

The average male drinks 2L of water when active outdoors (with standard deviation of 0.7L). You are planning a full day nature trip for 50 men & will bring 110L of water. What is the probability that you will run out?

$$\sigma = 0.7 \text{ L}$$



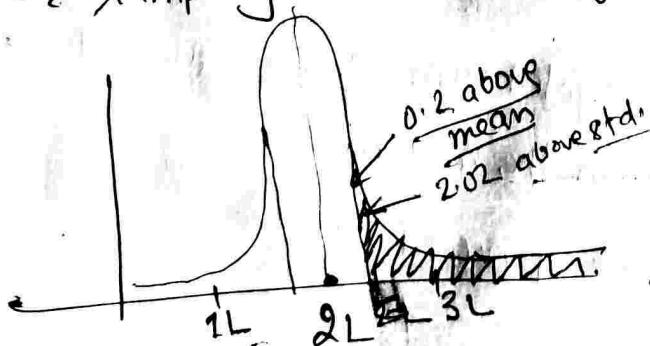
$\Rightarrow P(\text{runout})$

$$P(\text{use more than } 110 \text{ L})$$

$$P(\text{average water use per man}) > 2.2 \text{ L/m}$$

\Rightarrow Sampling distribution of the sample mean when $n=50$

$$\mu_{\bar{x}} = \mu = 2 \text{ L}; \quad \sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$$



$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \Rightarrow \underline{0.099}$$

$$\text{Z-score!} - \frac{2.2 - 2}{0.099} = \underline{2.020}$$

$P(\bar{x} \text{ will be more than } 2.020 \text{ std deviation above the mean})$

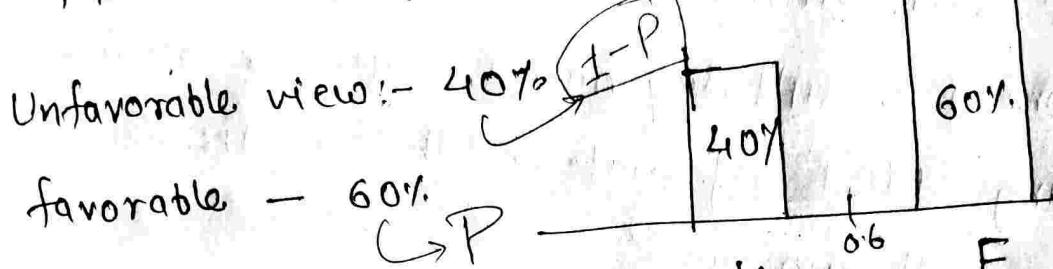
$$z(2.02) \rightarrow 0.9783 \rightarrow \underline{\text{area till } 0.9783 \Rightarrow}$$

$$1 - 0.9783 \Rightarrow \underline{2.17\%}$$

(Simplest Case of Binomial Distribution)

Bernoulli Distributions

In population survey, about president candidate



Pick a random member of that population & say what is the expected favorability rating of that member?

What would it be?

or
what is the mean of this distribution?

* And for discrete distribution like this, your mean or expected value is just going to be

probability weighted sum of the different values

Suppose $U=0, F=1$

$$\mu = 0.4 * 0 + 0.6 * 1 = \underline{\underline{0.6}}$$

Variance: $\sigma^2 =$ probability weighted sum of the squared distance from the mean.

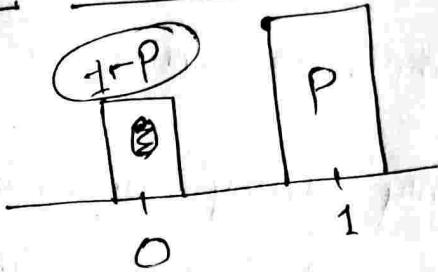
$$= 0.4(0-0.6)^2 + 0.6(1-0.6)^2$$

$$\sigma^2 = \underline{\underline{0.24}} ; \sigma = \underline{\underline{0.49}}$$

(8)

General formula for Bernoulli Distribution:

favorable $\rightarrow 60\% \quad P$
 Unfavorable $\rightarrow 40\% \quad 1-P$



$$\mu = (\pm P) \cdot 0 + P(1)$$

$$\boxed{\mu = P} \quad \text{Mean.}$$

$$\sigma^2 = (1-P) \cdot (0 - \frac{\mu}{P})^2 + P(1-P)^2$$

$$\sigma^2 = (1-P)(-P)^2 + P(1-P)^2$$

$$\sigma^2 = (1-P)(P)^2 + P(1-P)^2$$

$$\sigma^2 = P^2 - P^3 + P - 2P^2 + P^3$$

$$\boxed{\sigma^2 = P(1-P)} \quad \text{Variance}$$

$$\boxed{\sigma = \sqrt{P(1-P)}} \quad \rightarrow \text{standard deviation.}$$

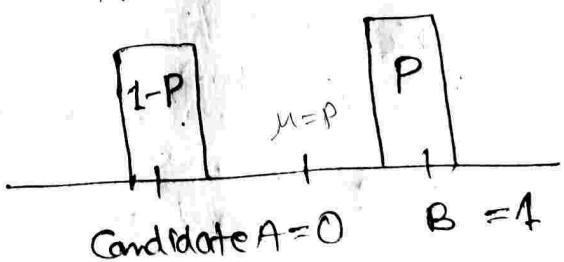
Margin of Error:-

from Bernoulli Distribution:-

$$\mu = P$$

Random Survey :

Presidential Election



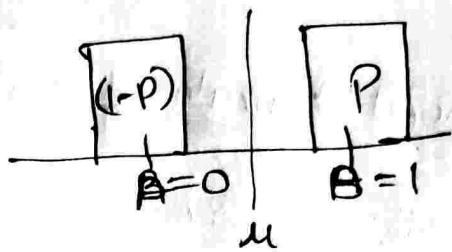
↳ to make it manipulatable
mathematically

$$A=0;$$

$$B=1$$

→ Margin of Error

↳ Sample of population



→ Randomly survey 100 people:

57 people → vote for A $\Rightarrow 0$

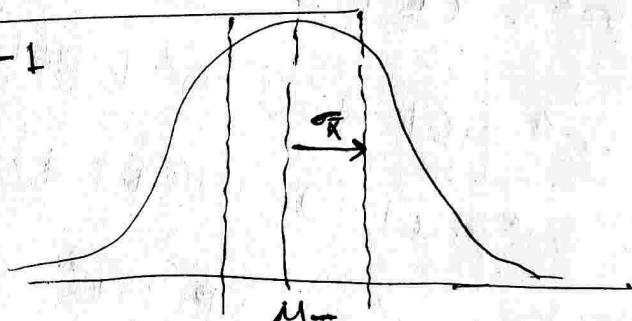
43 people → — → B $\Rightarrow 1$

Sample Mean: $\frac{57*0 + 43*1}{100} = 0.43$

$S^2: \frac{57(0 - 0.43)^2 + 43(1 - 0.43)^2}{100 - 1} = 0.2475$

$S = 0.50$

Sample standard deviation



sampling dist. of sample mean

$\mu_{\bar{x}} = \boxed{\mu} = P$

↳ standard deviation of the sampling ~~dist.~~ of sampling mean

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \rightarrow \text{sample size}$$

→ population std. deviation
to know this

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{100}}$$

↓
We need to actually survey 100 million people,

$$\sigma_{\bar{x}} \approx \frac{0.50}{10} \rightarrow \text{depend on sample}$$

↓
so, we will use our sampling std. deviation

$\boxed{S \approx 0.50}$

best estimate for σ

$$\sigma_{\bar{x}} \approx 0.05$$

↳ best estimation of

$$\sigma_{\bar{x}} \approx 0.05 \rightarrow 5\%$$

(9)

⇒ Now, come up with an interval around the sample mean where I'm reasonably confident using all my estimates that there's 95% chance that true mean is within two standard deviations

$$\downarrow \boxed{\mu = p}$$

⇒ ~~look at~~
 $P(\bar{x} \text{ is within } 2\sigma_{\bar{x}} \text{ of } \mu_{\bar{x}}) \approx 95\%$

↓
 random sample
 mean

$P(\mu_{\bar{x}} \text{ is within } 2\sigma_{\bar{x}} \text{ of } \bar{x}) \approx 95.4\%$

↓
 mean of sampling dist.

$P(p \text{ is within } \pm 0.05 \text{ of } \bar{x}) \approx 95\%$

$P(p \text{ is within } 0.10 \text{ of } \bar{x}) \approx 95\%$

$P(p \text{ is within } (0.43 \pm 0.10)) \approx 95\%$

95% Confidence interval of 33% to 53%

↑
P range

43% will vote for candidate B

57% will vote for A

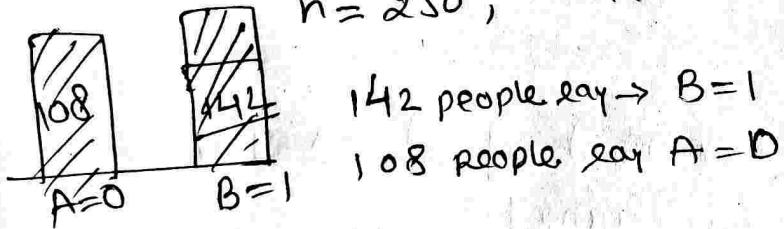
margin of error: - 10% ✓

exampleConfidence IntervalQues:

In a local teaching district a technology grant is available to teachers in order to install a cluster of four computers in their classrooms.

