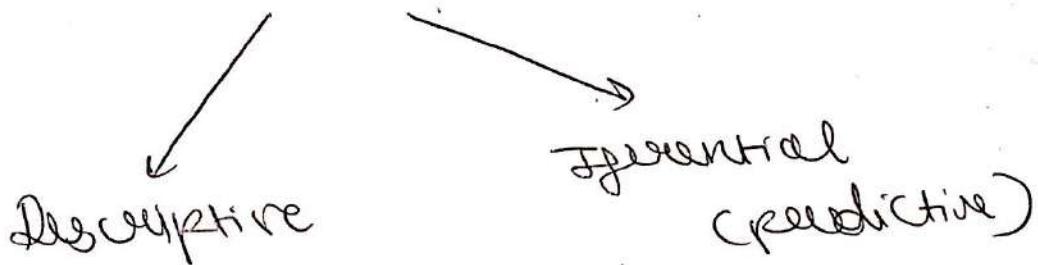


Statistics



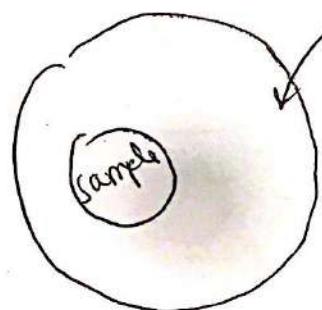
Descriptive statistics

Central tendencies of Data Set

- ① mean
- ② median
- ③ mode

} you should
know how to
calculate

Sample vs population mean :-



suppose I want to calculate
average height of people
of america, then, don't
take complete population

instead takes the mean
of sample.

$$\mu = \text{population mean} = \frac{\sum_{n=1}^N x_n}{N}$$

$$\bar{x} = \text{sample mean} = \frac{\sum_{n=1}^n x_n}{n}$$

$$\bar{M} = \frac{\sum_{i=1}^N x_i}{N} = \frac{x_1 + x_2 + \dots + x_N}{N} = \text{Population mean}$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n} = \text{Sample mean}$$

→ If the distribution is skewed then median is the best center of measure.

Dispersion :- (Variance)

2, 3, 3, 3

$$\bar{M} = \frac{2+2+3+3}{4} = 2.5$$

0, 0, 5, 5

$$\bar{M} = \frac{0+0+5+5}{4} = 2.5$$

They are more dispersed because numbers are farther away from mean

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

(Variance)

0, 0, 5, 5

#	2, 1, 3, 3			
i	x _i	M	x _i - M	(x _i - M) 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	5	2.5	-0.5	0.25

$$\sigma^2 = \frac{(2.5)^2 + (2.5)^2 + (2.5)^2 + (2.5)^2}{4}$$

$$\sigma^2 = \frac{6.25 + 6.25 + 6.25 + 6.25}{4}$$

$$\sigma^2 = 6.25$$

$$\sigma^2 = \frac{0.25 + 0.25 + 0.25 + 0.25}{4} = \frac{1}{4} = 0.25$$

We have two data sets, mean of both data sets are equal but variance of data sets differ.

$$\sigma^2 = 0.25$$

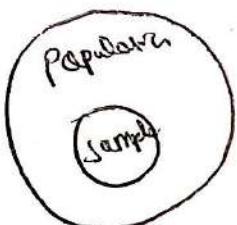
$$\sigma^2 = 6.25$$

→ Smaller variance means numbers in the data sets are much farther away than the mean than the numbers in the second data set which have less variance.

Sample variance :-

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

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



sample size

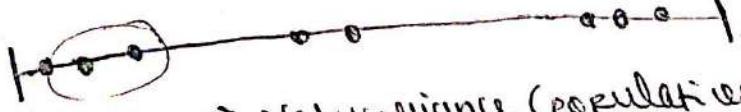
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

sample variance

We choose $\sqrt{s^2}$ instead of s

→ sometimes it may be possible that data set which is very skewed / distorted

$\bar{x}_1 < \bar{x}_2 < \bar{x}_3$



→ Net variance (population) lies in between data set ... than the actual

population variance. Population variance
 जब वर्ग से होता है। लेकिन यदि समूह
 एवं विचलन विवरण तो उनकी विवरण
 विवरण के साथ विवरण के विवरण
 विवरण के साथ विवरण के विवरण
 के विवरण तो उनकी विवरण के विवरण
 के विवरण के विवरण के विवरण

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad \begin{array}{l} \text{Sample variance} \\ \downarrow \end{array} \quad \rightarrow \text{unbiased estimator}$$

Standard deviation :-

$$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N}}$$

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

Alternate variance formulas :-

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

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

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

$$= \frac{\sum_{i=1}^N x_i^2}{N} - \frac{2\bar{x}(\bar{x})}{N} + \frac{\bar{x}^2 N}{N}$$

$$= \frac{\sum_{i=1}^N x_i^2}{N} - 2\bar{x}^2 + \bar{x}^2$$

Nominal variable :- categorical variable.

$$\sigma^2 = \frac{\sum_{i=1}^N x_i^2}{N} - \bar{x}^2$$

ordinal variable :-
They are also categorical variable.

variance.

$$\text{Interquartile range} = (Q_3 - Q_1)$$

Random Variables

(X)

NOTE:- Q_2 is median

→ we can solve the random variable

$$X = \begin{cases} 1, & \text{if rains} \\ 0, & \text{no rains} \end{cases}$$

→ Random variable is map of the event which maps us from random process to number.

Real value

number.

$$X = \begin{cases} 1, & \text{if heads} \\ 0, & \text{if tails} \end{cases}$$

→ Random variable, अतः वारिएबल
of value जिसके साथ हो सकती है

Random Variable

Discrete

- (i) roll a die
- (ii) toss a coin

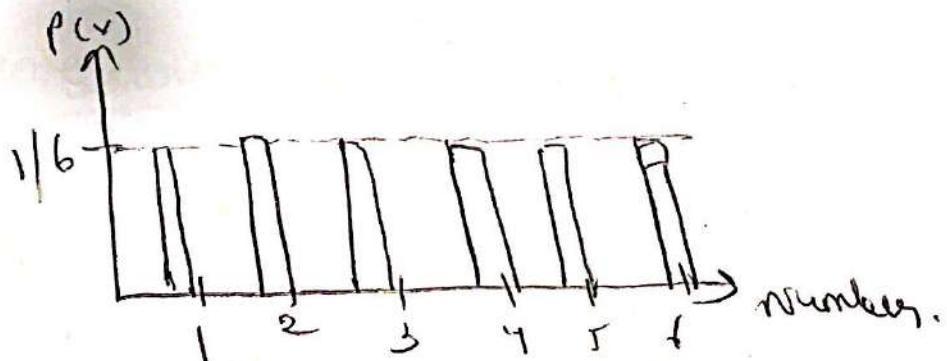
↓
we have countable
number of outcomes

continuous

↓
infinite number of
outcomes

↓
 $X = \text{exact amount of}\text{ rain in inches}$

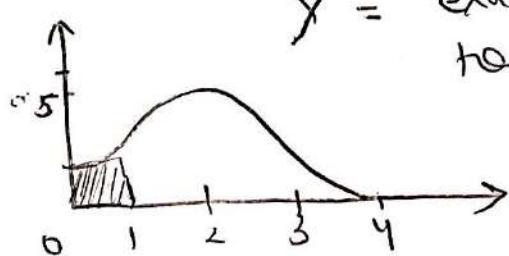
Probability distribution of $X = \# \text{ number facing up on the fair dice:-}$



$$X = \begin{cases} 1, & \text{if head} \\ 0, & \text{if tail} \end{cases}$$



Probability density functions



x = exact amount of rain tomorrow

$$P(X=2) = 0.5$$

$$P(|Y-2| < 1)$$

$$P(1.9 < Y < 2.1)$$

Q what is the probability that rain will occur less than 1 inches?

Ans:

$$\int_0^1 f(x) dx \quad \left\{ \begin{array}{l} \text{we integrate} \\ \text{the function} \\ \text{from 0 to 1} \end{array} \right.$$

→ Area under the curve (the complete area) is equal to 1.

Binomial distribution :-

Q coin flip 5 times.

X = No. of heads after 5 flips
probability of getting no head

$$P(X=0) = \left(\frac{1}{2}\right)^5 = \frac{1}{32}$$

$$P(X=1) = \frac{\text{Probability of exactly one head}}{\text{Total outcomes}} = \frac{5}{32}$$

XXXXX

$$P(X=2) = \frac{5C_2}{32} = \frac{5!}{3! \times 2!} \times \frac{1}{32}$$

~~$nC_r p^r q^{n-r}$~~

$$= \frac{5 \times 4 \times 3!}{2! \times 1 \times 3!} \times \frac{1}{32} = \frac{10}{32} = \frac{5}{16}$$

$$P(X=3) = \frac{5C_3}{32} = \frac{10}{32}$$

Probability exactly three heads

$$P(X=4) = \frac{5C_4}{32} = \frac{5!}{1! \times 4!} \times \frac{1}{32} = \frac{5}{32}$$

Probability exactly four heads.

$$P(X=5) = \frac{1}{32}$$

Probability getting all heads.
(5 heads).

Now,

$$P(X=0) = \frac{1}{32}$$

$$P(X=1) = \frac{5}{32}$$

$$P(X=2) = \frac{10}{32} = \frac{5}{16}$$

$$P(X=3) = \frac{10}{32} = \frac{5}{16}$$

$$P(X=4) = \frac{5}{32}$$

$$P(X=5) = \frac{1}{32}$$



$n=6$ (Basketball)

P of success = 30%

X = number of shots made

P of failure = 70%

$$P(X=0) = 0.7 \times 0.7 \times 0.7 \times 0.7 \times 0.7 \times 0.7 \\ = (0.7)^6$$

probability that i made no shots

$$P(X=1) = 0.3 \times (0.7)^5 \times 6 \\ \hookrightarrow P \text{ that i made 1 shot}$$

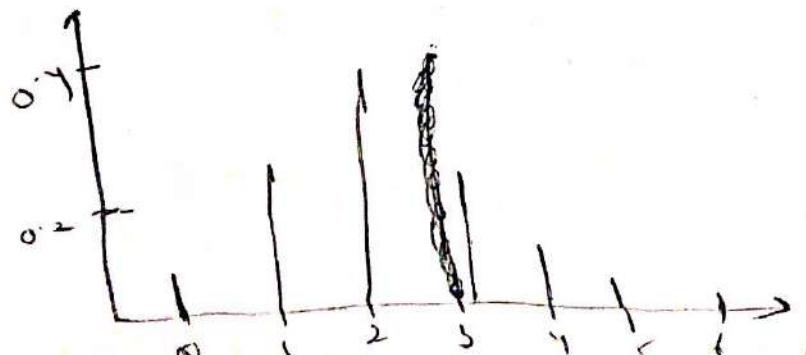
$$P(X=2) = 6 \times (0.3)^2 \times (0.7)^4$$

$$P(X=3) = 6 \times (0.3)^3 \times (0.7)^3$$

$$P(X=4) = 6 \times (0.3)^4 \times (0.7)^2$$

$$P(X=5) = 6 \times (0.3)^5 \times (0.7)$$

$$P(X=6) = 6 \times (0.3)^6 \times (0.7)^0 \\ = (0.3)^6$$



Expected value of a Random Variable :-

3 3 3 4 5

$$\frac{3+3+3+4+5}{5} = \frac{18}{5} = 3.6.$$

$$\frac{3(3)+1(4)+1(5)}{5} = \frac{1}{5} (3 \cdot 3 + 1 \cdot 4 + 1 \cdot 5)$$

$$= \frac{3}{5} \cdot 3 + \frac{1}{5} \cdot 4 + 1$$

$$= 0.6 \times 3 + 0.2 \times 4 + 1$$

$$= 60\% \text{ of } 3 + 20\% \text{ of } 4 + 20\% \text{ of } 5$$

60% of numbers are 3

20% of numbers are 4

20% of numbers are 5

→ If we have the infinite population
then we can find out the arithmetic mean
of the infinite population by this technique

$$E(x) = (60\% \text{ of } 3 + 20\% \text{ of } 4 + 20\% \text{ of } 5)$$

→ expected value is same as
mean but often it's different

Probability of k successes in a binomial distribution is :-

$$P(X=k) = {}^n C_k p^k q^{n-k}$$

Now,

$E(X)$ (expected value) of # of successes with probability p after n trials

$$E(X) = n \cdot p$$

For eg:- No. of baskets I make after 10

shots. Probability of success is 40%?

Ans: $E(X) = 0.4 \times 10 = 4$

\uparrow
probability

Hence, I can make 4 baskets

throwing it 10 times.

$$\begin{aligned} E(X) &= \sum_{k=0}^n k {}^n C_k p^k (1-p)^{n-k} \\ &= np \sum_{k=1}^n \frac{(n-1)!}{(k-1)! (n-k)!} \times p^{(k-1)} \times (1-p)^{n-k} \end{aligned}$$

$$= n p \sum_{a=0}^b \binom{b}{a} p^a (1-p)^{b-a}$$

Since this will sum up to 1

$$E(X) = np \rightarrow \text{expected value of binomial distribution}$$

→ Expected value is n times p where n is number of trials and p is probability of success.

→ This is only true for random variable X whose probability distribution is the binomial distribution.

Poisson Process :-

X = No. of calls pass in an hour

$$E(X) = \lambda = r \cdot p \quad \begin{matrix} \text{success in} \\ \text{small interval} \end{matrix}$$

$$\lambda \text{ calls/hour} = 60 \text{ min/hour} \times \frac{\lambda \text{ calls}}{60 \text{ min}}$$

We assume
1 call passes
in 1 min

$$P(X=k) = 60C_k \times \left(\frac{\lambda}{60}\right)^k \times \left(1 - \frac{\lambda}{60}\right)^{60-k}$$

where we
assume that
1 call passes
in 1 sec.

$$P(X=k) = 3600C_k \times \left(\frac{\lambda}{3600}\right)^k \times \left(1 - \frac{\lambda}{3600}\right)^{3600-k}$$

→ We make it more accurate.

→ We make more granular.

So,

$$\lim_{K \rightarrow \infty} \left(1 + \frac{\lambda}{n}\right)^{nK} = e^\lambda$$

So, Now, we put limit to n as to assume infinite calls passes in an hour. (n → ∞)

$$P(X=K) = \lim_{n \rightarrow \infty} \binom{n}{K} \left(\frac{\lambda}{n}\right)^K \left(1 - \frac{\lambda}{n}\right)^{n-K}$$

Solving this up we
get

$$P(X=K) = \lim_{n \rightarrow \infty} \binom{n}{K} \left(\frac{\lambda}{n}\right)^K \left(1 - \frac{\lambda}{n}\right)^{n-K}$$

λ is mean.
 K is required success.
 $X = np$

I assume, that 9 calls passed in an hour? Find probability that exactly two calls pass in an hour?

$$P(X=2) = \frac{(9)^2}{2!} e^{-9}$$

$$= \frac{e^{-m} m^r}{r!}$$

m = mean
 r = required success

NOTE :- i) यह poison का use करें और binomial
and use करो, परन्तु यदि n की value बहुत
large हो जाए तो ef:- $n = 100$, इसके
लिए को 100, 200 वाले 3 cell (1.101), परन्तु अगर
then we use poison's in that case.

$$P(x \text{ success}) = \frac{e^{-m} \times m^x}{x!}$$

n = no of success trials

m = mean

x = required success.

$$m = np$$

ii) यह experiment repeat upto N times

Frequency of success = $N p(x)$

Expected value

$$E(x) = N e^{-m} \frac{m^x}{x!}$$

Q) Find the probability that at most 5 defective
jars will be found in a box of 200 jars
if experience shows that 2% of
such jars are defective?

Ans:

$$P = \frac{2}{100} = 0.02$$

$$q = 1 - 0.02 = 0.98$$

$$n = 200$$

$$\text{mean} = np = 0.02 \times 200 \\ = 4.0$$

$$x = 0, 1, 2, 3, 4, 5$$

Poisson distribution $p(x) = \frac{e^{-\lambda} \lambda^x}{x!}$

$$p(0, 1, 2, 3, 4, 5) = p(0) + p(1) + p(2) + p(3) \\ + p(4) + p(5) \\ = \frac{e^{-4} (4)^0}{1!} + \dots$$

If probability that an individual suffer a bad reaction from a certain injection is 0.001, determine the probability that out of 2000 individuals

- i) exactly 3
- ii) more than 2 individuals
- iii) None
- iv) more than 1 individual will suffer a bad reaction?

Ans: $P = 0.001, n = 2000$

$$\lambda \text{ or } m = 0.001 \times 2000$$

$$= 2$$

i) $x = 3$ $p(3) = p(3) = \frac{e^{-2} \times (2)^3}{3!}$

$$\begin{aligned}
 \text{i)} \quad P(X=3, 4, 5, \dots, 2000) &= P(3) + P(4) + \dots + P(2000) \\
 &= 1 - [P(0) + P(1) + P(2)]
 \end{aligned}$$

ii) $X=0$

$$\text{iv)} \quad P(X=2, 3, 4, \dots, 2000) = 1 - [P(0) + P(1)]$$

NOTE :- persons की नहीं वाले अंक
अस्ति न हो, वे वाले अंक जहाँ
होते हैं, बिंदुओं की विवरण होते हैं।

Law of Large Numbers :-

$$\bar{X}_n = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\bar{X}_n \rightarrow \mu$$

It says that if our sample size approaches $n \rightarrow \infty$, it approaches population mean if $n \rightarrow \infty$.

X = no. of heads after 100 tosses of fair coin

$$\begin{aligned}
 \text{Now, } E(X) &= n \times p \\
 &= 100 \times \frac{1}{2} = 50
 \end{aligned}$$

Suppose

I shake the box one time I get 55 ready, then I shake it again I get 65, & get 45

$$\bar{x}_n = \frac{55 + 65 + 45 + \dots + n}{n}$$

$$\bar{x}_n \rightarrow 50 \text{ as } n \rightarrow \infty$$

law of large numbers

Normal Distribution

(Gaussian distribution)

(Bell curve)

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

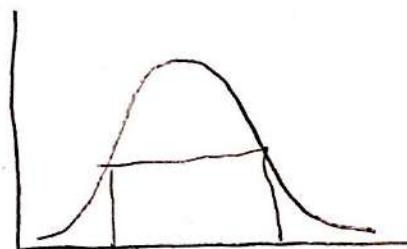
standard Z-score.

σ^2 = variance

σ = standard deviation

μ = mean.

x = value.



Note:- binomial / poisson \Rightarrow discrete random variable but Normal distribution our random variable is continuous.

$\Omega \rightarrow X \rightarrow$ lies between interval.

(here $x = \omega$ (it is continuous))

→ Random Variable here lies in a interval and it is continuous.

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

($-\infty < x < \infty$)

→ random variable infinite

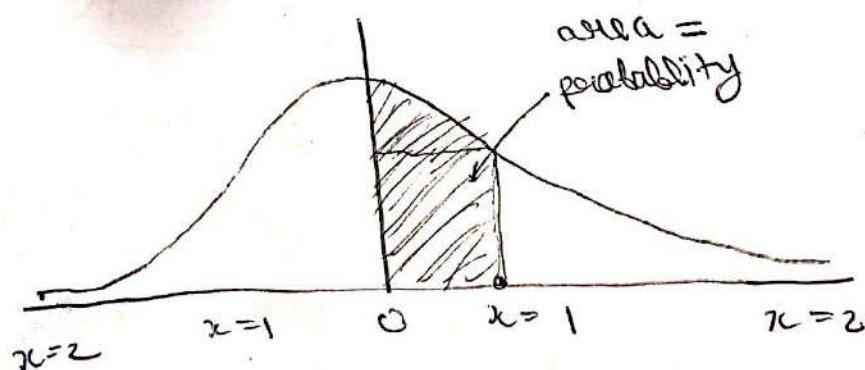
Standard Normal distribution :-

mean $\mu = 0$

standard deviation $\sigma = 1$

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-x^2/2}$$

Standard
Normal
distribution



→ Standard normal distributions
is also called Z -distribution.