From the 6250 teachers in the district, 250 were randomly selected and asked if they felt that computers were an essential teaching tool for their classroom. Of those selected, 142 teachers felt that computers were an essential teaching tool.

Sol: sample size = 250,
 $n = 250$, $N = 6250$



- 250
142
108 Questions:
- ① Calculate a 99% confidence interval for the proportion of teachers who felt that computers are an essential tool.
 - ② How could the survey be changed to narrow the confidence interval but to maintain the 99% confidence interval?

(10)

Sol:

$$\text{sample mean } \bar{x} = \frac{108(0) + 142(1)}{250} = 0.568.$$

$$\text{sample variance } s^2 = \frac{108(0 - 0.568)^2 + 142(1 - 0.568)^2}{250 - 1}$$

$$= \frac{0.322624 + 26.50}{249} = \frac{0.1077}{0.246}$$

$$\text{std deviation } s = \sqrt{0.3282}$$

Sample

~~Sampling distribution of sampling mean~~

$\sigma_{\bar{x}}$ ← population std deviation

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \leftarrow \text{sample size.}$$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{250}} \quad \sigma_{\bar{x}} \approx \frac{0.50}{\sqrt{250}} = \frac{0.031}{\sqrt{250}}$$

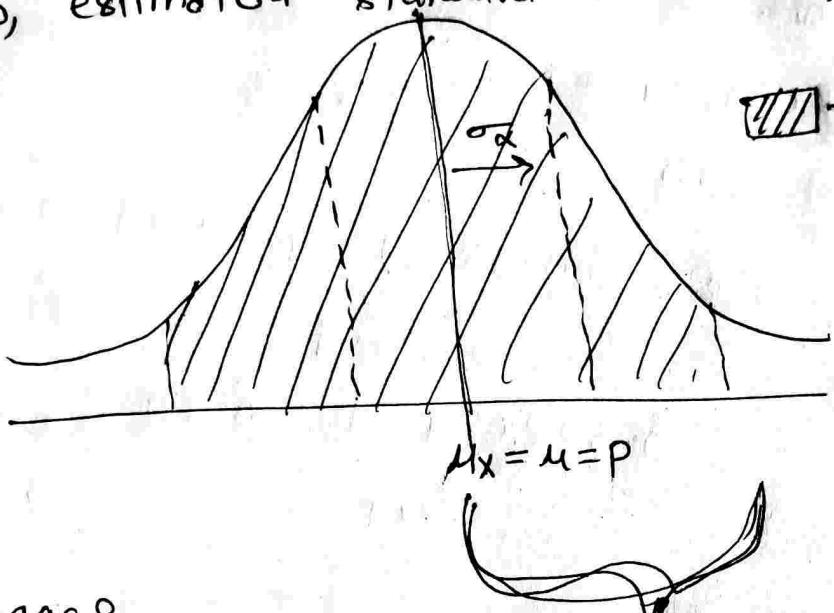
Estimated

~~Std. deviation~~
of Sampling dist. of
Sampling mean

$$\sigma_{\bar{x}} \approx 0.031$$

↳ How many standard deviations away from the mean do we have to be that we can be 99% confident that any sample from the sampling distribution will be in that interval?

so, estimated standard deviation $\sigma_{\bar{x}} \approx 0.031$



→ 99% confidence zone //
to basically
find its
boundary
in std. deviations
terms.

Check: 0.9958

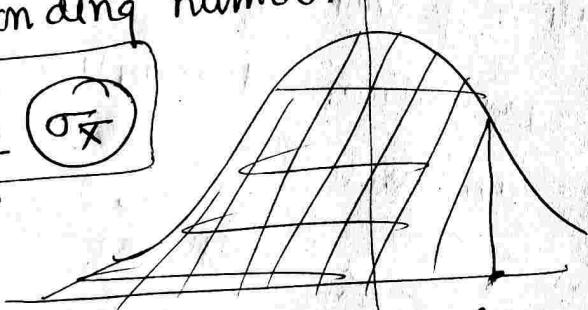
value in z-score table ←

$$\frac{99}{2} = 0.495$$

as it symmetric
around 0.

↓
& corresponding number

2.58 $\sigma_{\bar{x}}$



99% chance that a random \bar{x} is within $2.58 \sigma_{\bar{x}}$ of p

"Confident" 99% chance that \bar{x} is within

$$\frac{0.08}{\bar{x}} \text{ of } \sigma_{\bar{x}}$$

$$0.568$$

$$\rightarrow \bar{x} \pm 0.08$$

$$(0.488 \text{ to } 0.648)$$

$$48.8 \text{ to } 64.8 \%$$

⑥: take more samples.

example 7 patients' blood pressures have been measured after having been given a new drug for 3 months. They had blood pressure increases of 1.5, 2.9, 0.9, 3.9, 3.2, 2.1 & 1.9. Construct a 98% confidence interval for the true expected blood pressure increase for all patient in a population.

$$\text{So, sample size} = 7$$

$$\bar{x} = \frac{1.5 + 2.9 + 0.9 + 3.9 + 3.2 + 2.1 + 1.9}{7} = \frac{16.81}{7} = 2.342$$

$$S = 1.04$$

So, we estimate the true standard deviation of the population with our sample standard deviation. $\sigma \approx S = 1.04$

But here only 7-sample size, not so good estimate.
because n is small. ($n < 30$)

t-distribution \rightarrow



t-distribution

$$\sigma_x = \frac{\sigma}{\sqrt{n}} \approx \frac{S}{\sqrt{n}}$$

sample size is small

~~normal distribution~~

t-distribution ✓

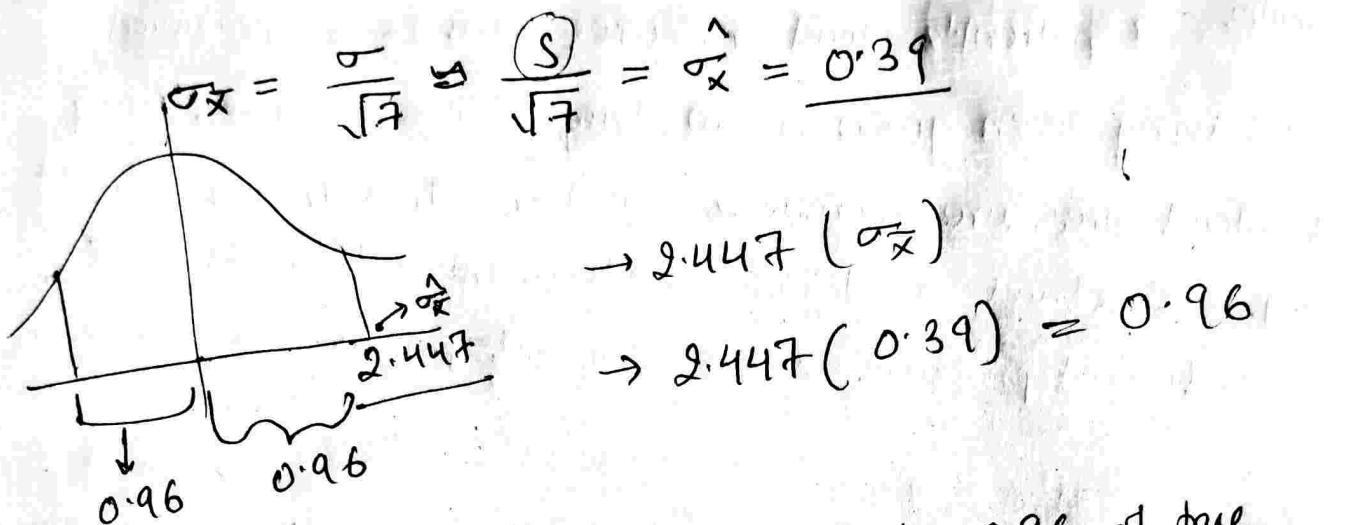
for t-distribution

use t-table!

$n=7$; 80 ($n-1$) degree of freedom

Two-sided row \rightarrow 95%, 6 degree

t-table \rightarrow $\frac{2.447 \text{ approximated standard deviation}}{\sigma_x}$



Confident 95% chance, $\bar{x} = 2.34$ is within 0.96 of true population mean. (μ).

\Rightarrow 95% chance μ is within 0.96 of 2.34

④ Hypothesis testing & P-values

example: A neurologist is testing a effect of a drug on response time by injecting 100 rats with a unit dose of the drug, subjecting each to neurological stimulus, and recording its response time.

mean response time for rats not injected with the drug is 1.2 seconds.

Mean of 100 injected rats' response times is 1.05 sec with sample standard deviation of 0.5 seconds.

Do you think that the drug has an effect on response time?