~~Step 1:~~ Convert x into standard normal values
by formula : $Z = \frac{x - \mu}{\sigma}$ —①

Step 2: Find the limits of Z corresponding
to the limits of x
when $x=a$ then $Z = \frac{a - \mu}{\sigma}$
put $\{x=a$
int ① $\}$

when $x=b$ then $Z = \frac{b - \mu}{\sigma}$

thus, the limits of Z are $\frac{a - \mu}{\sigma}$ to $\frac{b - \mu}{\sigma}$
when $x=a, x=b$

Step 3: $P(a < x < b) = P\left(\frac{a - \mu}{\sigma} < Z < \frac{b - \mu}{\sigma}\right)$
use Normal table.

Q The mean height of 500 students is 151 cm and the standard deviation is 15 cm. Assuming that the heights are normally distributed, find how many students have height between 120 and 155 cm.

Ans: → Height is a continuous variable

Given mean

$$\mu = 151$$

$$\text{Standard deviation} = \sigma = 15 \text{ cm}$$

Now,

$$x = 120 \Rightarrow z = \frac{120 - 151}{15} = \frac{-31}{15} = -2.0666 \approx -2.07$$

$$\text{when } x = 155 \Rightarrow z = \frac{155 - 151}{15} = \frac{4}{15} = 0.2666 \approx 0.27$$

$$P(120 < X < 155) = P(-2.07 < z < 0.27)$$

probability that ^{when X} lies between 120 to 155 is equal to $-2.07 < z < 0.27$

$$= P(-2.07 < z < 0) + P(0 < z < 0.27)$$

$$= 0.4808 + 0.1064 \quad \left. \right\} z \text{ value} \\ = 0.5872 \quad \text{Are we have standard } z \text{ table}$$

No. of students whose weight lies b/w 120 cm to 155 cm = $N \times P(120 < X < 155)$

$$= 500 \times 0.5872$$

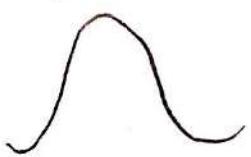
$$= 284 \text{ (approx)}$$

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

very imp.

Note :-

$$\textcircled{1} \quad \text{If } \sigma = 5$$



$$\sigma = 10$$

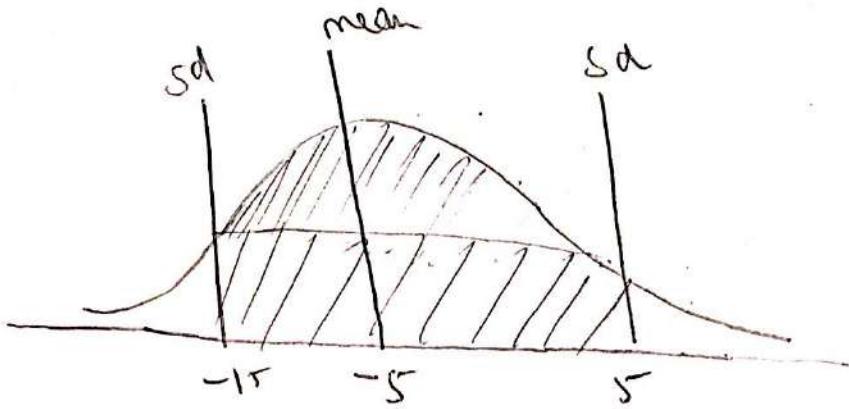


Cumulative density function :-

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

what is probability of getting -1 or less than $x = -1$





what is probability of standard deviation
of the mean? assuming normal distribution?

$$\int_{-15}^5 p(x) dx \rightarrow 68.3\%$$

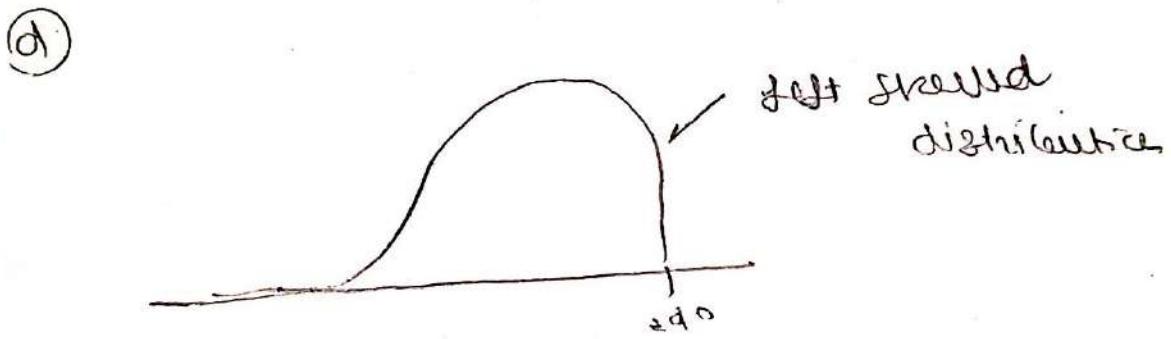
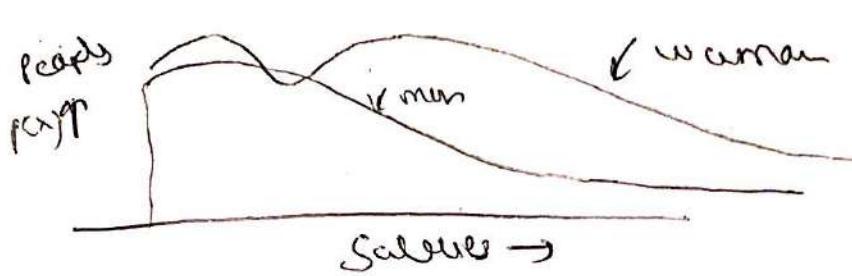
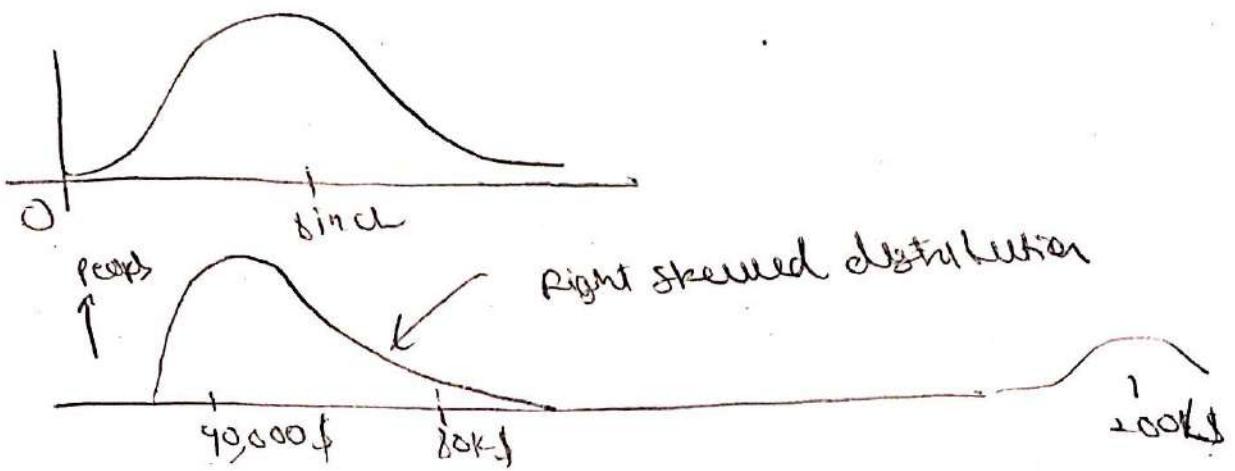
Q which of the following data sets is most likely to be normally distributed? For other choices, explain why they are not normally distributed?

① hand span (measured from tip of thumb to tip of the extended 5th finger) of a random sample of high school seniors.

② annual salary of all employees of large shipping company

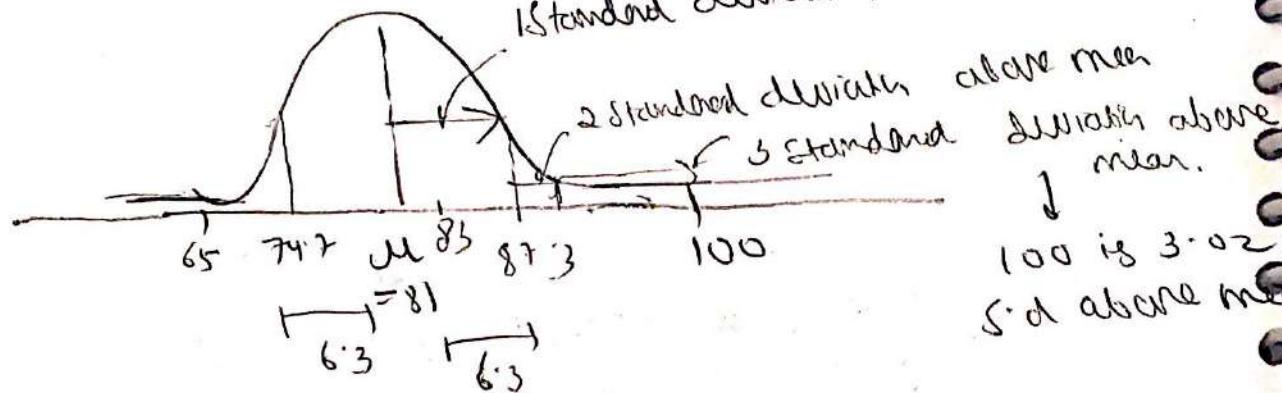
③ annual salaries of random sample of 50 CEO's of major companies, 25 women & 25 men.

④ date of 100 pennies taken from a cash drawer to a convenience store?



Grades on a statistics mid-term for a high school are normally distributed with $\mu = 81$ and $\sigma = 6.3$. Calculate z -scores for each of the following exam grades, draw & label sketch for each example.

~~100~~ z -score is how many standard deviations away from the mean (μ).



Q) 65 80, $\frac{65-81}{6.3} = \frac{-16}{6.3} = -2.54$

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

→ If means if 65 is (-2.54) ~~below~~
standard deviations away from
mean.

Q) 83, $Z = \frac{83-81}{6.3} = \frac{2}{6.3} = \frac{2}{6.3} = 0.32$

→ 83 is 0.32 standard deviation
away from mean.

→ We can say 83 is two grade
above the mean test mark in terms of stand-
ard deviation we have find out Z score, now
I can say 83 is 0.32 sd away from mean.

Q) ✓

Q) 100 $\frac{100-81}{6.3} = \frac{19}{6.3} = 3.02$

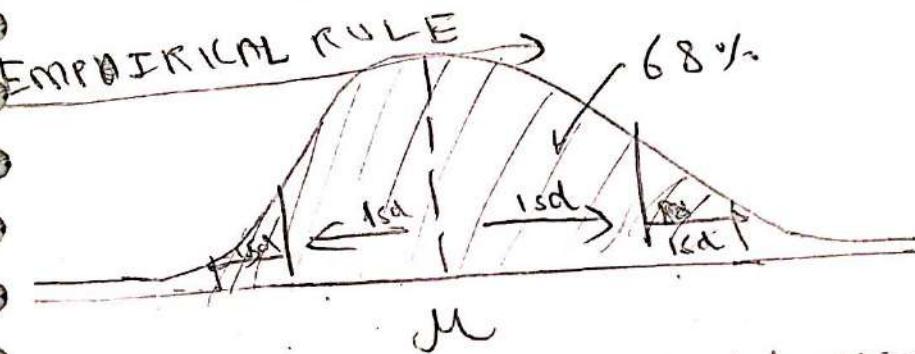
If means 100 is 3.02 standard deviation
above the mean.

Assume that mean weight of 1 year old girl in US is normally distributed with a mean of about 9.5 grams with standard deviation of approx 1.1 grams. Without using calculator, estimate the percentage of 1 year old girl in the US that meets the following conditions?

- ① less than 8.4 kg
- ② between 7.3 kg & 11.7 kg
- ③ more than 12.8 kg.

$$\sigma = 1.1 \text{ gram}$$

$$\mu = 9.5 \text{ gram}$$

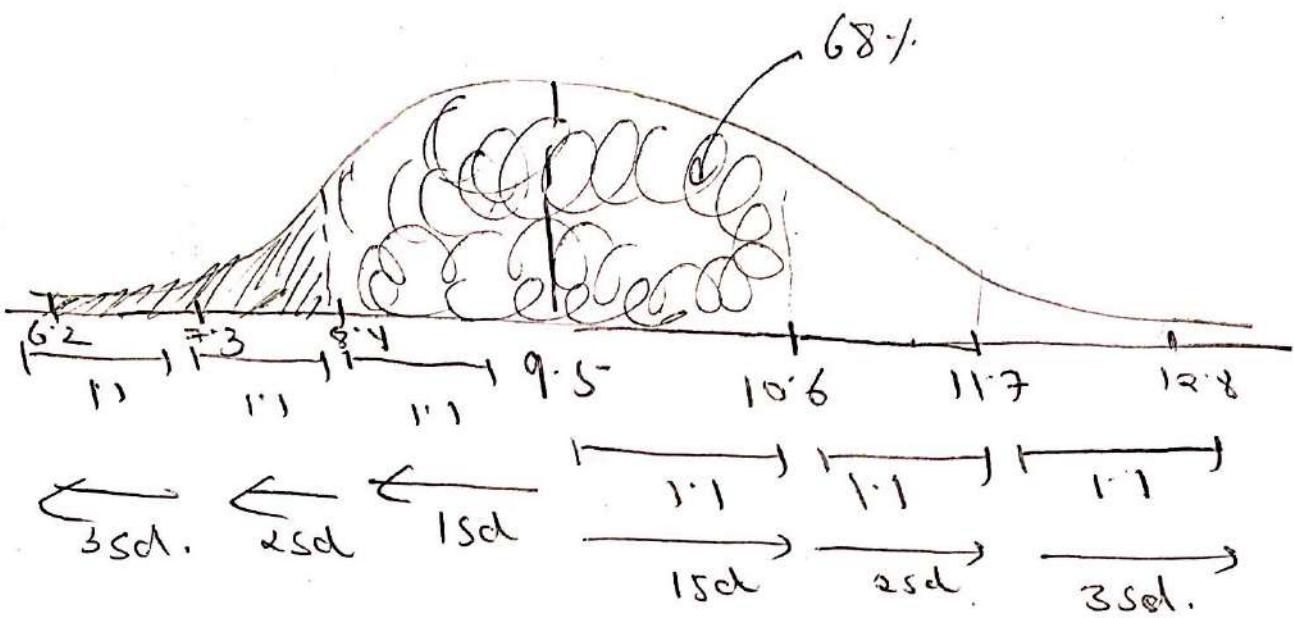


EMPIRICAL RULE

→ one standard deviation below mean & one standard deviation above mean then probability is 68%.

→ 2 sd below mean & 2 sd above mean then probability is 95%.

→ 3 sd below mean & 3 sd above mean then probability is 99.7%.



(a) below or less than 8.4

hence, area below or less than 8.4 is

$$\frac{100 - 68\%}{2} \\ = 16\%$$

(b) between 7.3 and 11.7
it is 95%.

(c) more than 12.8

$$= \frac{100 - 99.997}{2} \\ = \frac{3}{2} = 15\%$$

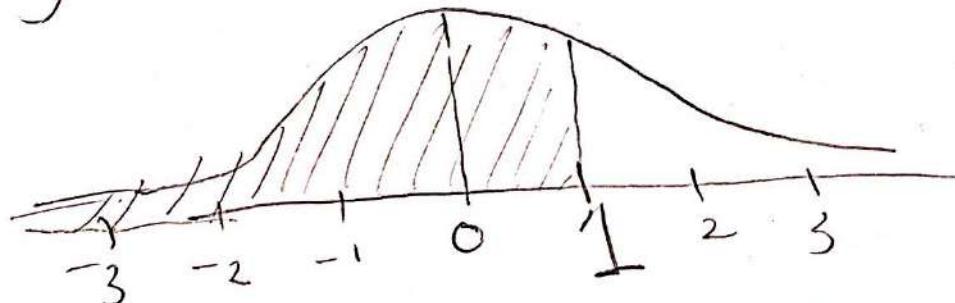
area between -3sd and +3sd = 99.7%

For standard Normal distribution, place the following in order from smallest to largest :-

- (a) percentage of data below 1
- (b) percentage of data below -1
- (c) mean
- (d) standard deviation.
- (e) percentage of data above 2.

Ans: For standard Normal distribution $\mu = 0$ and $\sigma = 1$

(a)



empirical rule $\rightarrow (68 - 95 - 99.7)$

$$50\% + \frac{68\%}{2}$$

$$50\% + 34\% = 84\% = .84$$

(b) Below (-1) area

$$\frac{100 - 68\%}{2} = 16\% = .16$$

(c) $\mu = 0$

(d) $\sigma = 1$

∴ ~~standard normal~~

Below -2 is

$$50\% + \frac{95\%}{2}$$

Above 2 is

$$\begin{aligned} & 100 - \left[50\% + \frac{95\%}{2} \right] \\ & = 100 - 50 - \cancel{\frac{95}{2}} \quad 47.5 \\ & = 2.5\% = .025 \end{aligned}$$

Now, we can arrange them,

$$c < b < a < d$$

- Q In the 2007 AP statistics examination scores were not normally distributed, with $\mu = 2.80$ and $\sigma = 1.34$. What is approximate z-score that corresponds to an exam score of 5 (Score range from 1-5)?

Ans:-

$$z\text{ score} = \frac{5 - 2.8}{1.34} = \frac{2.2}{1.34} = 1.64$$

(How many sd
you are away
from mean.)

Z score can be calculated for normal distribution as well as ~~non~~ not normal distribution too.

INFERENTIAL STATISTICS

Central Limit Theorem :-

→ As your sample size $\rightarrow \infty$ then we will get perfect normal distribution curve.

of sample mean of population.

→ Suppose ~~mean~~, sampling distribution of the sample mean is equal to the population mean.

~~population~~

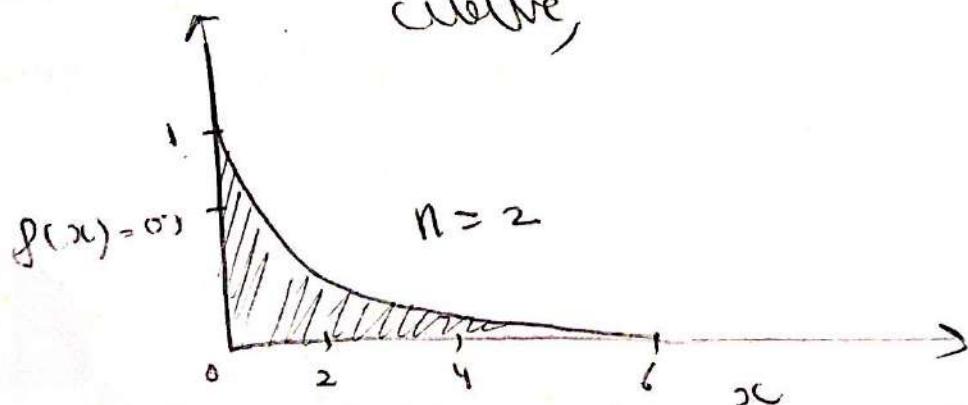
→ If population is normally distributed then our sample mean (\bar{x}) is also normally distributed.

→ But if population is not normally distributed then?