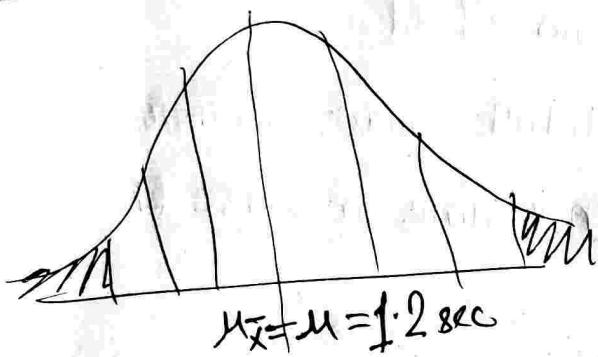
(12)

Sol:

H_0 : Drug has no effect $\Rightarrow \mu = 1.28\text{sec}$ (even w/o drug)

H_1 : Drug has an effect $\Rightarrow \mu \neq 1.25$ when the drug is given.

Assume H_0 is true:



$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{100}} \approx \frac{0.5}{\sqrt{100}} = \frac{0.5}{10}$$

$$\sigma_{\bar{x}} \approx 0.05$$

What is the probability of getting 1.05 seconds?

Or another way

is how many standard deviation away from this mean is 1.05 sec so that we can use empirical rule.

$$\begin{aligned} & Z\text{-score} \\ & \Rightarrow \frac{(1.2 - 1.05)}{0.05} \end{aligned}$$

$$\underline{\underline{Z = 3}}$$

3 standard deviation away from mean

probability of getting result (-ve direction) \rightarrow redem

probability of getting result (-ve direction) \rightarrow redem

3-deviation area $= 99.7\%$

$$1 - 99.7\% = \underline{\underline{0.3\%}} \quad \underline{\underline{.003}}$$

~~So~~ So, H_0 should be rejected as it has

Value 0.3%, probability of getting a result more extreme than this given the null hypothesis is called p-value.

$$\text{p-value (probability value)} = \underline{\underline{0.003}}$$

#



so, null hypothesis has no effect.

H_0 has an effect; but didn't know whether the drug would lower the response time or raise the response time

In this situation, where we're really just testing to see if it had an effect, whether an extreme positive effect, or an extreme negative effect, both considered as an effect. called a two-tailed test.

If we want to do one-tailed test.

$H_0 \rightarrow$ Drug has no effect $\mu_{\text{drug}} = 1.2 \text{ sec}$

$H_1 \rightarrow$ Drug lowers response $\mu_{\text{drug}} < 1.2 \text{ sec}$

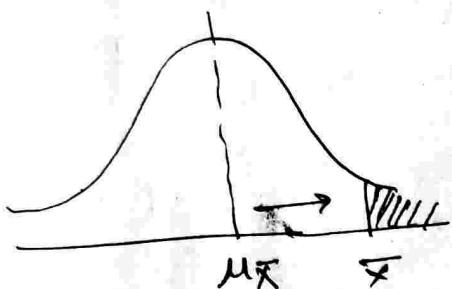
\hookrightarrow One about
One direction

$H_0 \rightarrow$ 0.15% \rightarrow 0.0015

→ does not lower

(13)

*
#

Z - statistic

$$\frac{\bar{x} - \mu_x}{\sigma_x} \Rightarrow \frac{\bar{x} - \mu_x}{\frac{\sigma}{\sqrt{n}}}$$

$$\sigma_Z = \frac{\sigma}{\sqrt{n}}$$

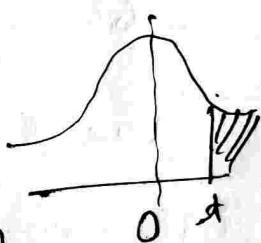
$$\Rightarrow Z \sim \frac{\bar{x} - \mu_x}{\frac{s}{\sqrt{n}}} ; \text{ 'ok' if } n > 30$$

↳ Z-table

normal distributed ↗

t-statistic:

$$\frac{\bar{x} - \mu_x}{\frac{s}{\sqrt{n}}} \text{ if } n \text{ is small}$$

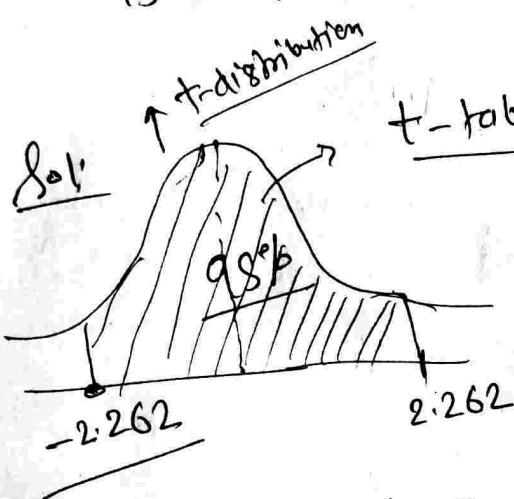
↳ t-table

↗ t-distribution

Type-I error :- Rejecting null hypothesis H_0 even though it is true

④ example: find 95% confidence interval $n=10$

15.6 16.2 22.5 20.8 16.4 19.4 16.6 17.9 12.7
13.9



t-table: two-sided row

$n=10$
degree = 9

t-value = 2.262

there is 95% chance, that if you pick random T-statistic
is going to be less than 2.262

$$-2.262 < t < 2.262$$

$$19.3 > \mu > 15.04$$

\uparrow 95% chance

sample mean 6 17.17 $\Rightarrow \bar{x}$

$$S = 2.98$$

t-statistic \Rightarrow

$$\frac{\bar{x} - \mu}{\frac{S}{\sqrt{n}}} < 2.262$$

$$-2.267 < \frac{17.17 - \mu}{\frac{2.98}{\sqrt{10}}} < 2.268$$

$$-2.267 * 0.9423 < 17.17 - \mu < 2.268 * 0.942$$

$$-2.13 < 17.17 - \mu < 2.13$$

$$2.13 < \mu - 17.17 < -2.13$$

$$2.13 + 17.17 > \mu > -2.13 + 17.17$$

multiple (-)

Hypothesis test

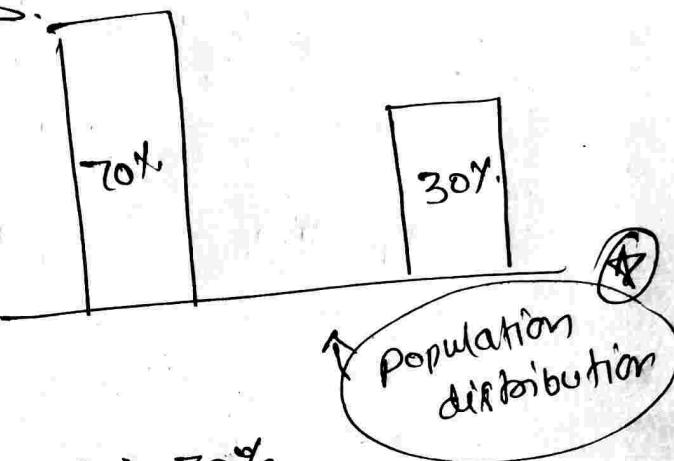
(14)

We want to test the hypothesis that more than 30% of U.S. households have internet access (with a significance level 5%). We collect a sample of 150 households & find that 57 have access.

Solve

Bernoulli Distribution

$$\text{Variance: } \frac{\sigma^2 = P(1-P)}{\mu = P}$$

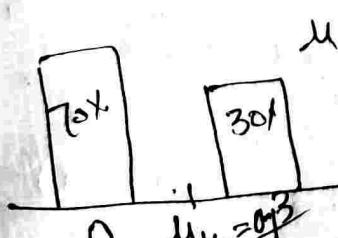


$$H_0 : P \leq 30\% \quad H_1 \Rightarrow P > 30\%$$

Assume P-value based on null hypothesis (proportion based on null hypothesis for population) & given that assumption what is the probability that 57 out of 150 of our samples actually have internet. And if that probability is less than 5% → reject H_0

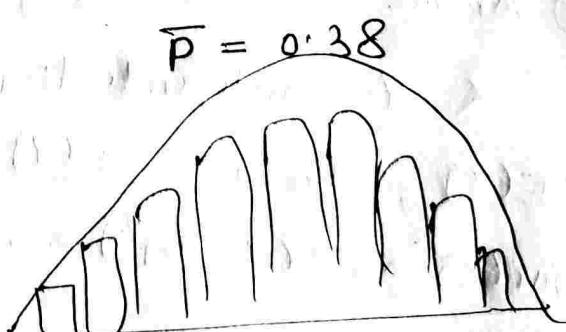
Assume H_0 true

$$P_{H_0} = 0.30$$



$$\begin{aligned} \mu &= P \\ \sigma_{H_0} &= \sqrt{P(1-P)} \\ &= \sqrt{(0.3)(0.7)} \\ &= \underline{0.42} \end{aligned}$$

distribution of sample proportion



Binomial distribution

if $np \geq 5$

$$n(1-p) \geq 5$$

we can assume,
distribution is normal

normal distribution

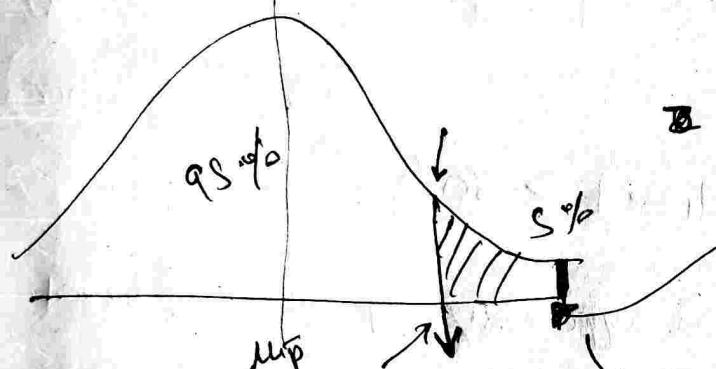
$$\sigma_{\bar{P}} = \frac{\sigma}{\sqrt{n}} = \frac{\sqrt{.21}}{\sqrt{150}} = 0.037$$