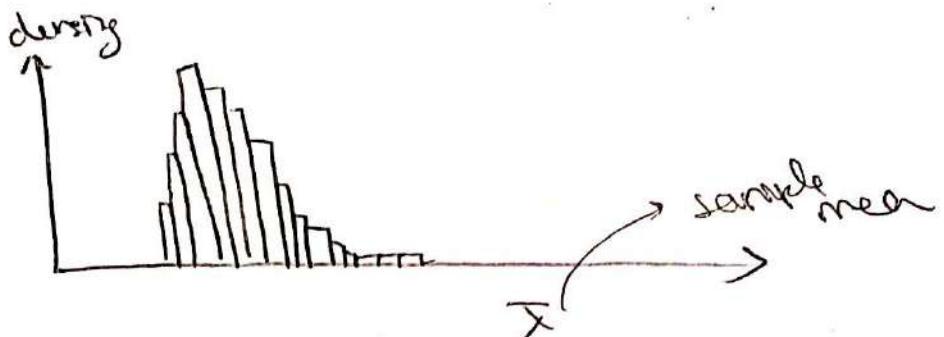
Central limit theorem

↙ The distribution of the sample mean tends towards the normal distribution as sample size increases, regardless the distribution from which we are sampling

Let suppose, we have not normally distributed curve,

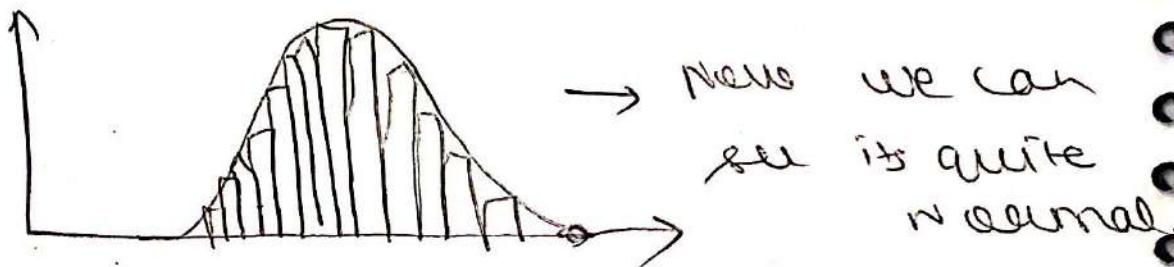


→ Let take two observations $n=2$ and take a mean out of it. Take 1000 of such observation and take mean.



If $n=2$ then sampling distribution of the sample mean is not normal.

→ Now what if I increase sample size as $n=4$ & plotted millions of sample means in histogram.



→ If $n \rightarrow \infty$ → ~~the~~ sampling distribution of sample mean tends towards Normal distribution as sample size increases.

→ Our sample size should atleast be 30 $n \geq 30$ for ~~an~~ normal distribution of sample.

NOTE:-

$$\bar{x} \text{ } n = 4$$



were for sample mean
were SD is more
(σ)

$$\bar{x} \text{ } n = 20$$



were for sample mean
, were SD is less
(σ)

Q Why is this important?

Ans: many statistics have distributions that are approximately normal for large sample sizes, even when we are sampling from a distribution that is not "normal".

→ we can often use well developed statistical inference procedures that are based on normal distribution, even if we are sampling from a population that is not normal, provided we have large sample size.

~~Mean~~

~~Standard Deviation~~

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

as $n \rightarrow \infty$

Suppose salaries at a very large corporation have mean of \$ 62000 and $\sigma = \$32,000$? If a single employee is randomly selected, what is the probability salary exceeds \$ 66000?

$$\text{Ans: } \mu = 62000$$

$$\sigma = 32000$$

$$P(X > 66000)$$

Now

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

$$= P\left(Z > \frac{66000 - 62000}{32000}\right)$$

$$= P(Z > 1.25)$$

→ This question cannot be answered accurately with the given information because salaries are not normally distributed, hence random variable X will not have the normal distribution.

Now, suppose we change the question,

If 100 employees are randomly selected, what is the probability their average salary exceeds \$66000?

$$\text{Ans: } \mu = 62000$$

$$\sigma = 32000$$

~~Normal distribution~~

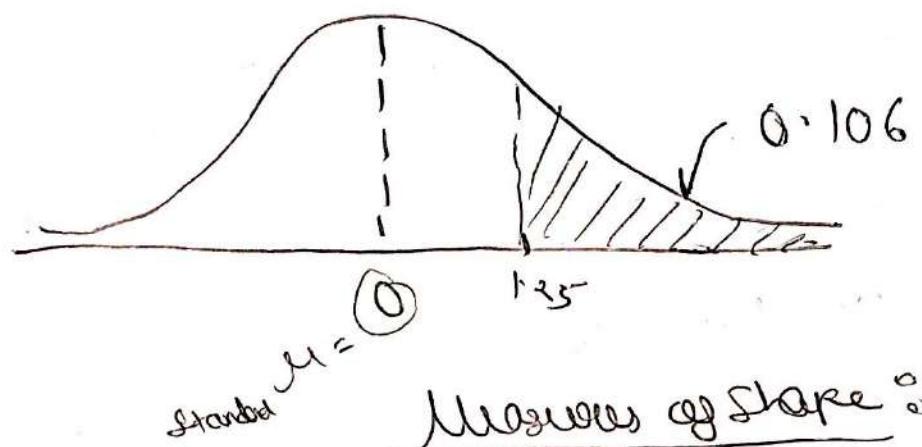
say that \bar{X} is Normally distributed by central limit theorem.

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

Now, we can

$$P\left(Z > \frac{66000 - 62000}{32000 / \sqrt{100}}\right) = P(Z > 1.25)$$

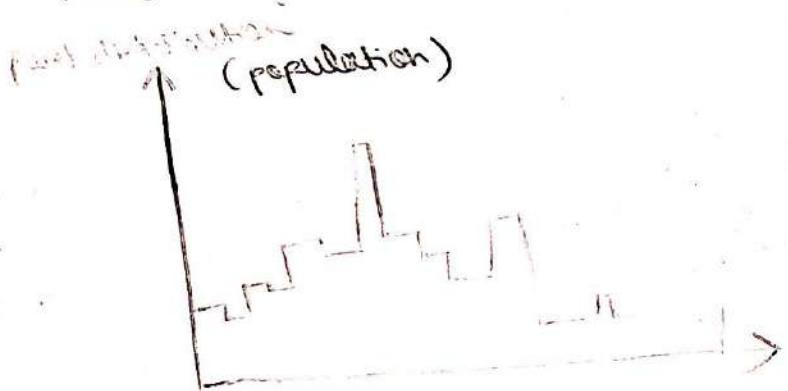
Now $Z=1.25$, Now curve is Normally distributed



Measures of Shape :-

Kurtosis \rightarrow
Skewness \rightarrow

Sampling distribution of a Sample mean :-



$$\text{mean} = 14.45$$

$$\text{median} = 11.00$$

$$sd = 9.78$$

$$\text{kurtosis} =$$

$N = 4$ \rightarrow $4^{\sqrt{2}}$ numbers
we have taken 10,000 data then 30 \sqrt{n} mean
 $n = 10000$ (10,000)
(we have taken 10,000 observations)



$$\text{mean} = 14.44$$

$$sd = 4.40$$

$$\text{kurtosis} = -0.26$$

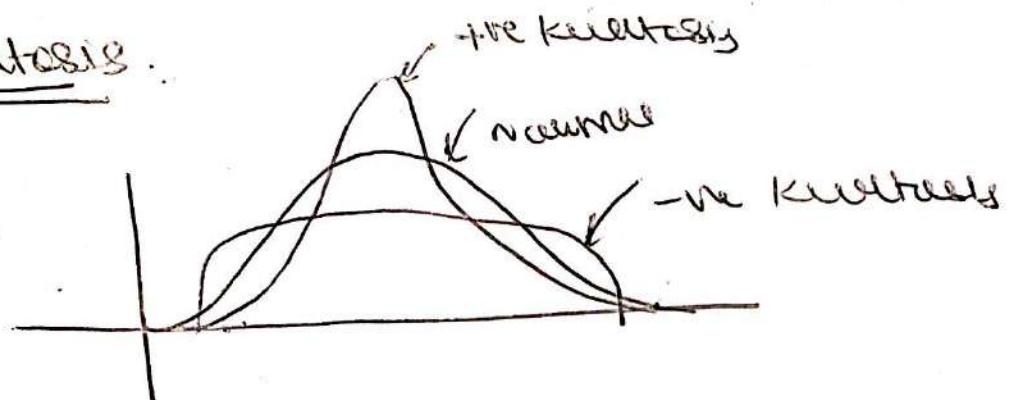
$$\text{if } N=5$$

sample size

$$= -6 \quad (4 \rightarrow n=25) \quad n=25$$

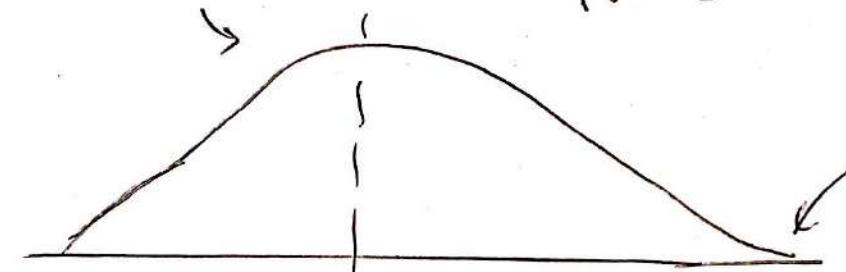
~~Note~~ :- population or mean or samples mean when $n \rightarrow \infty$ is almost same.
→ n.o. of observations

Kurtosis

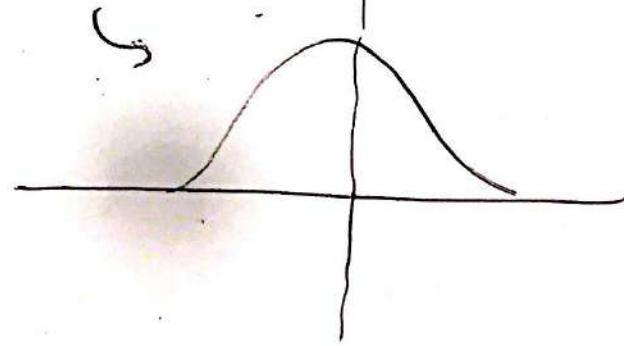


When

$$\text{kurtosis} = -0.26$$



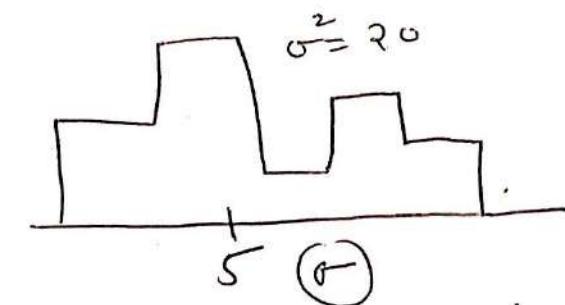
$$\text{kurtosis} = -0.04$$



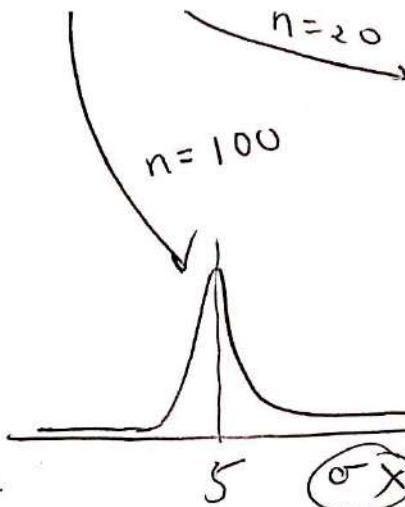
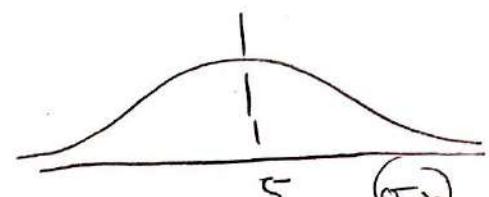
→ When N is larger (No. of sampling distribution) then mean is same in both the cases but ~~SD of each~~ SD is small for those whose N is larger.

$$\text{variance} \propto \frac{1}{N} \rightarrow \text{inversely proportional}$$

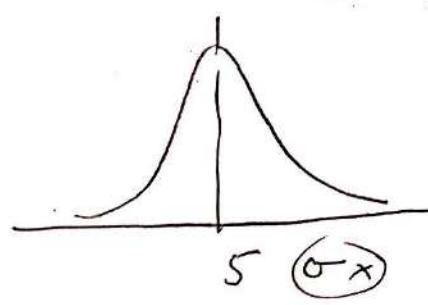
Standard error of mean :-



$$\sigma = 1.0$$



$$n = 100$$



$$\sigma_y^2 = \frac{\sigma^2}{n}$$
$$\sigma_x = 1$$

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$$

$$\sigma_{\bar{x}}^2 = \frac{2.0}{100}$$

$$\sigma_{\bar{x}} = \frac{1}{5}$$

$$\sigma_{\bar{x}} = \frac{1}{\sqrt{5}}$$

~~second~~

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

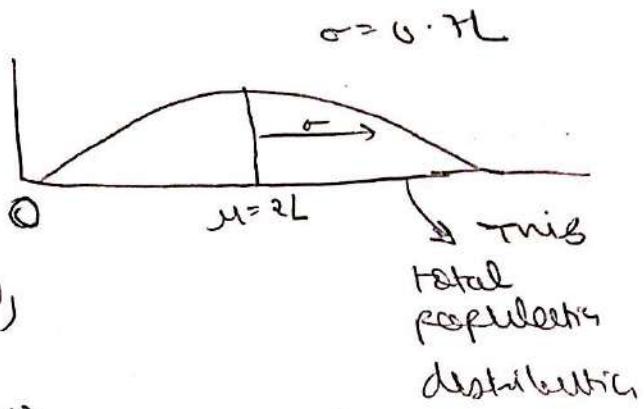
standard deviation of the
sampling distribution
of the sample mean

OR

standard error of
the mean.

Q Average male drinks 2L of water when active outdoor (~~which~~ with $\sigma = 0.7L$) ? You are planning a full day nature trip for 50 men and will bring 110L of water? What is probability you will run out?

Ans:

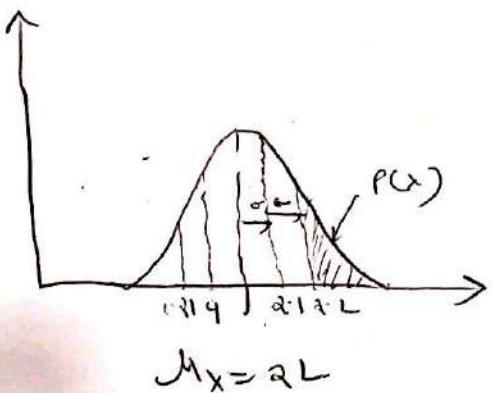


Now,

Now,

Sampling distribution of sample mean when

$$n = 50$$



$$P(X) = ?$$

$P(X) \approx$ It is, two standard deviations ~~more~~ away from the mean.

$$z = \frac{2.2 - 2}{0.099} = 2.02$$

It is ~~approx~~ 2.02 sd above the mean.

$$P(\text{run out})$$

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

~~total population~~ ~~water~~
= $P(\text{average water used per man is greater than } 2.2 \text{ L})$

$$\frac{110}{50} = 2.2$$

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

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

$$\sigma_{\bar{x}} = \frac{0.7}{\sqrt{50}} \approx 0.099$$

$$\sigma_{\bar{x}} \approx 0.099$$

According to z-table,

$$z = 0.9783$$

So the z-table gives the area below a particular value, hence it has given the area below 2.02.

hence,

$$\begin{aligned} P(x) &= 1 - 0.9783 \\ &= 0.0217 \\ &= 2.17\% \end{aligned}$$

Ques:

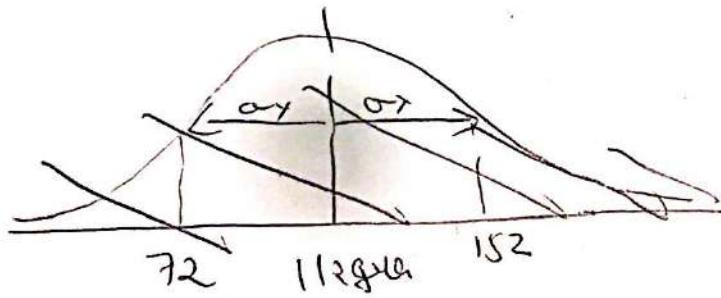
You sample 36 apples from your farm's harvest of over 200,000 apples? The mean weight of sample is 112g/cm (with a 40 g/cm sample standard deviation)?

Ans:

$$n = 36$$

$$\mu_x = 112 \text{ g/cm}$$

$$\sigma_x = 40 \text{ g/cm}$$



①

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

$$\Rightarrow \sigma_x = \sigma \sqrt{n}$$

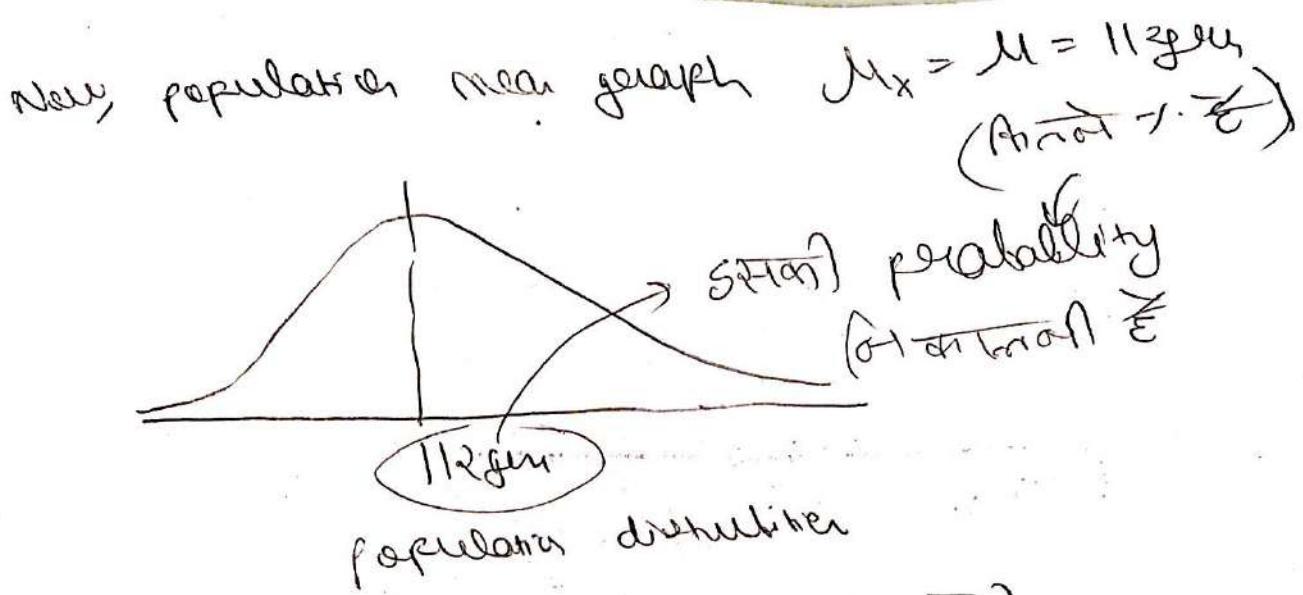
$$\Rightarrow \sigma = 4.0 \times \sqrt{36}$$

$$\Rightarrow \sigma = 40 \times 6 = 240$$

$$\sigma = 40$$

$$\sigma_x = \frac{40}{\sqrt{36}} = \frac{40}{6}$$

$$\sigma_x = 6.67$$

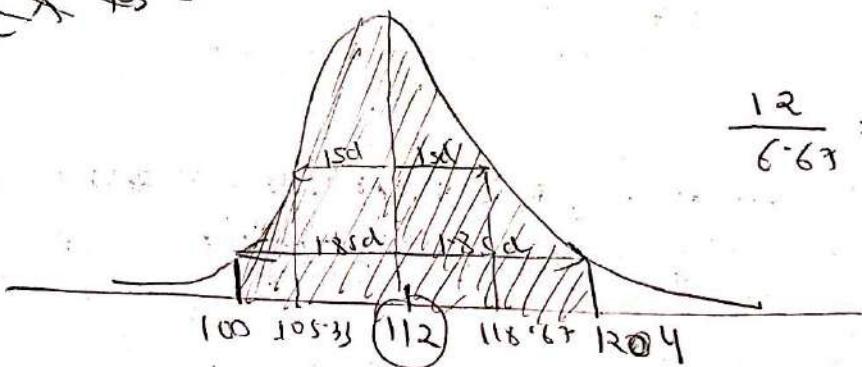


$P(\text{population mean is within 12 of } \bar{x})$
 (n)

$$= P(\bar{x} \text{ is within 12 of } \mu)$$

$$= P(\bar{x} \text{ is within 12 of } M_x) \quad [\text{as } \mu = M_x]$$

$\Rightarrow P(\bar{x} \text{ is within }$



Sample distribution of \bar{x} for $n=36$

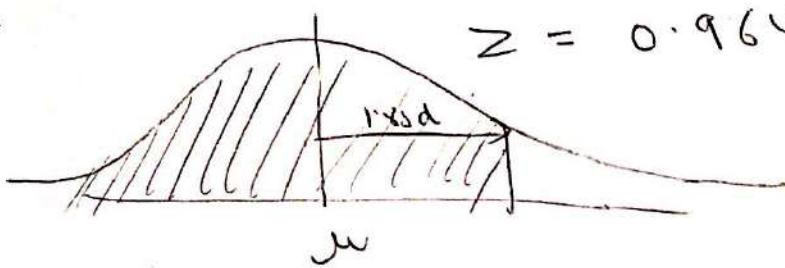
$$\sigma_x = 6.67 \text{ g/m}$$

Go to z-table,

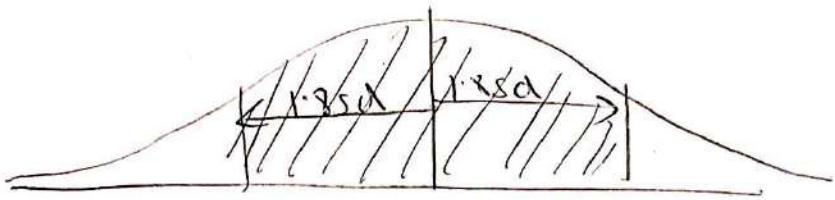
you will see z for $1.8 = 0.95$

It means:

$$z = 0.9641$$



But you have to find out this,



hence,

$$\begin{aligned} P(\text{required}) &= (0.964 - 0.5) \times 2 \\ &= \frac{0.964}{2} \times 2 \\ &= \frac{0.96}{2} \\ &= 0.9282 \\ &= 92.82\% \end{aligned}$$

~~probability distribution~~ only when an individual trial has
Bernoulli distribution :- only two possible outcomes
It is Bernoulli random variable

(discrete probability distribution)

QUESTION

* Toss a fair coin once? what is the distribution of number of heads?

$$\rightarrow p(\text{success}) = p$$

$$p(\text{failure}) = (1-p)$$

for $x=1$ if success occurs, $x=0$ if failure occurs

Then x has Bernoulli distribution:

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

$$P(X=2) = p^2(1-p)^{1-2}$$

for $x=0, 1$

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

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

$$\text{Mean of Bernoulli } \mu = p$$

$$\text{Variance of Bernoulli } \sigma^2 = p(1-p)$$

Q Approximateately 1 in 200 American adults are lawyers? one american adult is randomly selected? what is the distribution of the number of lawyers?

Ans: $P(X=2) = \left(\frac{1}{200}\right)^2 \left(1 - \frac{1}{200}\right)^{1-2}$
for $x=0, 1$

$$P(X=1) = \frac{1}{200} \quad P(X=0) = \frac{199}{200}$$

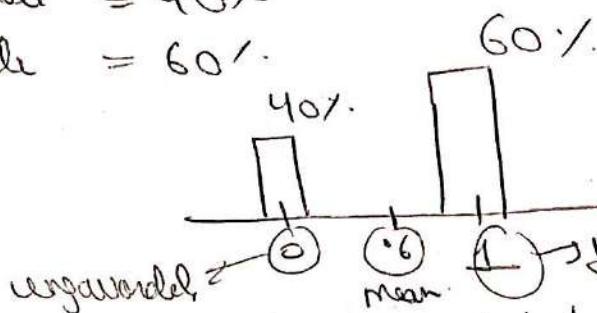
↓
we get one lawyer ↓
we get no lawyer

Now,

Suppose we are taking 2 trials

$$\text{unfavorable} = 40\%$$

$$\text{favorable} = 60\%$$



Mean of this distribution would be weighted value
(probability weighted sum)

$$\mu = 0.4 \times 0 + 0.6 \times 1$$

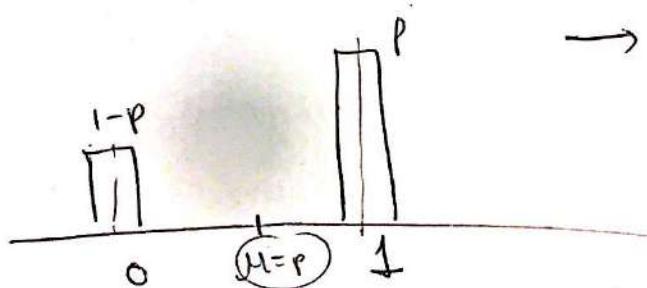
$$\mu = 0.6$$

Note, $\text{mean} = \text{probability of success} = p$

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

$$\sigma^2 = 0.24$$

$$\sigma = 0.49$$



→ ~~Bernoulli~~ Bernoulli distribution is the most simplified form of Binomial distribution.

$$\mu = (1-p) \times 0 + p \times 1 = p$$

$$\mu = p$$

mean of Bernoulli

$$\sigma^2 = (1-p)(0-p)^2 + p(1-p)^2$$

$$\sigma^2 = p - p^2 = p(1-p)$$

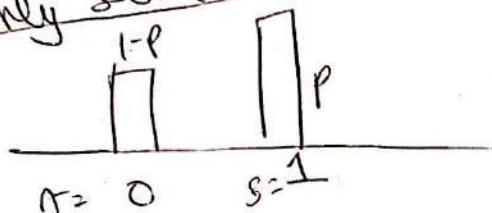
$$\boxed{\sigma^2 = p(1-p)}$$

Variance of Bernoulli

Margin of error :-

New Randomly surveyed 100 people from population :-

Let suppose



$$p = 43\%$$

$$(1-p) = 57\%$$

$$\bar{X} = \frac{px_0 + (1-p)x_1}{n}$$

sample mean

$$\bar{X} = \frac{57x_0 + 43x_1}{100} = 0.43$$

$$S^2 = \frac{s^2}{n} = \frac{(0 - 0.43)^2 + 43(1 - 0.43)^2}{(100 - 1)}$$

sample variance
of population

$$S^2 = 0.2475$$

$$S = 0.50 \quad \sigma_x = \frac{\sigma}{\sqrt{100}} = \frac{\sigma}{10} =$$

standard deviation

$$\sigma_x = \frac{0.50}{10} = 0.05$$

Margin of error :-

Confidence Interval :-

→ Confidence intervals are based on sample data, and give range of plausible values for a parameter.

e.g.:-

- we may be 95% confident that μ lies in interval $(-0.2, 3.1)$
- we may be 99% confident that μ lies in the interval $(2.5, 13.4)$

MARGIN OF ERROR :-

→ It tells you how far off your estimate is likely to be. or how much confidence you can in your estimate.

what you need for margin error :-

- Sample size = N
- Sample mean
- Standard deviation of population
- Z score for level of confidence.
 $(Z = 1.96)$
↳ for 95% confidence interval

How to calculate margin of error :-

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

standard error

$$\text{margin of error} = \text{standard error} \times z\text{ score}$$

Eg:- $N=49$

mean = 78

$\sigma = 5$

$z\text{ score} = 1.96$ (for 95% confidence interval)

Ans:-

$$\sigma_x = \frac{5}{\sqrt{49}} = \frac{5}{7} = 0.714$$

$$\text{standard error} = 0.714$$

$$\begin{aligned}\text{margin of error} &= 0.714 \times 1.96 \\ &= 1.40\end{aligned}$$

→ It means that 95% of the time, the population mean will be within 1.40 of the mean.

$$\text{confidence Interval} = \text{range } (\text{mean} - ME \text{ to } \text{mean} + ME)$$

confidence interval =

$$78 - 1.40 \text{ to } 78 + 1.40 \\ 76.6 \text{ to } 79.4$$

To get better confidence :-

- i) Have larger sd.
- ii) use larger sample size.
- iii) Accept broader confidence interval.

NOTE :-

→ confidence grows as the σ

margin of error increases.

* E is symbol that we use margin of err.

* C is confidence level.

$$E \uparrow \leftrightarrow C \uparrow$$

$$E \downarrow \leftrightarrow C \downarrow$$

Eg:- To estimate μ , we may use $\bar{x} \pm E$

where, margin of error.

$$E = z \times \frac{\sigma}{\sqrt{n}}$$

A May 2000 survey found that 38% of random sample of 1012 adults said that they believe in ghost? What is the margin of error for a 90% confidence interval for this poll?

Ans: success = 384 people
failure = 627 people

who believe in ghost
who don't

sample distribution :-

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

$$\sigma = \sqrt{pq}$$

$$\sigma_x = \sqrt{\frac{pq}{n}} = \sqrt{\frac{0.38 \times 0.62}{1012}}$$

$$\sigma_x = 0.0153$$

$$E = z \times \sigma_x$$

$$\left\{ z = 1.65 \text{ for } 90\% \text{ confid.} \right.$$

$$E = 1.6 \times 0.153$$

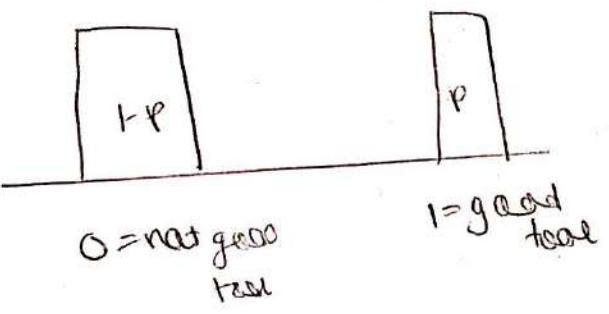
$$E = 0.252 \\ = 2.5\%$$

$$\text{margin of error} = 2.5\%$$

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

i) calculate 99% confidence interval for the proportion of teachers who felt that computers are essential tool.

ii) How could the survey be changed to narrow the confidence interval but to maintain 99% of confidence interval?



$$\bar{X} = \frac{1 \times 142 + 0 \times 108}{250}$$

~~$$\bar{X} = \frac{142}{250}$$~~

$$p = \frac{142}{250}$$

$$\bar{X} = p = \frac{142}{250} = 0.568$$

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

$$\sigma = \sqrt{\frac{142}{250} \times \frac{108}{250}} = \text{approx} 0.5$$

~~$$\sigma_{\text{sample std.}} = \sqrt{\frac{0.568 \times 0.426}{250}} = 0.010$$~~

$$\sigma_x = \frac{\sigma}{\sqrt{n}} = \frac{0.50}{\sqrt{50}}$$

$$\sigma_x = 0.031$$

$$E = 0.031 \times 2.58$$

$$E = 0.07$$

\rightarrow we are 99% confident, since that \bar{x} is within 0.08 of the population proportion (p)

$$\text{confidence interval} = 0.568 \pm 0.08$$

$$0.488 \text{ to } 0.648$$

- (ii) you can take small samples
small sample size confidence interval:
 patients blood pressure have been measured
 after having been given a new drug for
 3 months? They had BP increases of
 1.5, 2.9, 0.9, 3.9, 3.2, 2.1 and 1.1?
 construct 95% confidence interval for
 true expected BP increase for all
 patients in population,

Ans: sample mean

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

$$\bar{x} = 2.34$$

$$s^2 = (1.5 - 2.34)^2 + (2.9 - 2.34)^2 + (0.9 - 2.34)^2 + \dots \text{ (find it)}$$

$$s = 1.04$$

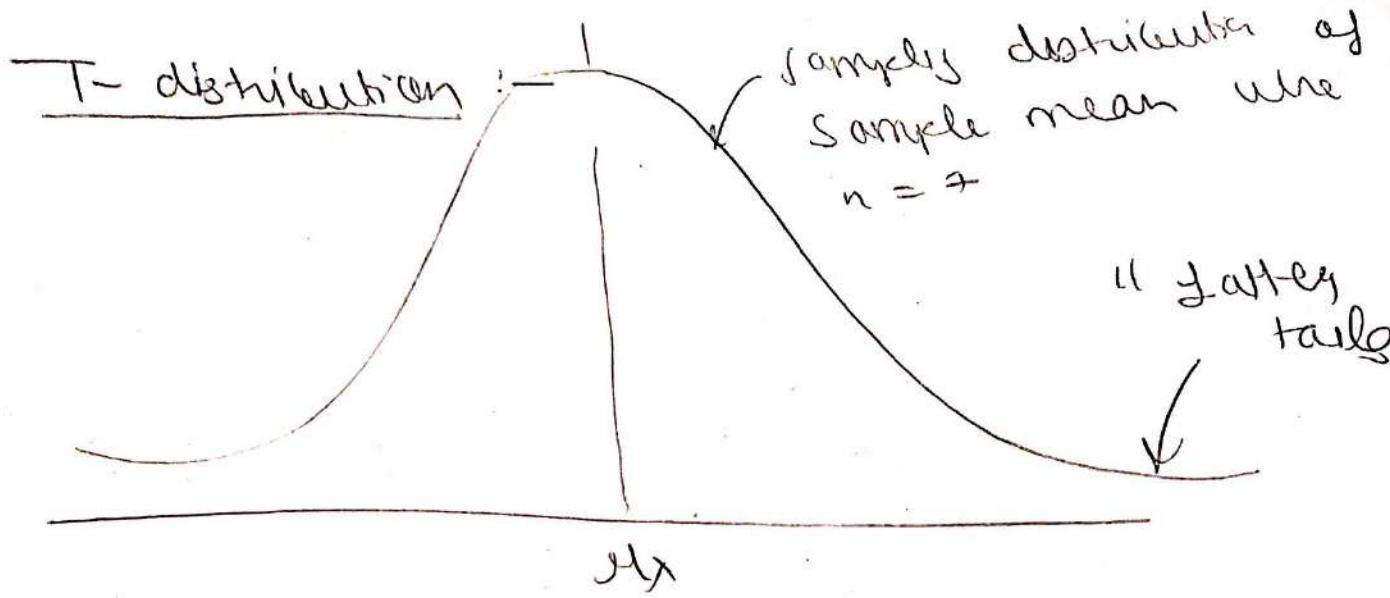
→ standard deviation

$$\sigma \approx s = 1.04$$

→ not so good because
n is small
 $n = 7$.

→ So instead of making the normal distribution ~~graph~~ graph of the sample mean we make t-distribution graph because our n is small
 $n = 7$.

→ when we have small sample size then,



$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{7}} \approx \frac{s}{\sqrt{7}}$$

$$\sigma_{\bar{x}} = \frac{1.04}{\sqrt{7}} = 0.39$$

this is the bad estimator we have taken $\sigma = s$, when n is small
 population sd sample sd

→ → T distribution is used when n is small, here

$$E = T \times \sigma_{\bar{x}}$$

→ T score gives area towards the right in confidence

$$E = 2.447 \times 0.39$$

Z score gives

$$E = \text{approx} 0.96$$

area towards left.

"confident"

→ 95% chance \bar{x} is within 0.96 of

→ 95% chance μ is within 0.96 of \bar{x}
 $2.34 - 0.96, 2.34 + 0.96 \Rightarrow 1.38 \text{ to } 3.32$

Descriptive - Statistics

→ Descriptive statistics is the set of numbers that describe the data.

Two categories :-

- 1) Measures of central tendency
- 2) Measures of Dispersion

Measure of central tendency :-

Average ← ① mean → = Average (c_1, c_2, \dots, c_{100})
middle observation ← ② median → = Median (n_1, n_2, \dots)
that divides data into two parts
Q: When we should report median and when else should report mean?

Ans: Seven employees in firm with following salaries:-

28000 \$
33000 \$
33000 \$
34000 \$
37000 \$
40000 \$
400000 \$

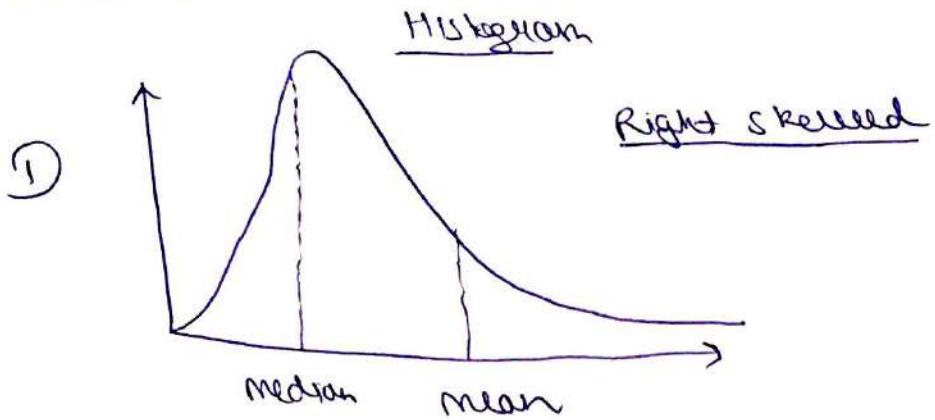
$$\text{Mean} = 86000 \$$$

$$\text{Median} = 34000 \$$$

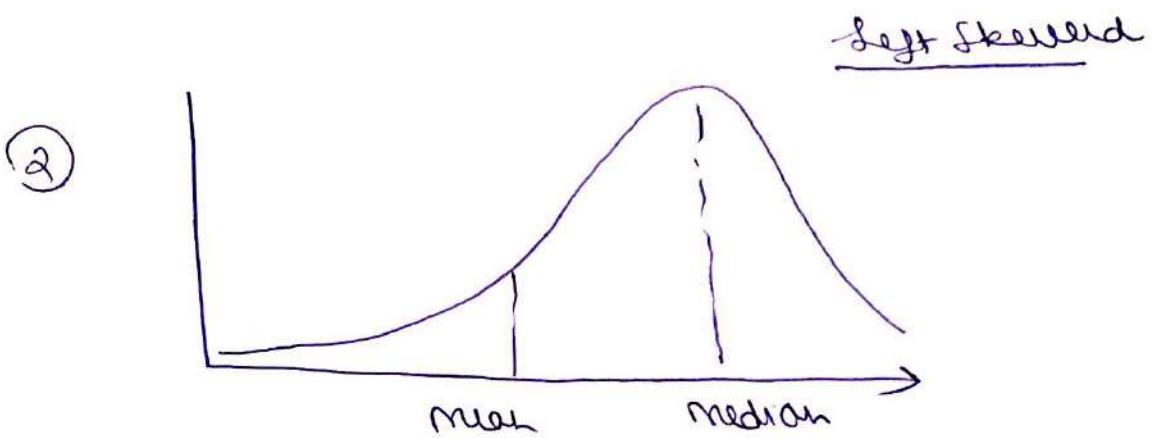
→ Hence we report median over here, since it is more appropriate.

→ Mean is influenced to a greater extent by outliers.

Mean vs Median :-



mean > median.



median > mean.