population

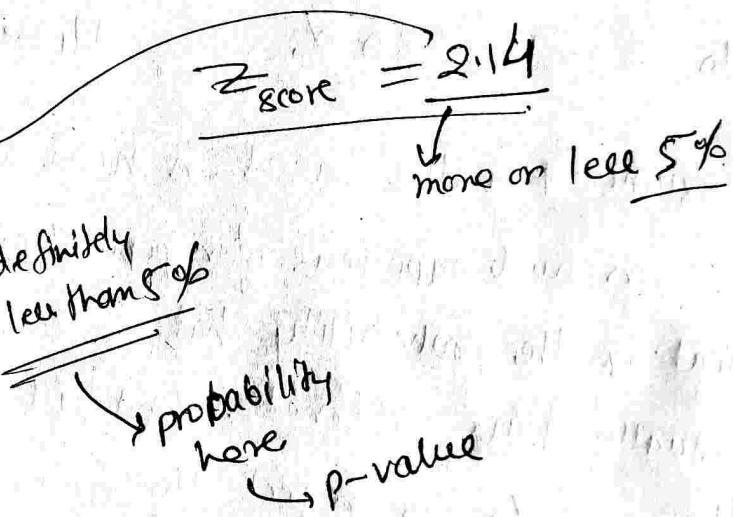
$$\mu_{\bar{P}} = 0.3$$

how much away from standard

$$\frac{\bar{P} - \mu_{\bar{P}}}{\sigma_{\bar{P}}} = \frac{0.38 - 0.3}{0.037} = \frac{0.08}{0.037}$$



what z value which has
ask



Variance of differences of random variables &

X Y independent random variables.

$$E(X) = \mu_X ; E(Y) = \mu_Y ;$$

$$\text{Var}(X) = E((X-\mu_X)^2) = \sigma_X^2 ; \text{Var}(Y) = E((Y-\mu_Y)^2) = \sigma_Y^2$$

third random variable

$$Z = X + Y$$

$$E(Z) = ? = E(X+Y) = E(X) + E(Y)$$
$$\mu_Z = \mu_X + \mu_Y$$

(15)

$$\text{variance } (Z) = \text{Var}(X) + \text{Var}(Y)$$



$$\sigma_Z^2 = \sigma_{x+y}^2 = \sigma_x^2 + \sigma_y^2$$

$$\& \sigma_{x-y}^2 = \sigma_x^2 + \sigma_y^2$$

$$\sigma_A^2 = \sigma_{x-y}^2 = \sigma_{x+(-y)}^2 = \sigma_x^2 + \sigma_y^2$$

same ~~as it~~
is square

$$\sigma_y^2 = \text{Var}(-Y) = E((-Y - E(-Y))^2)$$

$$= E((-1)^2 (Y + E(Y)^2))$$

$$\therefore E(-Y) = -E(Y)$$

$$E((Y - E(Y))^2) = \sigma_Y^2$$

#

~~say~~ Difference of sample means distribution!

random \bar{X} :

$$\sigma_X^2$$

$$\bar{Y}$$

$$\sigma_Y^2$$

Sampling
distribution
size = n

$$\mu_{\bar{X}} = \mu_X$$

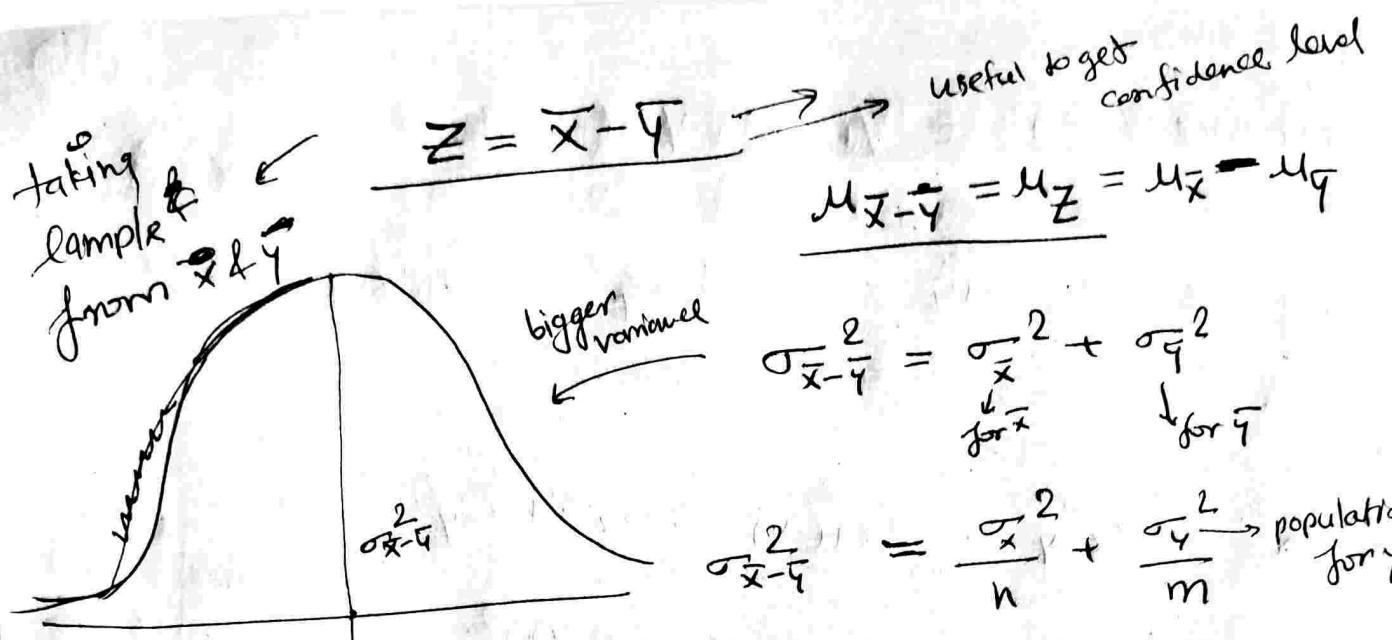
Sampling
distribution of \bar{Y}

$$\text{size} = m$$

$$\mu_{\bar{Y}} = \mu_Y$$

$$\sigma_{\bar{X}}^2 = \frac{\sigma_X^2}{n} ; \quad \text{population variance}$$

$$\sigma_{\bar{Y}}^2 = \frac{\sigma_Y^2}{m}$$

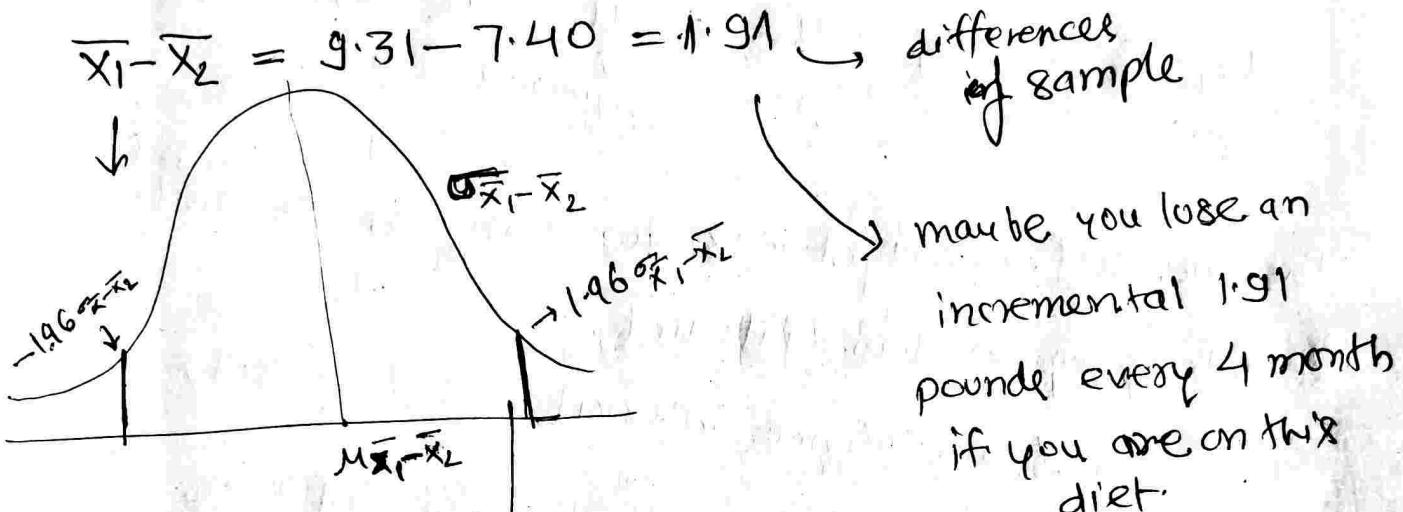


$$\sigma_{\bar{x}-\bar{y}} = \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}$$

example We're trying to test whether a new, low fat diet actually helps obese people ~~lose~~ lose weight.
 100 randomly assigned obese people are assigned to group 1 and put on the low fat diet. Another 100 randomly assigned obese people are assigned to group 2 and put on diet of approximately the same amount of food, but not as low in fat. After 4 months, the mean weight loss was 9.31 lbs. for group 1 ($s = 4.67$) & 7.40 lbs. ($s = 4.04$) for group 2.