Mode :- frequency of data.

$$= \text{Mode} \cdot \text{SNGL}(n_1, n_2, \dots)$$

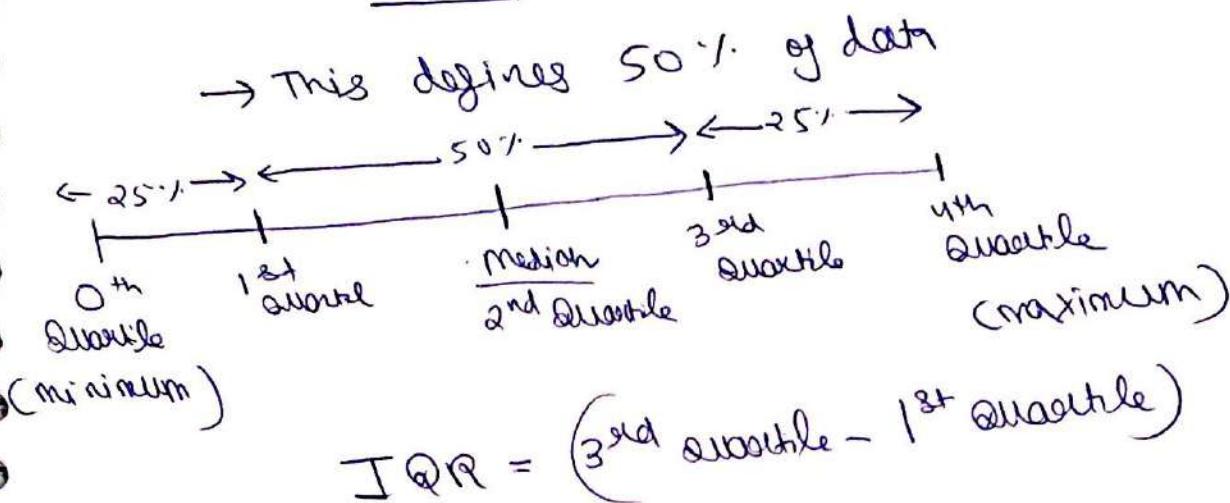
measures of spread :-

① Range :- Max - Min.

→ Higher range means greater spread in data.

Inter Quartile Range :-

IQR



$$IQR = (3\text{rd quartile} - 1\text{st quartile})$$

→ We prefer Inter Quartile range over range taking because in statistics we prefer sample of data.

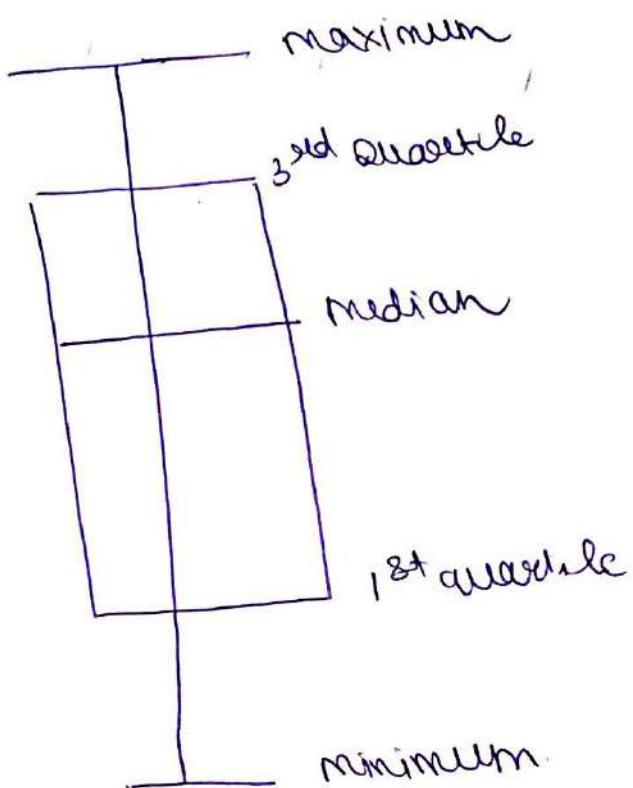
Ex:-

=QUARTILE.INC (array, quart)

$IQR = \text{QUARTILE.INC}(\text{array}, 3) - \text{QUARTILE.INC}(\text{array}, 1)$

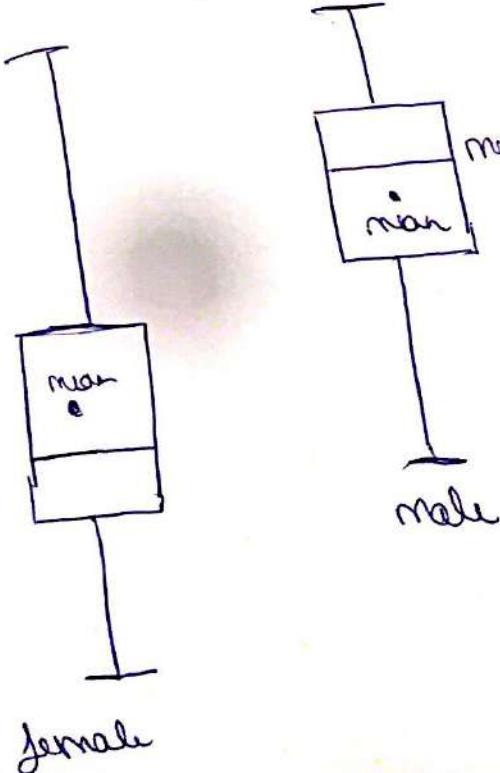
Box Plot :-

visual representation of the mean and various quartiles of data.



Eg:-

Hollywood stars earning



* we can infer that spread of females are more than spread of males.

* max. salary of female is equivalent to max. salary of male.

* Female is right skewed hence mean > median.

* male is left skewed hence median > mean.

③ Standard deviation

→ describes dispersion in data.

→ It calculates mean & then how far other points are deviated from mean.

$$= \text{STDEV.P}(\text{number1}, \text{number2}\dots)$$

→ If we have data which was sample from some larger population of data we will use.

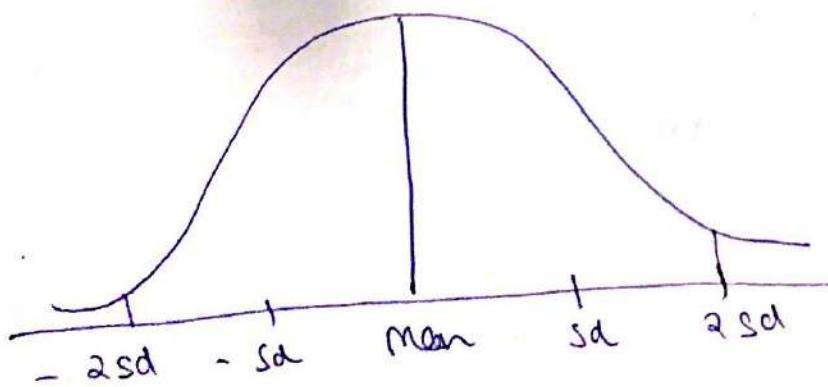
$$= \text{STDEV.S}(\text{number1}, \text{number2}\dots)$$



④ Variance :-

$$\text{Variance} = (\text{std})^2$$

NOTE :- 68% of data lie within one standard deviation & 95% lie within 2 standard deviation from mean.



Assumption :- distribution should be Normal.

Chebyshen's theorem :-

→ At least $\left(1 - \frac{1}{K^2}\right)^{\text{th}}$ of data lie within $\pm K$ standard deviations from the mean regardless of shape of distribution.

→ Specifically, It says that at least 75% of all values are within ± 2 std from mean regardless of shape of distribution.

measures of Association :- $\begin{cases} \text{covariance} \\ \text{correlation} \end{cases}$

Covariance :-

relationship between two variables.

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

$$= \text{COVARIANCE} \cdot S(\text{range}_1, \text{range}_2)$$

→ If covariance is +ve then it indicates positive correlation.

NOTE :- Covariance measure is susceptible to the unit of measurement, we can arbitrarily inflate or deflate covariance by choice of units.

Correlation :-

→ hence, covariance is not very appropriate.

→ For that correlation corr,

$$\text{correlation} = \frac{\text{covariance}(X, Y)}{\text{std}(X) \text{ std}(Y)}$$

$$= \text{corr}(\text{range}_1, \text{range}_2)$$

NOTE :-

Covariance:

* Range : $-\infty$ to $+\infty$

* Affected by units of measurement

Correlation:

* Range : -1 to $+1$

* No effect by unit of measurement

Causation

: -

smoking causes cancer

cancer causes smoking

To establish Causation:

* Two variables must be correlated

* causing variable must occur before caused variable.

* external variables should be ruled out.

Probability and Random Variables :-

Statistical distributions :-

Random Experiment → Random Variable

multiple possibilities → "Salary"
e.g. CEO salaries

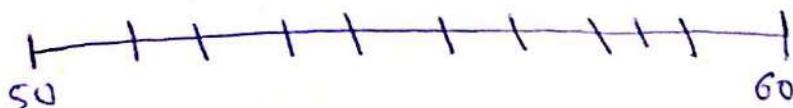
→ statistical distribution is the tool
that will help us model the salary for CEO
across small firms. Benefit is that we can
make prediction across different set of firms
not included in the data from which the
histogram is created.

Statistical distributions :-

- ① Beta
- ② Binomial
- ③ Gamma
- ④ Poisson
- ✓ ⑤ Normal
- ⑥ T distribution

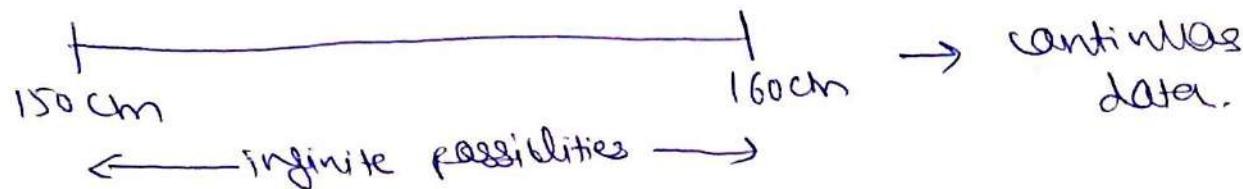
Discrete data vs Continuous data :-

- ① Number of Students :-



finite
(discrete)

③ Heights of Man and women



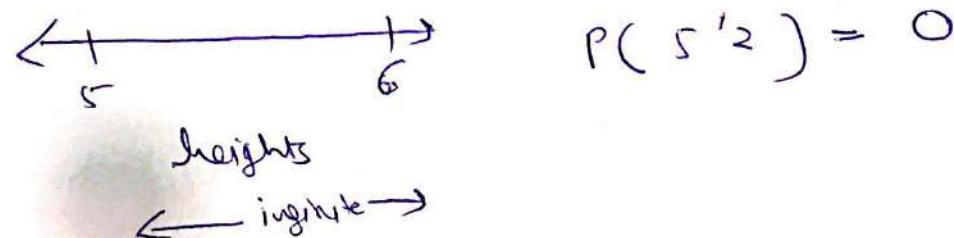
Assumption :- Many times, even though data may be strictly discrete, for eg :- annual salaries of employees, revenues earned by companies etc. we will still end up using continuous distribution to model the data.

Probability density function :-

→ It is the rule that assigns probabilities to various possible values that random variable takes.

Note:- probability of a particular outcome is always zero in continuous distribution.

Eg:-



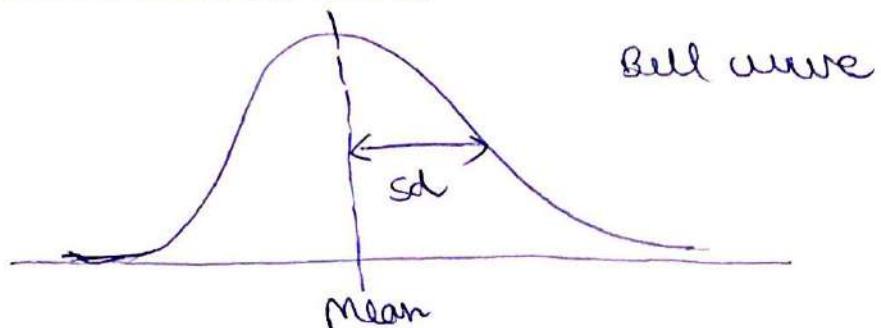
hence this is the reason that we use range of continuous data.

Eg:- $P(\text{height} > 5'2 \text{ and less} < 5'5) = ?$

$P(\text{height} < 5\text{ft}) = ?$

→ For discrete distribution we know how to calculate probability
 → But for continuous distribution we will use ranges.

Normal Distribution :-



→ Mean can have any value between

- ∞ to $+\infty$.

$$f(x) = \frac{1}{\sqrt{2\pi \times \text{std}^2}} e^{-(x-\text{mean})^2 / 2(\text{std})^2}$$

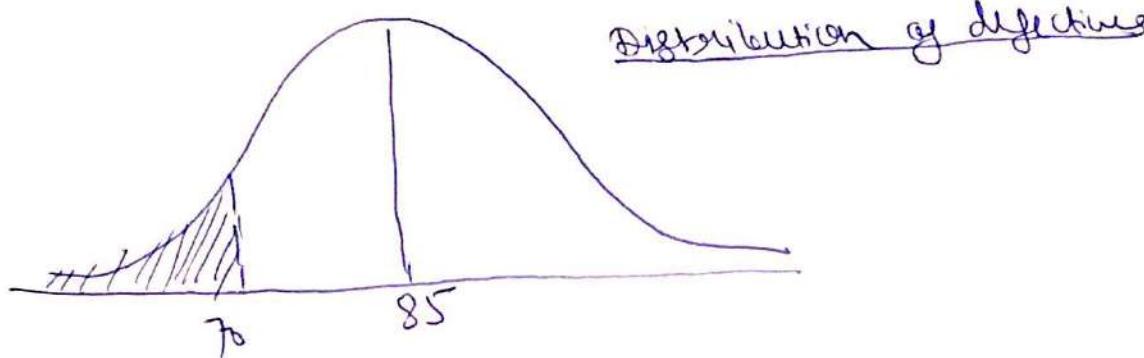
NOTE :-

- Q A bread producing company produces whole wheat loaves of bread in its factory. It observes that on average every day 85 loaves of bread get discarded on account of being defective. Std is 9 loaves ? Find probability that less than 70 loaves of bread will get discarded on being defective.

Ans:

Defective \sim Normal (85, 9)

Prob (Defective < 70) = ?



Syntax :-

= NORM.DIST (x , mean, std, TRUE)

$$= 0.047$$

$\downarrow 70$
 $\downarrow 85$

= 4.7% chance

we use
parameter
never use
FALSE

this will
give area under
curve to left side

NOTE:-

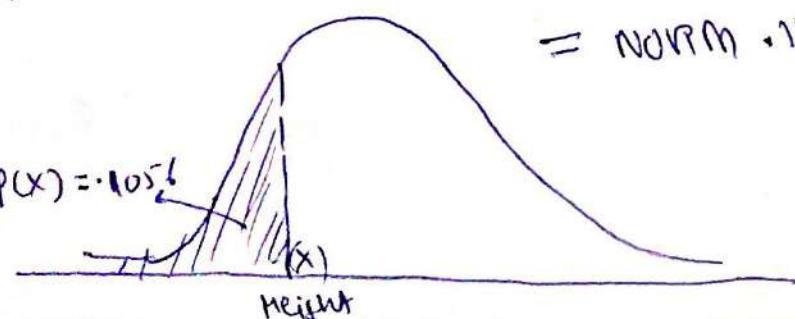
Prob (Defective < 75) = Prob (Defective ≤ 75)

Now,

For eg:- you are given probability value
and you need to find x-value that is height
value , you will use inverse function

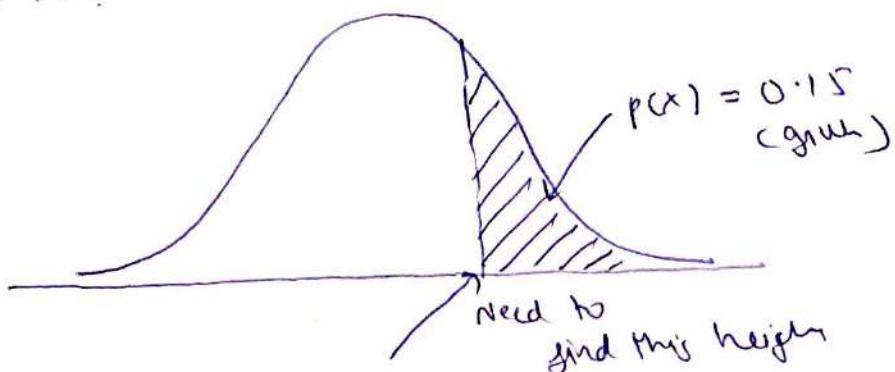
= NORM.INV (probability, mean, std)

$$P(x) = 0.95$$



Q Find out a height such that probability Soner's height is greater than this particular height is 0.15?

Ans:



$$\begin{aligned} P(\text{left}) &= 1 - 0.15 \\ &= 0.85 \end{aligned}$$

$$= \text{NORM.} \cdot \text{INV}(0.85, \text{mean}, \text{std})$$

\downarrow \downarrow
5.5 0.4

Q John can take two roads to airport from his home.

Road A : mean = 54 minutes , std = 3 mins

Road B : mean = 60 mins , std = 10 mins

Q which road should he chose if on midday Sunday he must be at airport within 50 minutes to pick up his spouse?

Ans:

Road A

$$\text{Prob}(\text{Time} < 50) = \text{NORM.DIST}(50, 54, 3, \text{TRUE})$$

$$= 0.09$$

Road B

$$\text{Prob}(\text{Time} < 50) = \text{NORM.DIST}(50, 60, 10, \text{TRUE})$$

$$= 0.15$$

Binomial Distribution :-

Two popular Discrete distribution :-

- ① Binomial
- ② Poisson

Bernoulli Process :-

Situation where random variable have only two possible outcomes.

Eg :- pass / fail

→ NOTE :- Repeated trials to the Bernoulli process gives rise to binomial distribution.

Eg -- Probability of winning at least 4 times in 10 rolls of the dice?

win (6) loose (0)

$$P(X) = nC_x p^x (1-p)^{n-x}$$

$$= \text{BINOM. DIST} (x, n, p, \text{FALSE/TRUE})$$

↓ ↓ ↓
 No. of n. of P of
 success trials success

↓
 (success = x)

Q Probability that you win 3 times in 10 rolls of dice?

$$P(X=3) = \text{BINOM.DIST}(3, 10, 0.1667, \text{FALSE}) \\ = 0.1551$$

Q Probability that you win at most five times in 10 rolls of dice?

$$P(X \leq 5) = \text{BINOM.DIST}(5, 10, 0.1667, \text{TRUE}) \\ = 0.98$$

Q Probability that you win at least three times in 10 rolls of dice?

$$P(X \geq 3) = 1 - P(X \leq 2) = 1 - \text{BINOM.DIST}(2, 10, 0.1667, \text{TRUE})$$

$$= 0.22$$

Mean of binomial distribution = np
 $sd = \sqrt{npq}$

Poisson distribution :-

- * No. of fatal + traffic accidents in a city in week
- * No. of customers arriving at store during given hour.

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!} = \text{POISSON.DIST}(\lambda, x, \text{FALSE/TRUE})$$

→ In binomial distribution the random variable can take the value from 0 to till trial ends. but here in poisson random variable can take (0-∞)

NOTE :- $= T\text{-INV}(\text{probability, df})$

↳ we are given probability ~~and~~ to left and
we need to find the value of \bar{x}

→ T distribution is symmetric about zero.

Confidence interval :-

Interval with some confidence

Eg:- Take elections.

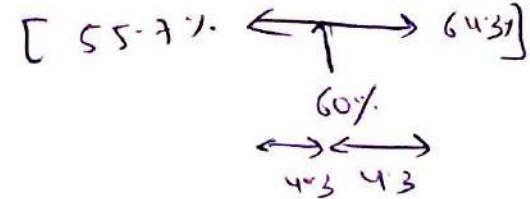
(A) (B)

300 votes 200 votes

$P(A) = 60\%$, but this 60% is actually the probability from the sample, will candidate A get 60% of all votes in actual elections?

True confidence interval comes.

95% confidence interval



$$z\text{-statistic} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

7.43 is margin of error

$$T\text{-statistic} = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

Whenever we know population std. dev.
use z statistics and when we have data
about the sample we use T statistics.

Another continuous distribution

T distribution

→ Similar to normal distribution.
→ However standard deviation of this distribution depends on single parameter called degree of freedom.

✓ → As degree of freedom increases it tends towards normal distribution.

→ It is used as an tool to calculate confidence intervals and hypothesis testing.

$=T.DIST$

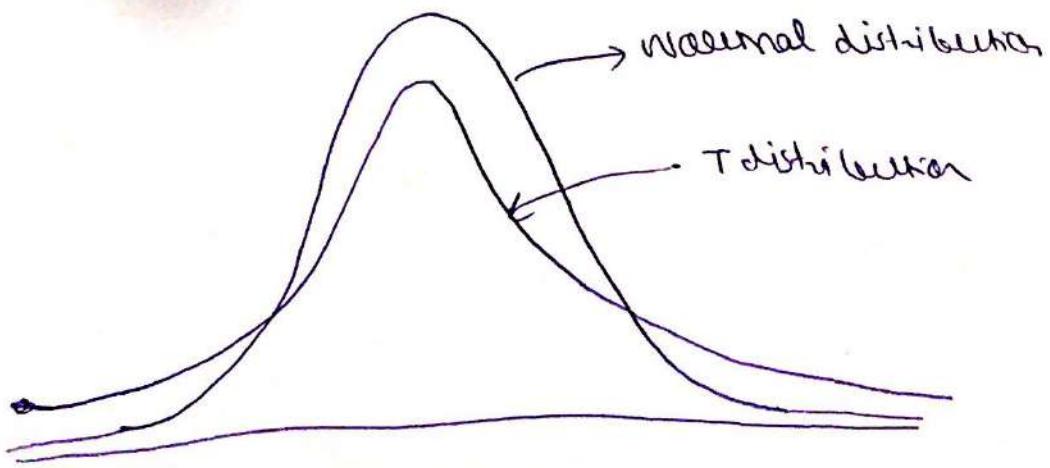
and $=T.INV$

$$P(x) = T.DIST (x, df, TRUE)$$

↳ we get area under curve left of x.

* degree of freedom is the parameter of t distribution.

→ larger set of data would have greater degrees of freedom.



Now,

To find $z_{\alpha/2}$

$$z_{\alpha/2} = \text{NORM. INV} (\alpha/2, 0, 1)$$

mean
standard deviation

$$\text{margin of error} = |z_{\alpha/2}| * \frac{\sigma}{\sqrt{n}}$$

$$= \text{CONFIDENCE-NORM} (\alpha, \sigma, n)$$

when population standard deviation (σ) is not known

* we replace it by sample standard deviation.

* z statistics ($z_{\alpha/2}$) get replaced by t statistics. ($t_{\alpha/2}$)

$$\bar{x} - |t_{\alpha/2}| \frac{s}{\sqrt{n}} < \mu < \bar{x} + |t_{\alpha/2}| * \frac{s}{\sqrt{n}}$$

$$\text{here margin of error} = \text{CONFIDENCE-T} (\alpha, s, n)$$

↑
significance level
↓
std
↓
sample size

→ if population standard deviation

is not known we use z statistics and if

sample standard deviation is not known we use

t statistics.

Confidence Interval for population proportion :-

- Q A consultancy firm surveyed a randomly selected set of 210 CEO of fast growing small companies in US and Europe. Only 51% of these executives had management succession plan in place, remaining did not have one.

Use this information to compute a 90% confidence interval to estimate the proportion of all 'fast growing small companies' that have management succession plan?

$$\hat{p} - |z_{\alpha/2}| \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

margin of error.

$$\hat{p} - |z_{\alpha/2}| \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} < p < \hat{p} + |z_{\alpha/2}| \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

actual proportion
population which
is unknown to us
& we wish to guess
it with some
information we
have.

Sample Size :-

Q How big should our sample is?

A Quality control manager at battery manufacturer wants to estimate the average number of defective batteries contained in a box shipped by a company? How many boxes does she needs to open to figure out average number of defective batteries contained in a box?

Given, margin of error = ± 0.3 batteries

Confidence Interval = 95%

Population Std = 0.9

Ans:- margin of error = $0.3 = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

$$0.3 = \left| \text{NORM-INV}(0.05/2, 0, 1) \right| \times \frac{\sigma}{\sqrt{n}}$$

$$0.3 = 1.96 \times \frac{0.9}{\sqrt{n}}$$

$$\sqrt{n} = 5.88$$

$$n = 35$$

Q Sample size calculation when we have population proportion greater than mean

Q A pollster wanting make a prediction about a particular candidate.

How many voters should the pollster survey?

Pollster in a prediction wants to have margin of error $\pm 3\%$ with confidence level 95%?

Ans:

$$\text{margin of error} = |z_{\alpha/2}| * \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$0.3 = |\text{NORM.INV}(0.05/2, 0, 1)| * \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$\downarrow 1.96$

$$n = \frac{(1.96)^2}{(0.03)^2} * \hat{p}(1-\hat{p})$$

\hat{p} is the sample proportion who would vote for that particular candidate.

Now,

we use a conservative estimate of \hat{p} .

$$\hat{p}(1-\hat{p})$$

$\underbrace{\quad}_{\text{For this to be maximum}}$

$$\hat{p} = 0.5$$

Hypothesis testing and p-value :-

Neurologist is testing the effect of drug on response time by injecting 100 rats with a unit dose of drug, subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with drug is 1.2 seconds. Mean of 100 injected rats response time is 1.04 seconds with sample $s.d = 0.5$ seconds. Do you think that drug has effect on response time?

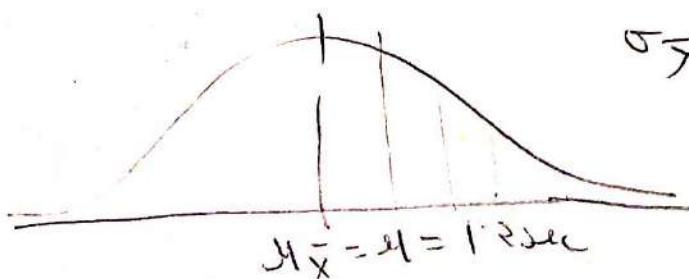
Null Hypothesis H_0 :

↳ drug has no effect
on response time

H_0 : drug has no effect $\Rightarrow \mu = 1.25$ sec even with drug

H_1 : drug has an effect $\Rightarrow \mu \neq 1.2$ sec when drug is given

Let assume H_0 is true



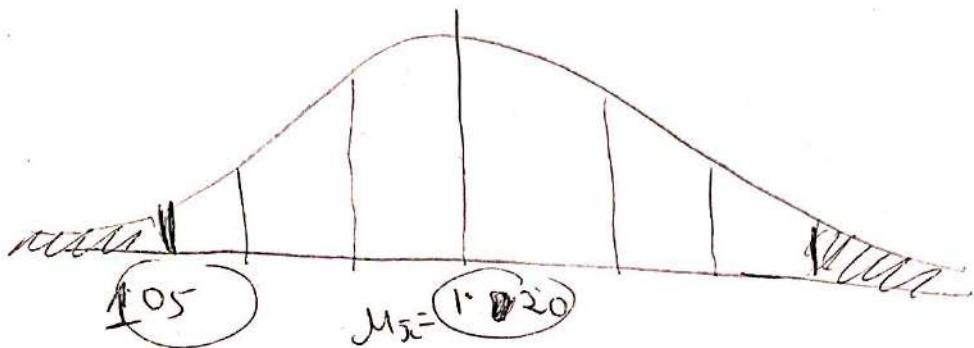
$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{100}} \approx \frac{0.5}{\sqrt{100}} = \frac{0.5}{10} = 0.05$$

$$\sigma_{\bar{X}} = 0.05$$

$$Z = \frac{1.2 - 1.05}{0.05}$$

$$\Rightarrow Z = 3$$

\rightarrow mean



$$\begin{aligned} \text{To find this area, } &= 100 - 99.7 \\ &= 0.3\%. \end{aligned}$$

\rightarrow Hence if we assume that drug has no effect probability of getting the sample is only 0.3% .. ~~has occurred~~
because if null hypothesis is being rejected because if null hypothesis is true then there is only 0.3% by chance of getting this result if null hypothesis is true.

$p \text{ value} = 0.003$

\rightarrow Hence, drug has some effect

$$n = 100 \quad c = 0.99 \quad \alpha = 1 - 0.99 \\ \alpha = 0.01$$

2 doctors believe that average teen sleeps on average no longer than 10 hours / day ?
A researcher believes that teens on average sleep longer ? write H_0 and H_1 ?

Ans: $H_0: \mu \leq 10 \text{ hours / day}$ ✓
 $H_1: \mu > 10 \text{ hours / day}$ ✓

school board claims that at least 60% of students bring phone to school ? Teacher believes this number is too high & randomly samples 25 students to test at level of significance of 0.02 ? write H_0 and H_a ?

Ans: $H_0: p \geq 0.60$ $\alpha = 0.02$
 $H_1: p < 0.60$ $c = 0.98$
 $= 98\%$
 $n = 25$

Large sample

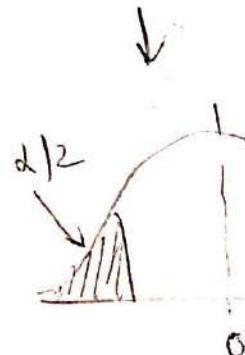
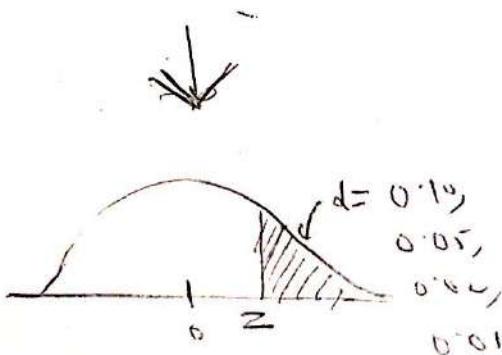
Hypothesis testing for means & large samples

- Here we will take large sample as $n \geq 30$. If $n < 30$, then it is small sample & we will use T distribution there. but, if it is large sample we will be using Normal distribution.

$$\text{Degree of freedom} = n - 1$$

→ T distribution approaches Normal distribution when you have large number of samples.

<u>c</u>	<u>d</u>	<u>one tail test</u>	<u>two tail test</u>
0.90	0.10	$t_{2.8}$ r ² value	± 1.645 r ² value
0.95	0.05	1.645	± 1.96
0.98	0.02	2.05	± 2.33
0.99	0.01	2.33	± 2.57



Hypothesis - Testing

In hypothesis testing, we turn a question of interest into hypotheses about the value of a parameter or parameters.

We create: Sometimes referred to as the status quo hypothesis

- A null hypothesis, denoted by H_0
- An alternative hypothesis, denoted by H_a

← Sometimes referred to as the research hypothesis

H_0 is the method of statistical inference. It is an assumption that can be tested on a statement

We calculate an appropriate test statistic (based on sample data), and determine how much evidence there is against the null hypothesis.

If the evidence is strong enough, (if it meets a certain *significance level*), we can reject the null hypothesis in favour of the alternative hypothesis.

Is the average weight of a certain type of candy bar different from the desired 58 grams?

$$H_0: \mu = 58$$

$$H_a: \mu \neq 58$$

- ▶ Is there significant evidence against the null hypothesis?
(Two approaches: the rejection region approach, and the p -value approach.)
- ▶ What is an appropriate conclusion to the problem?

Z Test for one Mean :-

Normal distribution

> help(pharm)

> ?mean # $P(X <= 70)$

> pharm ($N = 70$, mean = 75, $sd = 5$, lower.tail = T, digits = 6) 0.158

Hypothesis Tests for the Population Mean μ

(When sampling from a normally distributed population)

> # $P(X \geq 85)$

> pharm ($N = 85$, mean = 75, $sd = 5$, level, tail = F) more N samples

Two scenarios:

Rare in practice (but a useful starting point)

► σ is known. Z test

► σ is not known. t test

Much more common

less N samples

T-distribution :-

These can be used to find p-values or critical values for constructing confidence intervals for statistic that follows t distribution.

> help(pt)

> ?pt

t-statistic = 2.3, one sided.

we want to find $p(t > 2.3)$

> $p_t(v = 2.3, df = 25, \text{lower.tail} = F)$
0.0150

Is there strong evidence that the population mean μ is different from some value that is of interest to us?

Is it different from some *hypothesized* value?

two sided p values

> $p_t(v = 2.3, df = 25, \text{lower.tail} = F) + p_t(v = -2.3, df = 25, \text{lower.tail} = T)$

Population mean Hypothesized value

The null hypothesis is $H_0: \mu = \mu_0$ ↙ e.g.

$H_0: \mu = 10$

We choose one of the possible alternatives:

- ▶ $H_a: \mu < \mu_0$ } One-sided alternatives
- ▶ $H_a: \mu > \mu_0$ } (one-tailed tests)
- ▶ $H_a: \mu \neq \mu_0$ } Two-sided alternative
(two-tailed test)

find 1 for 95% confidence

value of t with 2.5% in each tail

> $qt(p = 0.025, df = 25, \text{lower.tail} = T)$
-2.06

parametric T test and confidence intervals
They are used to examine difference in means
in the population. Also to find relationship between
categorical variable and a continuous variable

The appropriate choice of alternative hypothesis depends on the problem at hand, and *should not be based on the current sample's data.*

You should not use the same data that suggests a hypothesis to test that hypothesis.

Suppose we have a simple random sample of n observations from a normally distributed population where σ is known.

- > help(t.test)
- > t.test
- > boxplot(lungcap ~ smoking)
 H_0 : mean lungcap of smokers = mean non-smokers
 H_a : \neq (two-sided)
- > t.test(lungcap ~ smoke, mu=0, alt="two.sided")
 $(\text{csg} = 0.95,)$
p-value = 0.00039
confidence interval = (-1.3 - 0.4)

To test $H_0: \mu = \mu_0$, we use the test statistic:

$$\frac{Z}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

where $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$

If H_0 is true, the Z test statistic will have the standard normal distribution

Example:

Suppose a supplier to a sushi restaurant claims their bluefin tuna contains no more than 0.40 ppm of mercury on average.

The owner fears the supplier's claim is incorrect, and the average mercury level is higher.

μ is the true mean mercury content
of bluefin tuna from this supplier

We will test the null hypothesis:

$$H_0: \mu = 0.40$$

$$H_a: \mu > 0.40$$

In a random sample of 16 pieces of tuna
from this supplier, there was an average of
0.74 ppm of mercury.

Does this yield strong evidence that the true
mean mercury content is greater than 0.40
ppm?

(Suppose $\sigma = 0.08$ ppm)

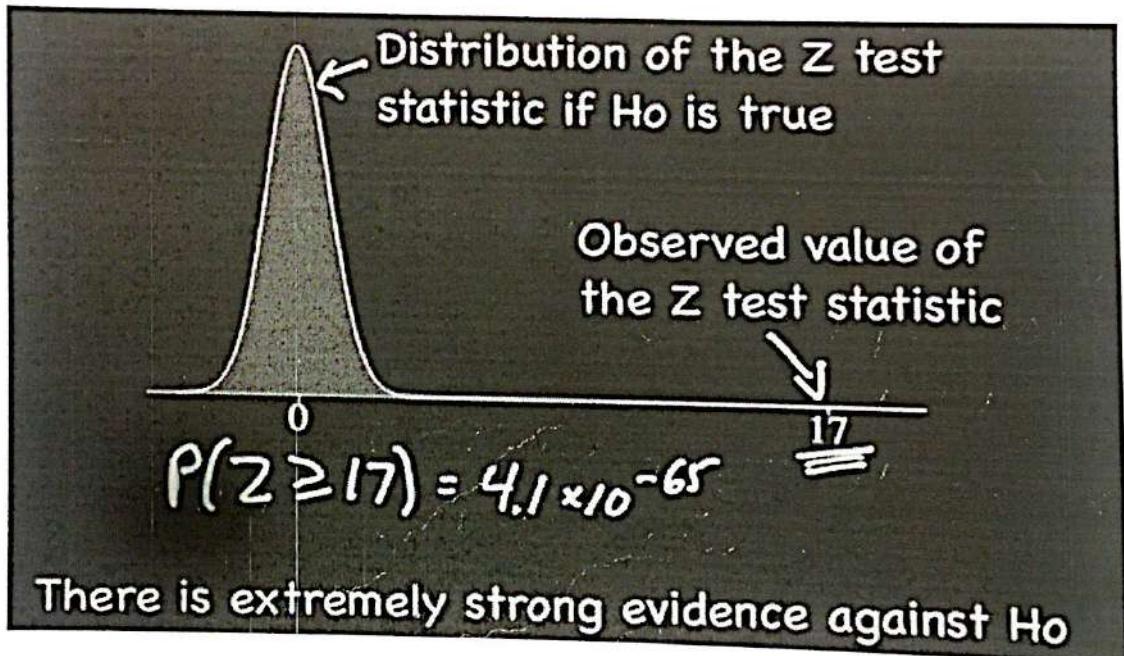
$$\begin{array}{c} H_0: \mu \leq 0.40 \\ H_a: \mu > 0.40 \end{array}$$

$H_0: \mu = 0.40$, $H_a: \mu > 0.40$

$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}} = \frac{0.74 - 0.40}{0.01}$$

$$\underline{\bar{X} = 0.74}, \underline{\sigma = 0.08}, \underline{n = 16} \quad = 17$$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} = \frac{0.08}{\sqrt{16}} = 0.02$$



Here the evidence against the null hypothesis is overwhelming.

(We would almost never see what was observed if the null hypothesis were true.)

But how far out does the test statistic need to be before we can say the evidence is *significant*?

Before we can *reject* the null hypothesis in favour of the alternative?

To judge whether the evidence against the null hypothesis is significant, we use one of two approaches:

- ▶ The rejection region approach
- ▶ The p -value approach

Rejection Region Approach :-

Z Tests for the Population Mean μ

The Rejection Region Approach

Suppose we are sampling from a normally distributed population, and σ is known.

To test $H_0: \mu = \mu_0$, we use the test statistic:

$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}} \quad \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

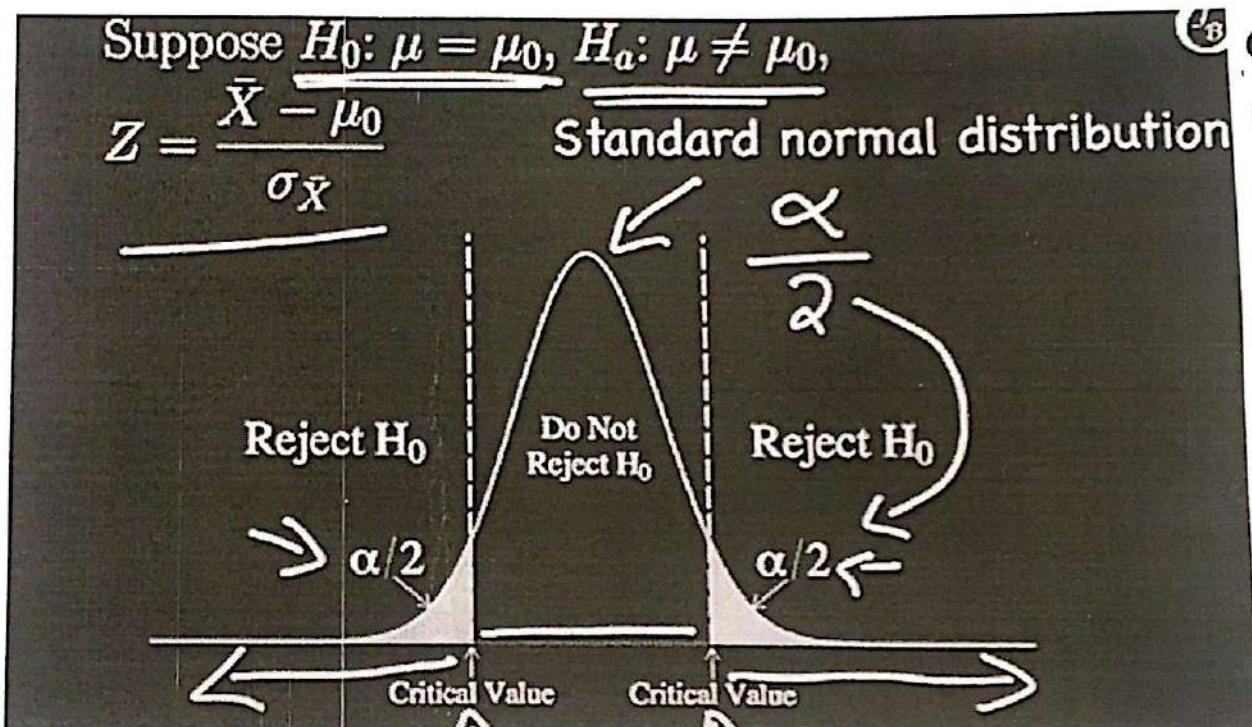
If H_0 is true, the Z test statistic will have the standard normal distribution

The rejection region approach:

- ▶ Choose a value for α , the significance level of the test.
 $(\alpha$ is the probability of rejecting the null hypothesis if it is true.)