(16)

Low-fat! $\bar{x}_1 = 9.31 \text{ lbs}$ { given $\frac{\text{Control!}}{\bar{x}_2} = 7.40$
 $s_1 = 4.67$ } $s_2 = 4.04$



$$\sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}$$

create 95%
interval

↓
how many
std. deviations

95% chance that
1.91 is within
 $1.96 \sigma_{\bar{x}_1 - \bar{x}_2}$ of the
 $M(\bar{x}_1 - \bar{x}_2)$

↓
look \approx table

↓
gives cumulative
value

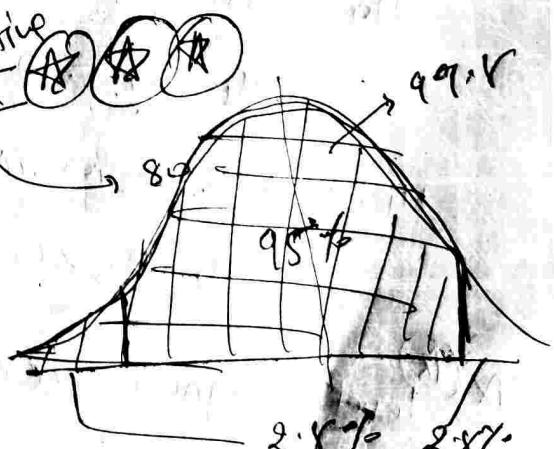
1.96

So get a 95% confidence interval around this number

do we lose weight?

$$\begin{aligned}\sigma_{\bar{x}_1 - \bar{x}_2} &= \sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{m}} \\ &= \sqrt{\frac{(4.67)^2}{100} + \frac{(4.04)^2}{100}} = 0.617\end{aligned}$$

So $1.96 \times 0.617 = 1.21$
So 95% level $\bar{x}_1 - \bar{x}_2 \in [1.91 \pm 1.21]$ ✓



so find Z for
95 + 2.5% = 97.5%

Confidence 95% confidence interval for $\mu_{\bar{x}_1} - \mu_{\bar{x}_2}$

$$1.91 \pm 1.21$$

$$\mu_{\bar{x}_1} - \mu_{\bar{x}_2} = \mu_{\bar{x}_1} - \mu_{\bar{x}_2} = \boxed{\mu_1 - \mu_2}$$

95% confidence level of this

The true diff. in weight 1088 b/w going on the low-fat diet & not going on low-fat diet.

95% confidence interval

$$0.7 - 3.12$$

means 95% confidence that you will lose some weight %

Hypothesis test for above example:

H_0 : low-fat does nothing

$$\mu_1 - \mu_2 = 0 \Rightarrow \mu_{\bar{x}_1} - \mu_{\bar{x}_2} = 0 \Rightarrow \mu_{\bar{x}_1} - \mu_{\bar{x}_2} = 0$$

H_1 : improvement

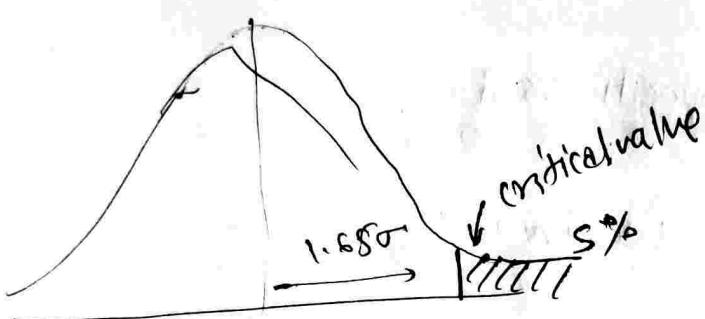
$$\mu_1 - \mu_2 > 0 \Rightarrow \mu_{\bar{x}_1} - \mu_{\bar{x}_2} > 0$$

Level of significance (threshold level or alpha) \Rightarrow ~~0.05~~

significance level of 95% $\underline{\alpha = 0.05}$ 5%

$$\Rightarrow \underline{0.05}$$

(17)



$$\mu_{\bar{X}_1 - \bar{X}_2} = 0$$

step ① critical z-value for 95% \rightarrow critical value.

from Z-table 1 \Rightarrow 1.65

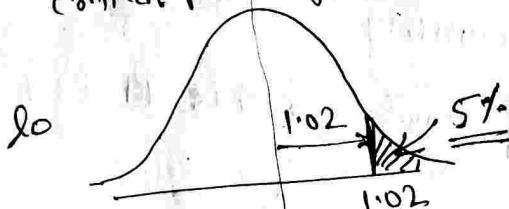
step ② $\sigma_{\bar{X}_1 - \bar{X}_2} = ?$

$$\sigma_{\bar{X}_1 - \bar{X}_2} = 0.617 \text{ (already calculated)}$$

step ③ so, $Z \text{ value} * \sigma_{\bar{X}_1 - \bar{X}_2}$

$$1.65 * (0.617) = \underline{\underline{1.02}}$$

critical point from mean is 1.02 far



step ④ $\bar{X}_1 - \bar{X}_2 = 1.91$; where it lies; 95% or 5% interval

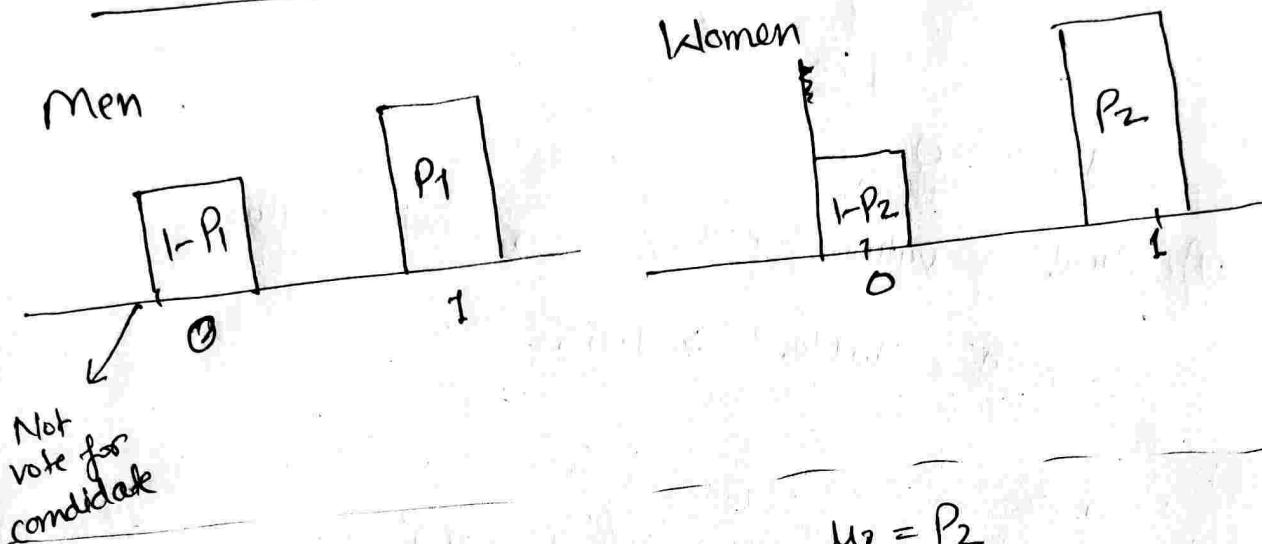
it's in 5% interval

so H_0 has less than 5% probability!

less than α .

H_0 (rejected)

④ Comparing population proportions 1 :-



$$\mu_1 = P_1$$

$$\sigma_1^2 = P_1(1-P_1) =$$

$$\mu_2 = P_2$$

$$\sigma_2^2 = P_2(1-P_2)$$

? Get the meaningful difference b/w the way that the men will vote & women will vote.