People
often choose $\alpha = 0.05$

- ▶ Find the appropriate rejection region.
- ▶ Reject the null hypothesis if the test statistic falls in the rejection region.



Suppose $H_0: \mu = \mu_0$, $H_a: \mu \neq \mu_0$, $\alpha = 0.05$

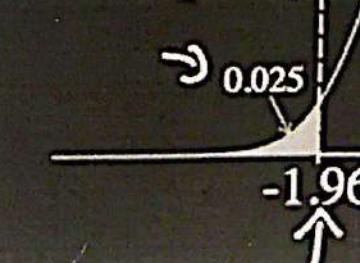
$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}}$$

Reject H_0 if

$$Z \geq 1.96 \text{ or}$$

$$Z \leq -1.96$$

Reject H_0



$$\frac{0.05}{2} = 0.025$$

Reject H_0

$$Z_{0.025} = 1.96$$

Suppose $H_0: \mu = \mu_0$, $H_a: \mu < \mu_0$, $\alpha = 0.05$

$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}}$$

Reject H_0 if

$$Z \leq -1.645$$

Reject H_0

→ 0.05

Do Not Reject H_0

-1.645



Suppose $H_0: \mu = \mu_0$, $H_a: \mu > \mu_0$, $\alpha = 0.05$

$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}}$$

Reject H_0 if

$$Z \geq 1.645$$

Reject H_0

0.05

Do Not Reject H_0

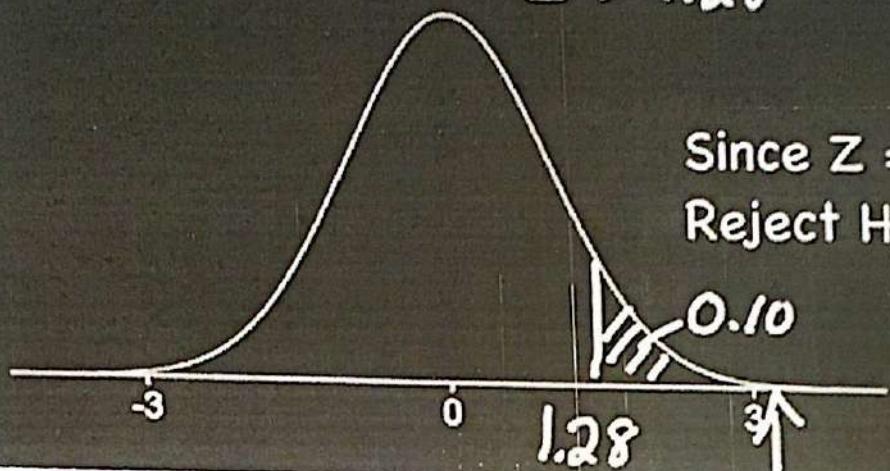
1.645

Suppose $H_0: \mu = \mu_0$, $H_a: \mu > \mu_0$, $\alpha = 0.10$

$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}} = 3.1.$$

Reject H_0 if

$$Z \geq 1.28$$



Since $Z = 3.1$,
Reject H_0

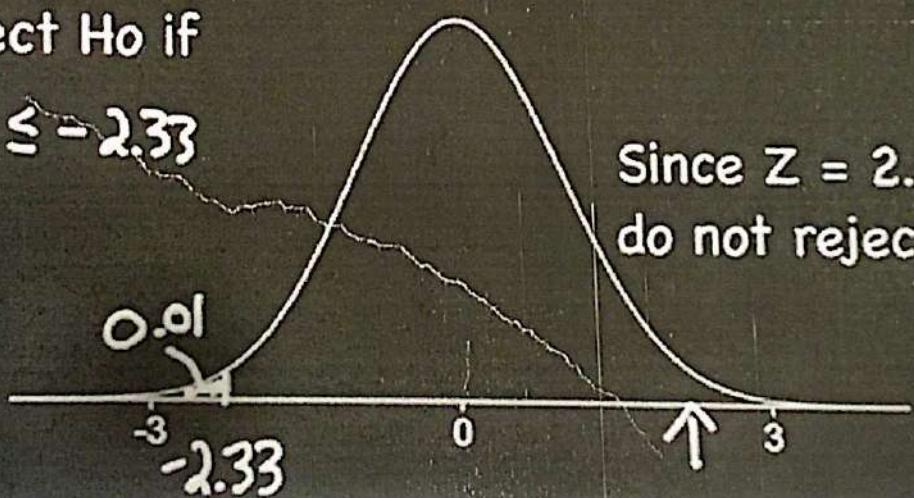
Suppose $H_0: \mu = \mu_0$, $H_a: \mu < \mu_0$, $\alpha = 0.01$

$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}} = 2.6.$$

Reject H_0 if

$$Z \leq -2.33$$

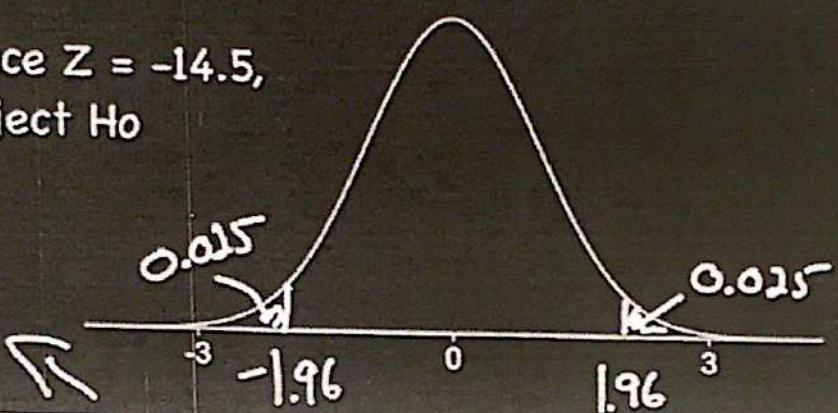
Since $Z = 2.6$,
do not reject H_0



Suppose $H_0: \mu = \mu_0$, $H_a: \mu \neq \mu_0$, $\alpha = 0.05$

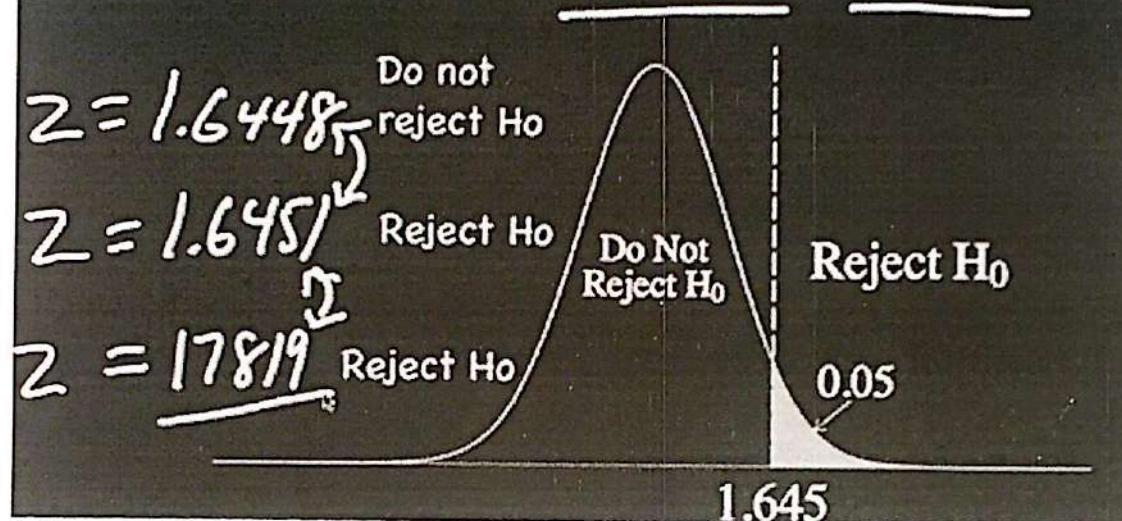
$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}} = -14.5.$$

Since $Z = -14.5$,
Reject H_0



But in some ways the rejection region approach is a little silly.

Suppose $H_0: \mu = \mu_0$, $H_a: \mu > \mu_0$, $\alpha = 0.05$



Many people prefer the *p-value* approach.

The *p-value* is a measure of the strength of the evidence against the null hypothesis.

P Value Approach

Z Tests for the Population Mean μ

The p-value

The definition:

The *p*-value is the probability of getting the observed value of the test statistic, or a value with even greater evidence against H_0 , if the null hypothesis is true.

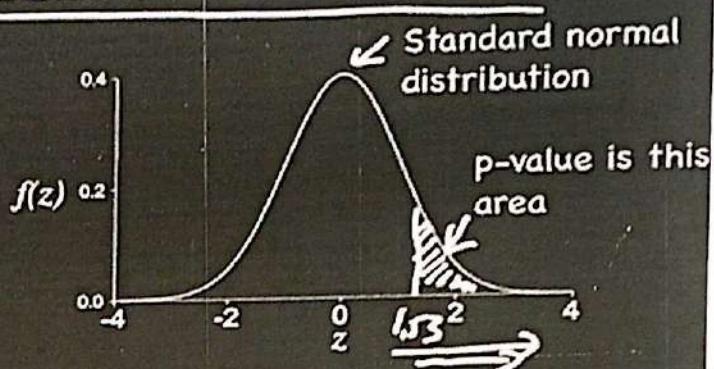
Example:

Distribution of Z if H_0 is true:

$$\begin{aligned} \rightarrow H_0: \mu &= \mu_0 \\ H_a: \mu &> \mu_0 \\ Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} & \end{aligned}$$

Suppose

$$z = 1.53$$



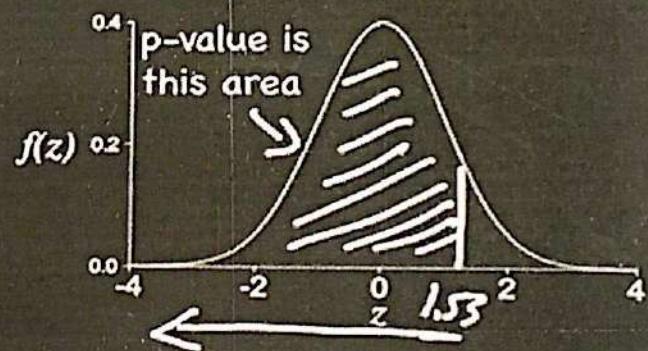
$$\begin{aligned} \text{p-value} &= P(Z \geq 1.53) \\ &= 0.063 \end{aligned}$$

Example:

Distribution of Z if H_0 is true:

$$\begin{aligned} H_0: \mu &= \mu_0 \\ H_a: \mu &< \mu_0 \\ Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} & \end{aligned}$$

$$z = 1.53$$



$$\begin{aligned} \text{p-value} &= P(Z \leq 1.53) \\ &= 0.937 \end{aligned}$$

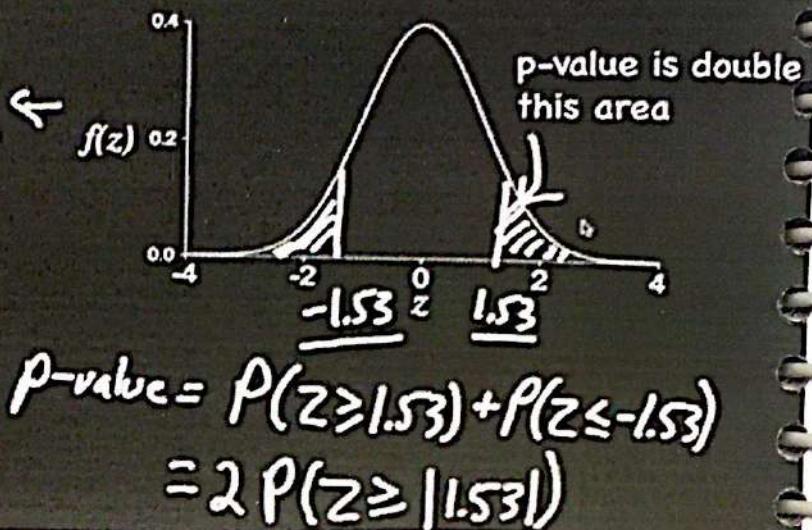
Example: Distribution of Z if H_0 is true:

$$H_0: \mu = \mu_0$$

$$H_a: \mu \neq \mu_0 \leftarrow$$

$$Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}}$$

$$z = 1.53$$



The smaller the p -value, the greater the evidence against the null hypothesis.

If we have a given significance level α , then:

Reject H_0 if $p\text{-value} \leq \alpha$

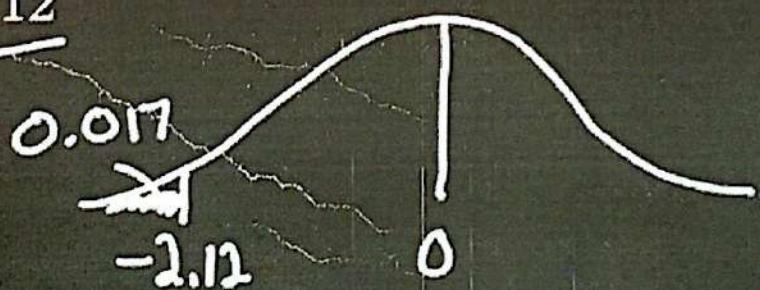
Cutoff level
for significance

(If $p\text{-value} \leq \alpha$, the evidence against H_0 is significant at the α level of significance.)

Example:

Suppose $H_0: \mu = 10$, $H_a: \mu \neq 10$, $\alpha = 0.05$,

$$Z = -2.12$$



$$p\text{-value} = 2 \times 0.017 = 0.034$$

Here we reject our null hypothesis because our p value is less than alpha (level of significance) .

If we do not have a given significance level, then it is not as cut-and-dried.

Example of Z test .

A producer of cereal is producing boxes of cereal with a stated weight of 750 grams.

It is known from a large body of past experience that the standard deviation of the weights in their filling process is approximately 16 grams.

In order to ensure that not many boxes are underfilled, the producer sets the mean fill amount at 780 grams.

As part of the quality control process, they periodically draw a random sample of 25 boxes, measure the weights, and test the null hypothesis that the mean fill amount is 780.

$$H_0: \mu = 780 \quad H_a: \mu \neq 780 \quad \alpha = 0.05$$

For a sample of 25 boxes, $\bar{X} = 776$.

Does this give strong evidence the true mean weight differs from 780 grams?

$$H_0: \mu = 780, \quad H_a: \mu \neq 780, \quad \alpha = 0.05$$

$$\bar{X} = 776, \quad \sigma = 16, \quad n = 25.$$

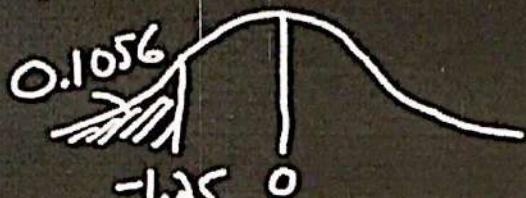
$$Z = \frac{\bar{X} - \mu_0}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}}$$



$$H_0: \mu = 780, H_a: \mu \neq 780, \alpha = 0.05$$

$$\bar{X} = 776, \sigma = 16, n = 25.$$

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} = \frac{776 - 780}{16/\sqrt{25}} = -1.25$$



The evidence against H_0 is not significant at the 0.05 level of significance

$$p\text{-value} = 2 \times 0.1056 = \underline{0.211}$$

There is little or no evidence ($p\text{-value} = 0.211$) that the true mean weight of the filling process differs from 780 grams.

Note: This does not imply there is strong evidence the mean is 780!

One sided Test vs Two sided Test

A company claims that their product contains no more than 2 grams of saturated fat on average.

You intend to test whether there is strong evidence the mean saturated fat content is greater than their claim.

$$H_0: \mu = 2 \quad H_a: \mu > 2$$

A production process creates silicon wafers. The desired thickness is $725 \mu\text{m}$.

You intend to test whether there is strong evidence the mean thickness differs from this amount.

$$H_0: \mu = 725 \quad H_a: \mu \neq 725$$

The choice between a one-sided or two-sided alternative hypothesis can be subject to debate.

It should not be based on the current sample's data!

We might strongly suspect that a newly developed drug decreases blood pressure.

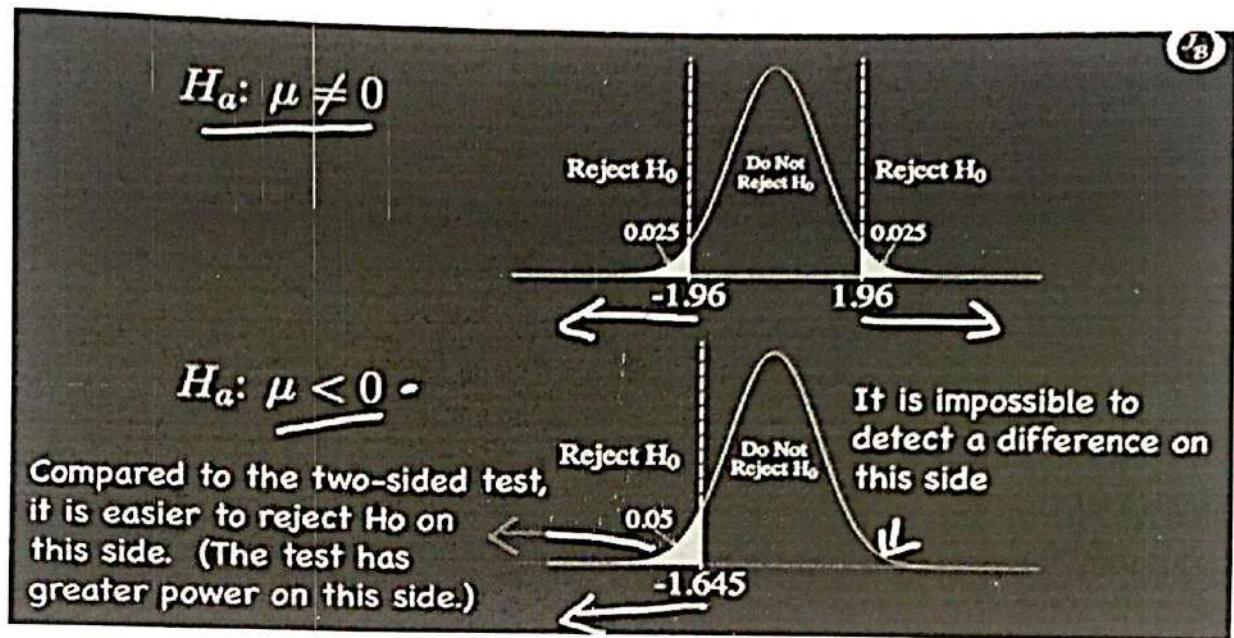
We may wish to test: $H_0: \mu = 0$ Population mean change
in blood pressure

But which is the appropriate alternative?

$$H_a: \mu \neq 0$$

$$H_a: \mu < 0$$

Suppose we are using a Z test at $\alpha = 0.05$.