$\frac{P_1 - P_2}{\sqrt{P_1(1-P_1) + P_2(1-P_2)}}$ is meaning & come up 95% confidence level for this $(P_1 - P_2)$

$$n = 1000 \text{ men}$$

$$642 \rightarrow \text{vote} \rightarrow \text{for } 1$$

rest $\rightarrow 0$

1000 women

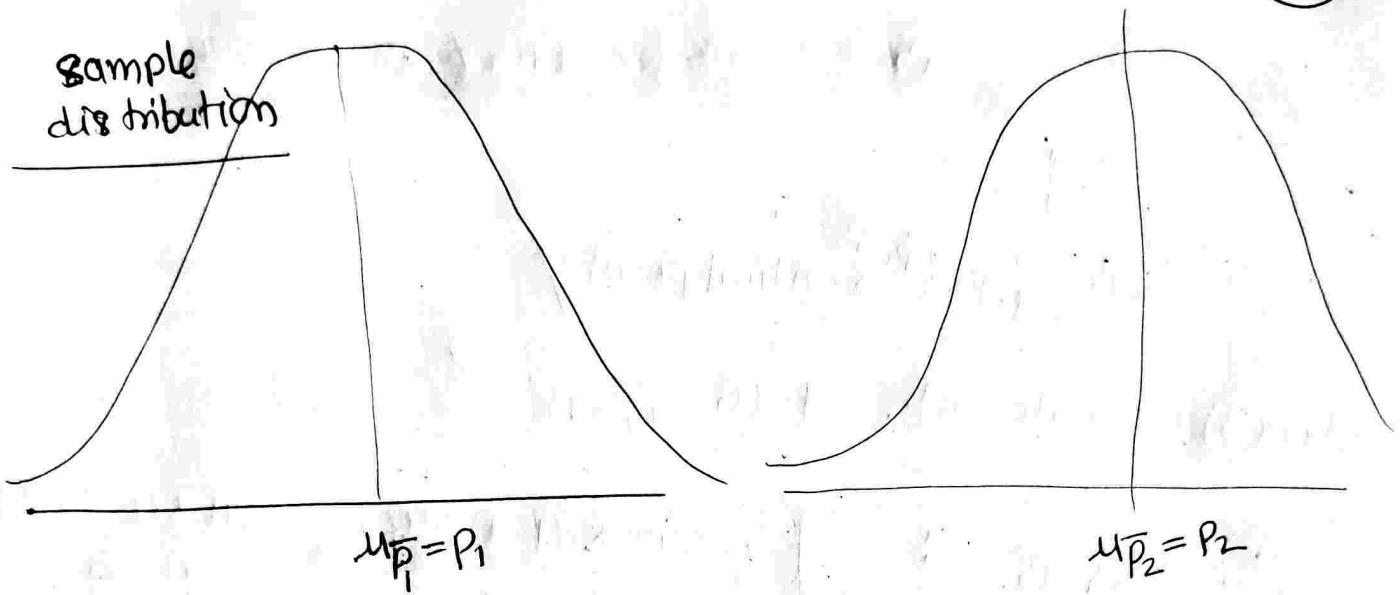
591 $\rightarrow 1$
rest $\rightarrow 0$

$$\bar{P} = 0.642$$

sample mean
or proportion

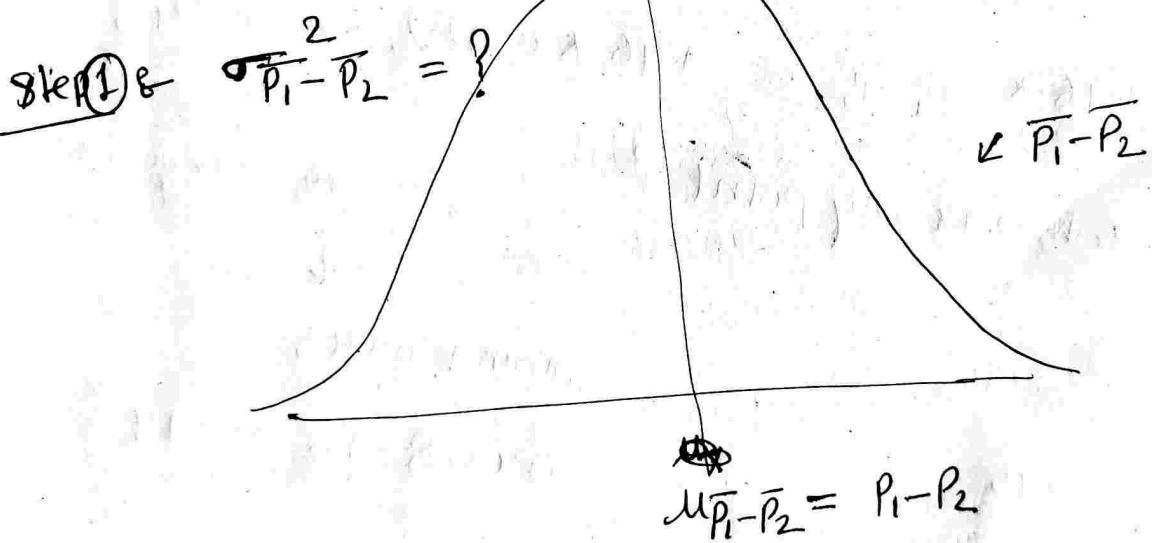
$$\bar{P} = 0.591$$

(18)



$$\sigma_{\bar{P}_1}^2 = \frac{P_1(1-P_1)}{1000}$$

$$\sigma_{\bar{P}_2}^2 = \frac{P_2(1-P_2)}{1000}$$

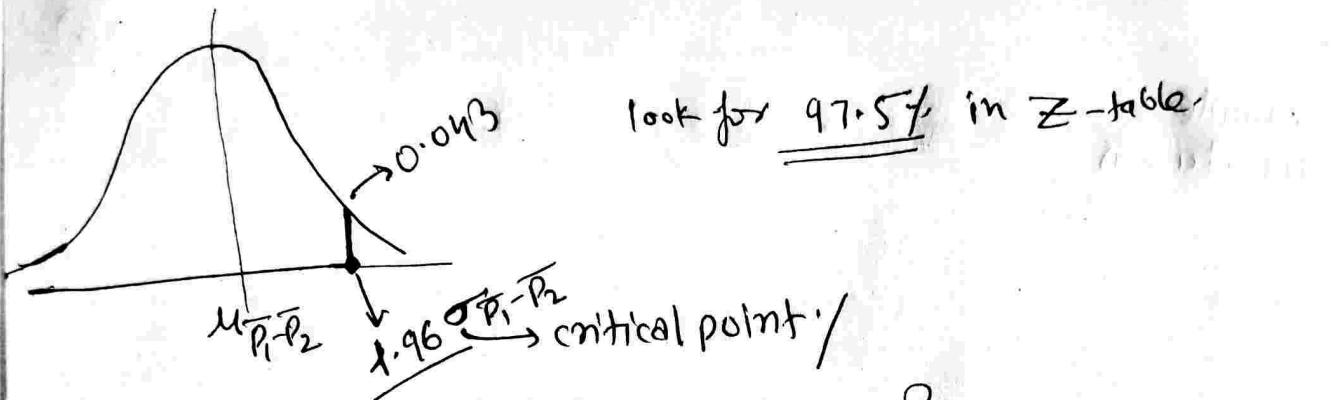


$$\sigma_{\bar{P}_1 - \bar{P}_2}^2 = \sigma_{\bar{P}_1}^2 + \sigma_{\bar{P}_2}^2$$

$$\sigma_{\bar{P}_1 - \bar{P}_2} = \sqrt{\frac{P_1(1-P_1)}{1000} + \frac{P_2(1-P_2)}{1000}}$$

step 2 critical point for 95% interval

$$\bar{P}_1 - \bar{P}_2 = 0.051 \leftarrow \text{our point}$$



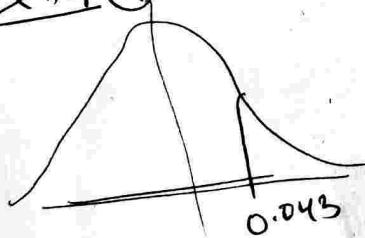
Step ③ calculate $1.96 \sigma_{\bar{P}_1 - \bar{P}_2}$?

$$\therefore \sigma_{\bar{P}_1 - \bar{P}_2} \approx \sqrt{\frac{0.642 * (1 - 0.642)}{1000} + \frac{0.591 * (1 - 0.591)}{1000}}$$

$$\sigma_{\bar{P}_1 - \bar{P}_2} \approx 0.022$$

~~∴~~ $1.96 * \sigma_{\bar{P}_1 - \bar{P}_2} = 1.96 * 0.022 = \underline{0.043}$

Step ④ where is my point?



mean ± 0.043

$0.008 \text{ to } 0.094$

$\& 0.051$ lie within

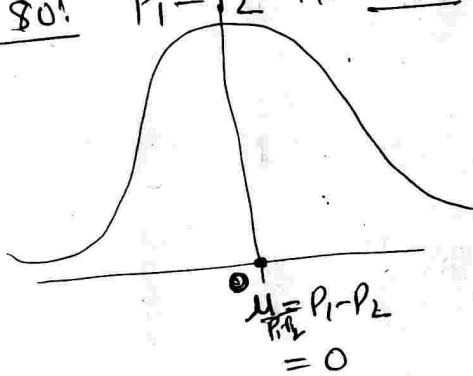
so 95% confidence ✓

95% chance men are more likely
to vote for candidate than women ✓

(19)

Hypothesis test: H_0 : No difference $P_1 - P_2 = 0$; H_1 : $P_1 - P_2 > 0$; $P_1 \neq P_2$ $\alpha = 5\%$

Assume H_0 is true; figure out probability for
 $P(\bar{P}_1 - \bar{P}_2 | H_0) < 5\%$. ~~(*)~~ ~~(*)~~
then reject H_0

So! $P_1 - P_2$ is zero

$$\bar{P}_1 = 0.642$$

$$\bar{P}_2 = 0.591$$

$$\bar{P}_1 - \bar{P}_2 = 0.051$$

my
our point
on graph

critical point

probability
to get this

$$Z\text{-score} = \frac{0.051 - 0}{\sqrt{\bar{P}_1 + \bar{P}_2}}$$

$$\sigma_{\bar{P}_1 - \bar{P}_2} = \sqrt{\frac{P_1(1-P_1)}{1000} + \frac{P_2(1-P_2)}{100}}$$

$$\text{in } H_0 \quad P_1 = P_2$$

$$\sigma = \sqrt{\frac{1}{1000} (P(1-P) + P(1-P))}$$

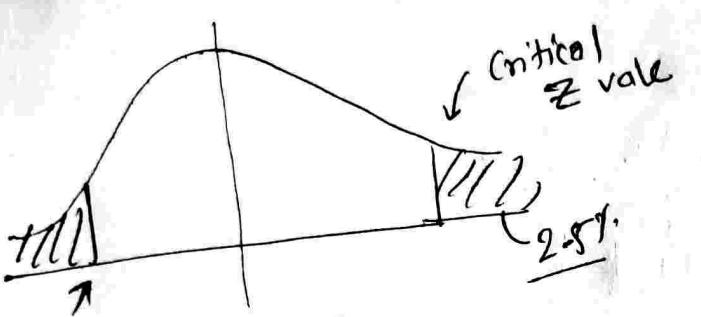
$$\sigma_{\bar{P}_1 - \bar{P}_2} = \sqrt{\frac{2P(1-P)}{1000}}$$

so our sample $1000 + 1000 = 2000$

$$Z = \frac{0.051}{0.0217} = 2.38$$

↓
std. deviat away

$$\bar{P} = \frac{642 + 591}{2000} = 0.6165$$

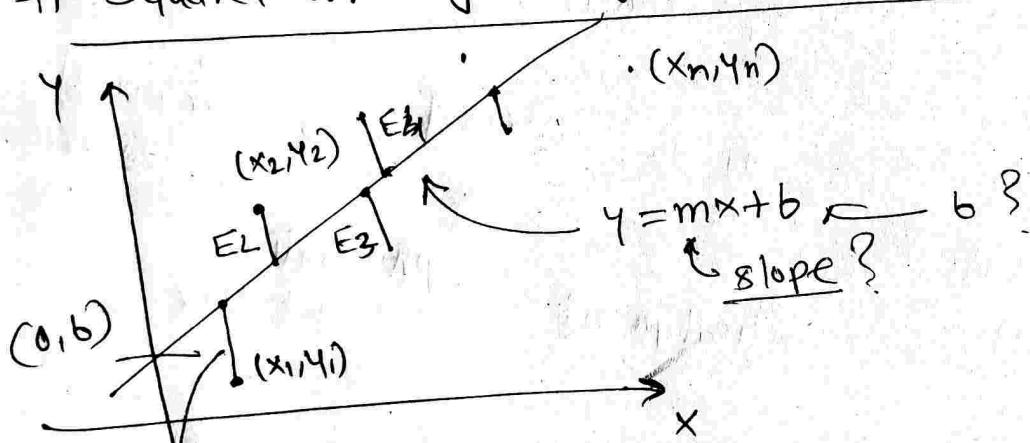


97.5% → $\frac{z}{1.96}$
 z value in
term of
std. deviation

$\Rightarrow 2.35$ is extreme right/left of 1.96
Hence probability is $< 5\%$

\Rightarrow Hence, reject H_0

Squared error of regression line &



$$\text{Error } 1 = y_1 - (mx_1 + b)$$

$$\epsilon_2 = y_2 - (mx_2 + b)$$

$$\epsilon_n = y_n - (mx_n + b)$$

$$SE_{\text{Line}} = (y_1 - (mx_1 + b))^2 + (y_2 - (mx_2 + b))^2 + \dots + (y_n - (mx_n + b))^2$$

find m & b that minimize

$\underline{SE_{\text{Line}}}$

(20)

$$SE_{\text{LINE}} = (y_1 - (mx_1 + b))^2 + \cancel{y_2} (y_2 - (mx_2 + b))^2 + \dots + (y_n - (mx_n + b))^2$$

 SE_{LINE}

$$= y_1^2 + 2y_1(mx_1 + b) + (mx_1 + b)^2 +$$

$$+ y_2^2 - 2y_2(mx_2 + b) + (mx_2 + b)^2$$

+ ~~y₃~~:

$$y_3^2 - 2y_3(mx_3 + b) + (mx_3 + b)^2$$

$$= y_1^2 - 2y_1mx_1 - 2y_1b + m_1x_1^2 + 2mx_1b + b^2$$

+ 2mx₁b + b²

$$+ y_2^2 - 2y_2mx_2 - 2y_2b + m_2x_2^2 + 2mx_2b + b^2$$

$$+ y_3^2 - 2y_3mx_3 - 2y_3b + m_3x_3^2 + 2mx_3b + b^2$$

$$= (y_1^2 + y_2^2 + \dots + y_n^2) - 2m(x_1y_1 + x_2y_2 + \dots + x_ny_n) - 2b(y_1 + y_2 + \dots + y_n)$$

$$+ m^2(x_1^2 + x_2^2 + \dots + x_n^2) + 2mb(x_1 + x_2 + \dots + x_n)$$

$$+ nb^2 \quad \text{--- (1)}$$

$$\frac{y_1^2 + y_2^2 + \dots + y_n^2}{n} = \bar{y}^2$$

$$SE_{\text{LINE}} = n \left(\frac{y_1^2 + y_2^2 + \dots + y_n^2}{n} - \bar{y}^2 \right)$$

--- (a)

$$\frac{x_1y_1 + x_2y_2 + \dots + x_ny_n}{n} = \bar{xy} \Rightarrow x_1 + y_1 + x_2 + y_2 + \dots + x_n + y_n = n\bar{xy}$$

--- (b)

similarly, we have $n\bar{y}$, $n\bar{x}^2$, $n\bar{x}$ --- (c)

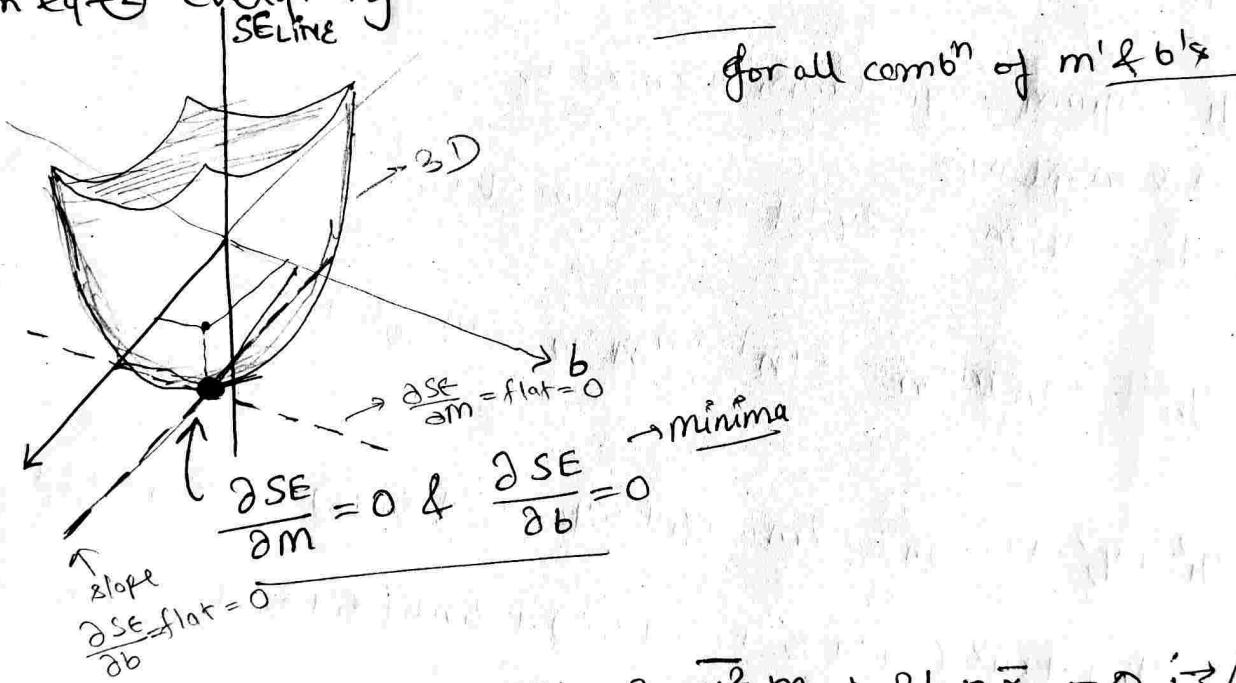
Plug ④⑥⑦ in ①.

$$SE_{LINE} = \overline{ny^2} - 2mn\overline{xy} - 2bn\overline{y} + m^2n\overline{x^2} + 2mbn\overline{x} + nb^2 \quad (2)$$

eq-② breaking into 3-D calculus

eq-② is a surface;

In eq-② everything is a constant except m^2 & b^2 .



$$SE_{LINE} \Rightarrow \frac{\partial SE}{\partial m} = -2n\overline{xy} + 2n\overline{x^2} \underline{m} + 2\underline{b}n\overline{x} = 0; \rightarrow 1/2$$

$$\frac{\partial SE}{\partial b} = -2n\overline{y} + 2\underline{m}\overline{x} + 2n\underline{b} = 0; \rightarrow 1/2$$

$$-\overline{xy} + \underline{m}\overline{x^2} + \underline{b}\overline{x} = 0 \quad \left\{ \begin{array}{l} \text{simplified} \\ (1) \end{array} \right.$$

$$-\overline{y} + \underline{m}\overline{x} + \underline{b} = 0 \quad (2)$$

add \overline{xy} both sides at eq-①② & adding \overline{y} both sides

$$\overline{y} - \overline{y} + \underline{m}\overline{x^2} + \underline{b}\overline{x} = 0 + \overline{y}$$

$$m\overline{x} + b = \overline{y}$$

21

$$m\overline{x^2} + b\overline{x} = \overline{xy}$$

$$m\bar{x} + b = \bar{y} -$$

divide \overline{x}

$$m \frac{\overbrace{x^2}^2}{\overbrace{x}^1} + b = \frac{\overbrace{x^4}^4}{\overbrace{x}^1}$$

$$y = mx + b$$

(\bar{x}, \bar{y}) lies on the
line

$\left(\frac{\bar{x}^2}{\bar{x}}, \frac{\bar{xy}}{\bar{x}} \right)$ on best fitting line

Sol (\bar{x}, \bar{y}) & $\left(\frac{\bar{x}^2}{\bar{x}}, \frac{\bar{xy}}{\bar{x}}\right)$ lie on $y = mx + b$

$$m\bar{x} + b = \bar{y}$$

$$m \frac{x^2}{x} + b = \frac{x^4}{x}$$

$$m\left(\bar{x} - \frac{\bar{x^2}}{\bar{x}}\right) = \bar{y} - \frac{\bar{xy}}{\bar{x}}$$

$$m = \frac{\frac{y_2 - y_1}{x_2 - x_1}}{\frac{x_2 - x_1}{x_2 - x_1}}$$

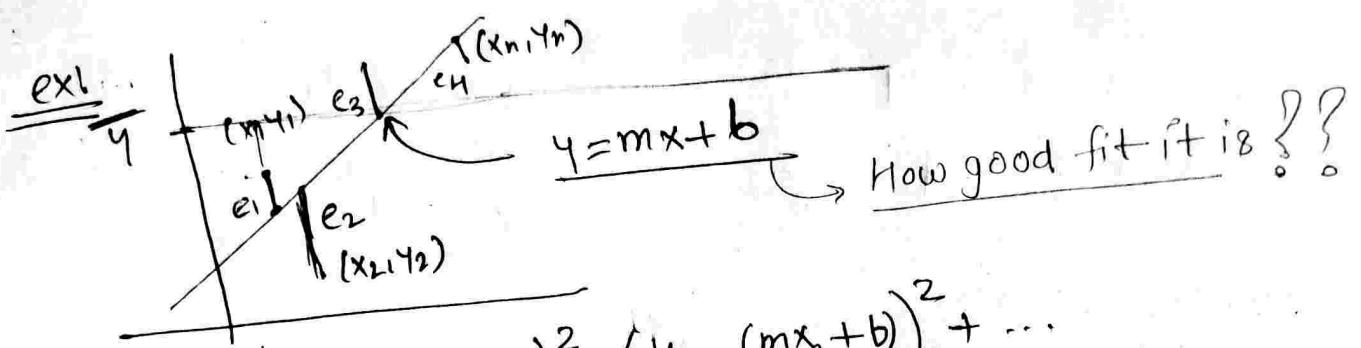
$$\frac{(\bar{x})}{(\bar{x})}$$

$$\frac{\overline{x^4} - \bar{x^4}}{(\bar{x})^2 - \bar{x^2}}$$

$$b = \bar{y} - m\bar{x}$$

$$m = \frac{\overline{x_4} - \overline{x_4}}{(\overline{x})^2 - \overline{x^2}}$$

$$\checkmark = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - (\bar{x})^2}$$



$$SE_{\text{LINE}} = \frac{\sum (y_i - (mx_i + b))^2}{(n-2)} \quad \text{--- (1)}$$

How much (what %) of the total variation in y is described by the variation in x .

$$\text{total variation in } y = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n}} \Rightarrow SE_{\bar{y}}$$

$$\text{variance} = \frac{\text{total variation}}{n}$$

How much of the total variation is not described by the regression line?

$$\frac{SE_{\text{LINE}}}{SE_{\bar{y}}}$$

tells what % of ~~the~~ total variation is not described by variation in x (or by the regression line)

then remainder will be described by the line.

$$1 - \frac{SE_{\text{LINE}}}{SE_{\bar{y}}} \quad \begin{array}{l} \text{blw actual & line} \\ \text{sqared error blw } y \end{array}$$

What % of total percentage variation described by line.

Coefficient of determination

$$r^2 = \boxed{1 - \frac{SE_{\text{LINE}}}{SE_{\bar{y}}}}$$

if squared is really small: \rightarrow line is good fit

$\rightarrow \frac{SE_{\text{LINE}}}{SE_{\bar{y}}}$ will be small

$\rightarrow r^2$ will be close to 1



\rightarrow if large $\rightarrow r^2$ will close to 0

examples for points $(-2, -3), (-1, -1), (1, 2) \& (4, 3)$

sol: best fit line: $\bar{y} = \bar{y} - m\bar{x}$
by formula: $m = \frac{\bar{y} - \bar{y}}{\bar{x}^2 - (\bar{x})^2}$

Step 1

$$\boxed{y = \frac{41}{42}x - \frac{5}{21}}$$

Step 2: sqrd from mean \bar{y}

total squared error from $\bar{y} \Rightarrow$

$$SE_{\bar{y}} = 22.75$$

total sqrd error \Rightarrow

$$SE_{\text{LINE}} = 2.74$$

Step 3

$$\frac{2.74}{22.75}$$

% of total variation

not explained by variation in x.

$$\Rightarrow r^2 = 1 - 0.12 = 0.88$$

88% of total variation is explained by variance in X.

Covariance the regression line

$$\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

$$= E[XY - XE[Y] - E[X]Y + E[X]E[Y]]$$

$$= E[XY] - E[X]E[Y] - E[XE[Y]] - E[E[X]Y]$$

$$\text{cov}(X, Y) = E[XY] - E[Y]E[X] + E[E[X]E[Y]] - E[X]E[Y] + E[X]E[Y]$$

$$\text{cov}(X, Y) = E[XY] - \underbrace{E[Y]E[X]}_{\frac{Y}{X}} \quad \checkmark$$

$$E[XY] \approx \bar{XY}$$

$$\text{cov}(X, Y) = \bar{XY} - \bar{Y}\bar{X} \quad \leftarrow \text{numerator of slope of line}$$

$$\hat{m} = \frac{\bar{XY} - \bar{Y}\bar{X}}{\bar{X}^2 - (\bar{X})^2} \quad \checkmark$$

$$\bar{X} \cdot \bar{X} - \bar{X} \cdot \bar{X}$$

$$\text{cov}(X, X)$$

$$\hat{m} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$$

covariance means:- how changes in X are associated to Y

Chi-square distributions

$X \rightarrow$ random variable; normally distributed

$$X \sim N(0, 1) \quad E[X] = 0; \quad \underline{\text{Var}(X) = 1.}$$

mean variance

$$Q_1 = X^2$$

the distribution for the random variable Q_1 will be
an example of chi-square distribution

$$Q_1 \sim \mathcal{X}^2_1 \leftarrow \text{degree of freedom}$$

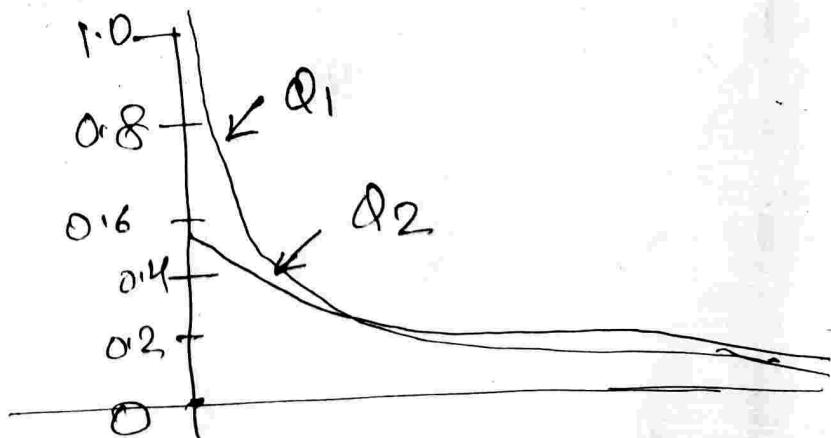
$$\cancel{Q_2 = X_1^2 + X_2^2}$$

chi-square distribution.

$$Q_2 \sim \mathcal{X}^2_2 \leftarrow \text{degree of freedom}$$

$$Q_3 = X_1^2 + X_2^2 + X_3^2$$

$$\cancel{Q_3 \sim \mathcal{X}^2_3}$$



What is the probability of $Q_2 > 2.41$?

$P(Q_2 > 2.41) \rightarrow$ look χ^2 -square table
against d (degree of freedom)

$$= \frac{0.30}{30\%}$$