If we choose the alternative hypothesis based on the direction observed in the sample, then the reported p -value will be half of what it should be.

(We will sometimes be reporting results as significant when they are not.)

Results different from what one expects can often be the most interesting results.

In two-sided tests we are still almost always interested in the *direction* of the difference.

Example: You are investigating possible differences in five-year survival rates for two types of treatment.

$$H_0: p_1 = p_2, H_a: \underline{p_1 \neq p_2}$$

Five-year survival
rate for Treatment 1

Five-year survival
rate for Treatment 2

In tests with two-sided alternatives, the null hypothesis is almost always wrong!

If in the end we simply reject it, without giving an indication of the direction of the difference, we just wasted our time.

In the end, the choice between one or two-sided alternative hypothesis does not need to matter much:

Report the *p*-value (usually two-sided), and let the reader make up their own mind.

Type 1 Errors and Type 2 Errors

Hypothesis Testing

*Type I errors, Type II errors, and
the Power of the Test*

A Type I error is rejecting H_0 when, in reality, it is true.

A Type II error is failing to reject H_0 when, in reality, it is false.

Suppose we test:

$$H_0: \mu = 10$$

$$H_a: \mu > 10$$

and we reject H_0 at $\alpha = 0.05$.

One of two things occurred:

- The null hypothesis is false and we made the correct decision.
- The null hypothesis is true, and we made a Type I error.

Suppose we test:

$$H_0: \mu = 10$$

$$H_a: \mu > 10$$

and we do not reject H_0 at $\alpha = 0.05$.

One of two things occurred:

- ▶ The null hypothesis is true and we made the correct decision.
- ▶ The null hypothesis is false, and we made a Type II error.

Possible outcomes of a hypothesis test:

		Underlying reality	
		H_0 is false	H_0 is true
Conclusion from test	Reject H_0	Correct decision	→ Type I error
	Do not reject H_0	→ Type II error	Correct decision



Known

Consider a criminal trial. We test the hypotheses:

- H_0 : The defendant did not commit the crime.
- H_a : The defendant committed the crime.

Type I error: Convicting a person who, in reality, did not commit the crime.

Type II error: Acquitting a person who, in reality, committed the crime.

Very Imp :-

- Reducing probability of Type I and Type II errors
- Probability of Type I error is zero if we increase sample size
- Probability of Type II error can be reduced by taking larger sample size

The probability of a Type I error, given H_0 is true, is called the *significance level* of the test (α).

$$P(\text{Type I error} | H_0 \text{ is true}) = \underline{\alpha}$$

The probability of a Type II error is represented by β .

The value of β depends on a number of factors, including the choice of α , the sample size, and the true value of the parameter.

The power of a test is the probability of rejecting the null hypothesis, given it is false.

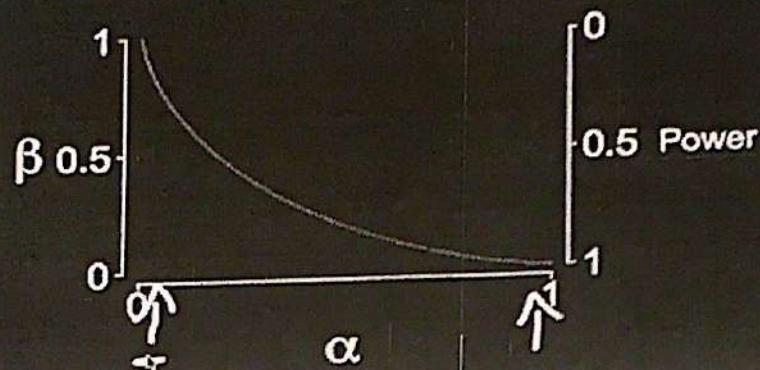
$$\text{Power} = 1 - P(\text{Type II error}) = 1 - \beta$$

The power of a test depends on several factors, including the choice of α , the sample size, and the true value of the parameter.

If we choose a very small value of α , we will be making it very difficult to reject the null hypothesis. (Type II errors will be common.)

If we choose a larger value of α , Type II errors will be less common.

The relationship between α and β for a test of $H_0: \mu = 0$:



Statistical Significance vs Practical Significance .

Hypothesis testing tests for *statistical* significance.

Statistical significance means the effect observed in the sample was unlikely to have occurred due to chance alone.

(It would be very unlikely to see what was observed in the sample if the null hypothesis is true.)

Suppose a call centre claims their average wait time is 30 seconds. We decide to test:

$$H_0: \mu = 30$$

$$H_a: \mu > 30$$

We find $\bar{X} = 30.6$, and a *p*-value of 0.002.

Statistical significance is strongly related to sample size:

If the sample size is large enough, even tiny differences from the hypothesized value will be found statistically significant.

Relationship b/w Hypothesis Test and Confidence Intervals .

The Relationship Between
Hypothesis Tests and
Confidence Intervals

Suppose we find a $(1 - \alpha)100\%$ confidence interval for μ using:

(A two-sided confidence interval) $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

And we wish to carry out a test of:

$$H_0: \mu = \mu_0 \leftarrow$$

$$H_a: \mu \neq \mu_0 \leftarrow$$

with a significance level of $\underline{\alpha}$.

$$Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \leftarrow$$

The confidence interval will be made up of all values of μ_0 for which we would not reject the null hypothesis.

$$\alpha = 0.10$$

Suppose a 90% confidence interval for μ is found to be $(13.1, 26.4)$. The hypothesized value of 25 falls within the 90% interval

What would be the results of a test of $H_0: \mu = 25$ against $H_a: \mu \neq 25$ at $\alpha = 0.10$?

The null hypothesis would not be rejected at a significance level of 0.10

The p-value of the test would be greater than 0.10

$$\alpha = 0.05$$

Suppose a 95% confidence interval for μ is found to be $(48.2, 87.6)$. The hypothesized value of 40 falls outside the 95% interval

What would be the results of a test of $H_0: \mu = 40$ against $H_a: \mu \neq 40$ at $\alpha = 0.05$?

The null hypothesis would be rejected at a significance level of 0.05

The p-value of the test would be less than 0.05

Why is this the case?

We would reject H_0 at $\alpha = 0.05$ if:

The Z test statistic

$$\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \leq -1.96 \text{ or } \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \geq 1.96$$

Isolating μ_0 , we would reject H_0 if:

$$\underline{\mu_0} \geq \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}} \text{ or } \overline{\mu_0} \leq \bar{X} - 1.96 \frac{\sigma}{\sqrt{n}}$$

The upper bound of the
95% confidence interval

The lower bound of the
95% confidence interval

If the hypothesized value falls outside the confidence interval we will reject the NULL hypothesis . If it falls within confidence interval then we will not reject NULL Hypothesis .

Correlation Analysis

$x, y \rightarrow$ degree of relationship

↓
coefficient of correlation
(ρ)

Types of Correlation :-

i) +ve correlation

$x \uparrow \Rightarrow y \uparrow$ or $x \downarrow \Rightarrow y \downarrow$

ii) -ve correlation

$x \downarrow \Rightarrow y \downarrow$ or $x \uparrow \Rightarrow y \downarrow$

iii) no correlation

$x \uparrow \Rightarrow y \downarrow$ or $x \downarrow \Rightarrow y \uparrow$ $\rho = 0$

iv) perfect correlation $\rho = 1$ (perfect +ve)

$\rho = -1$ (perfect -ve)

Methods to find correlation :-

i) coefficient of correlation (ρ)

ii) rank correlation

NOTE :- $|1 < \rho < 0|$

coefficient of correlation :-

(ρ)

$$(1) \rho = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$$

$$(2) \rho = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}}$$

where $\text{cov}(x, y) = \frac{1}{n} \sum (x - \bar{x})(y - \bar{y})$

Measures how two variables vary together.

$$\text{var}(x) = \frac{\sum (x - \bar{x})^2}{n}$$

$$\text{var}(y) = \frac{\sum (y - \bar{y})^2}{n}$$

Dr. D. N. Patel

Regression Analysis :-

$$y - \bar{y} = \rho \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$x - \bar{x} = \rho \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$m_1 = \rho \frac{\sigma_y}{\sigma_x}, \quad m_2 = \rho \frac{\sigma_x}{\sigma_y}$$

- i) Regression line always intersect at point (\bar{x}, \bar{y})
- ii) If two regression lines are given and θ be the acute angle b/w them then

$$\tan \theta = \left| \frac{1 - \rho^2}{\rho} \right| \sqrt{\frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}}$$

Note:-
 \rightarrow If $\rho = 1 \Rightarrow \tan \theta = 0 \Rightarrow$ two lines coincide

\rightarrow If $\rho = 0 \Rightarrow \tan \theta = \infty \Rightarrow \theta = 90^\circ$
 When there is no correlation b/w the lines then they are perpendicular.

- iii) coefficient of correlation is Geometric mean of Regression Coefficients.

$$\rho = \sqrt{m_{xy} \times m_{yx}}$$

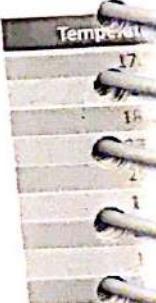
Correlation

Many times, we have
Such as Temperature

Temp

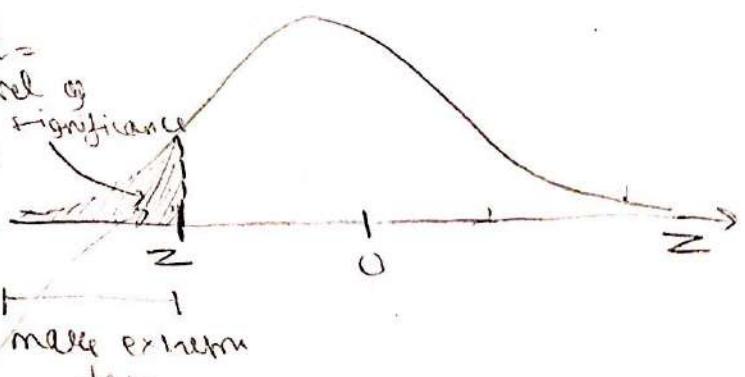
Correlation

To see how the

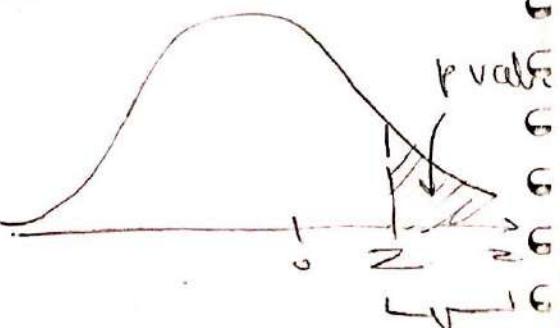


p-value :- probability of obtaining a sample "more extreme" than the one observed in your data, assuming null hypothesis is true.

left tail test :-



right tail test :-

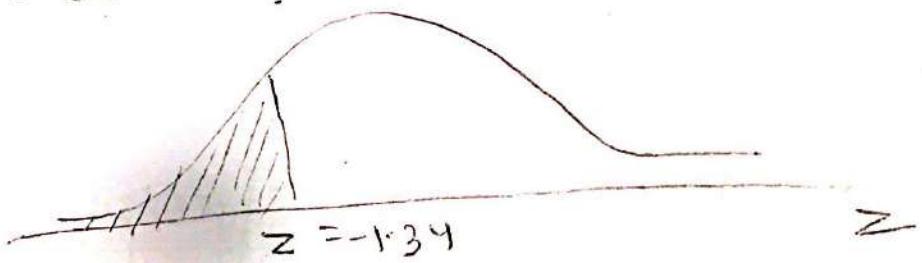


Eg. - left tail test :-

$$H_0 = \mu \geq 0.15$$

$$H_1 = \mu < 0.15$$

From data: $z = -1.34$



$$P(z < -1.34) = 0.0901$$

$$p = 0.0901$$

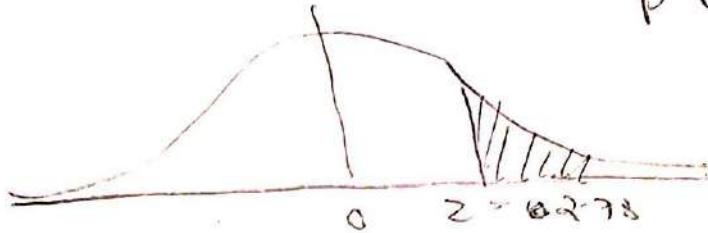
$\rightarrow z$ at value -1.34 , area 0.0901

area under the curve left side up to -1.34 is the p value

right tail test eg :-

$$H_0: \mu \leq 0.43$$

$$H_1: \mu > 0.43$$



given data
 $z = 2.78$

$$P(z > 2.78)$$

$$= (1 - 0.9971)$$

$$= 0.0027$$

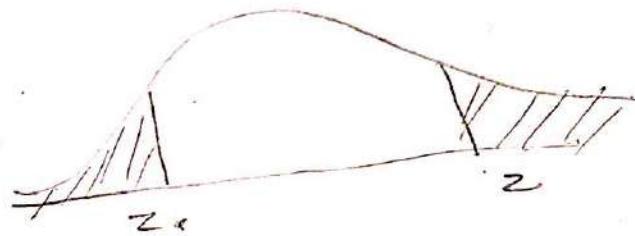
$$p\text{ value} = 0.0027$$

Two tail test:-

$$H_0: \mu = 0$$

$$H_1: \mu \neq 0$$

→ 2 tail
test



If we have two means then,

$$H_0: \mu_1 = \mu_2 = 0$$

$$H_1: \mu_1 - \mu_2 \neq 0$$

→ 2 tail test

Now Here Test statistic:-

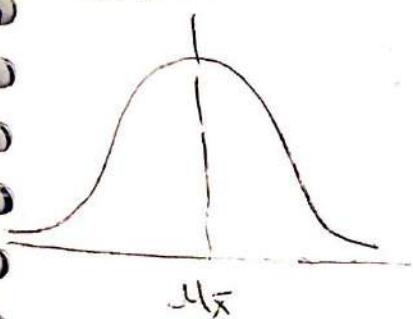
$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

z score
involving
two means

Z statistic vs T statistic

Z statistic

Normal.

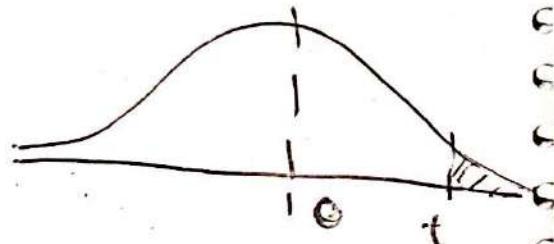


T statistic

T distribution.

$$= \frac{\bar{x} - M_x}{\frac{s}{\sqrt{n}}}$$

[if s is small]



$$= Z = \frac{\bar{x} - M_x}{\frac{s}{\sqrt{n}}} \quad [\text{if } n > 30]$$

Normal distribution more

Type - I error / Type II error :-

Q) School district A states that its high school have 85% passage rates on high school exam? A new school was recently opened in the district, and it was found that sample of 150 students had passage rate of 83%

with standard deviation $\sigma = 1$, does school
bus new school have passage rate than
rest of school district H_1

Null hypothesis:

$$H_0; \mu = 85\%$$

Alternative hypothesis:

$$H_1; \mu \neq 85\%$$

$$\alpha = 0.05 = 5\%$$



→ If statistics we calculate (for example
=) is within 95% of range, we will
conclude that what we are testing (usually
the mean) is right where we expect
it to be, so we will obtain the
null hypothesis

→ If statistics we calculate is outside of
that range, we will conclude that what
we are testing is not where we expect it to
be, we reject H_0 & accept H_1 .

→ we haven't measured the entire school we only measured sample of students. The decision we make may or may not accurately reflect reality.

Outcome 1

we reject Null Hypothesis when in reality it is false → GOOD

Outcome 2

we reject Null hypothesis when in reality, it is true → Type I error.

Outcome 3

we retain Null Hypothesis when in reality, it is false → Type II error,

Outcome 4 :-

we retain Null hypothesis when in reality, it is true → (GOOD).

Small sample hypothesis test :-

If mean emission of all engines of a new design needs to follow 20 ppm of design to meet new emission requirements. = 10 engines are manufactured for testing purpose & emission level of each is determined. Emission data is :-

$H_0: \mu \geq 20$	15.6	16.2	22.5	20.5	16.4	19.4	16.6
$H_1: \mu < 20$	17.9	12.7	13.9				
$H_0: \mu < 20$							

Does the data supply sufficient evidence to conclude that this type of engine meets new standard? Assume we are willing to make a type - I error with $\alpha = 0.01$.

$\bar{x} = 17.17$ $H_0: \mu \geq 20$

~~Given~~ $s = 2.98$ $H_1: \mu < 20$

$H_0: \mu \leq 20$
 $H_1: \mu > 20$

$$\left. \begin{array}{l} H_0: \mu \geq 20 \text{ ppm} \\ H_1: \mu < 20 \text{ ppm} \end{array} \right\}$$

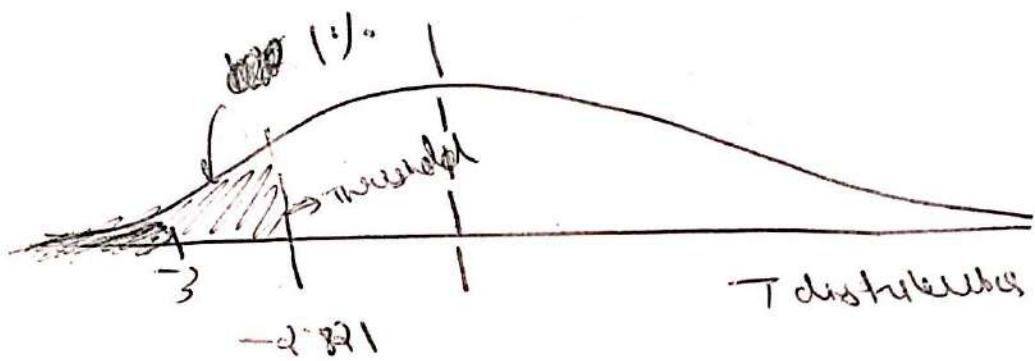
Assume H_0 is true.

We will reject H_0 if $P(\bar{x} = 17.17 \mid H_0 \text{ true}) < 1\%$

We will follow T statistic.

$$\text{as } n < 30$$

$$t = \frac{17.18 - 20}{\sqrt{\frac{2.98}{10}}} = -3.00$$



असा असा ($p = 0.01$)

Hence $\frac{p = 1.1^{\circ} \text{ के लिए } T \text{ distribution}}{\text{graph देता है } -2.821 \text{ का असा}}$

गौणि calculate \bar{x} और t तो $t = -3$

$t = -3$ से उमंगवाया कलन
इस प्रायत्तिकी ने नहीं थी असा -3
 -2.8 के लिए असा तो, हमें हमें
reject null hypothesis &
 $H_0 \rightarrow \checkmark$

Answers

$$n = \frac{1.96^2}{0.03^2} \times 0.5(1-0.5)$$

$$n = 106.8$$

Hypothesis test :-

→ All hypothesis follow a basic logic:

① you made an assumption.

② If your data contradicts the assumption

then you conclude that assumption must be wrong.

→ It is the scientific method for decision making. It is widely used procedure to test

variety of claims.

Eg:- manufacturer claims that average bottle contains 200 mL of water. Now you come and tested 10 bottles and average comes out to be 170 mL. Will you reject? Probably yes? Suppose average comes out to be 205 mL, probably you won't have rejected the ~~opposite~~ claim.

Hypothesis testing allows us to tell that till where we can reject ~~our~~ our claim. You can do that by your gut feeling also, but hypothesis testing is the scientific & statistical way of doing that.

NOTE :-

Any distribution can be turned to standard normal distribution by doing this:

$$\bar{x} \sim \text{Normal}(\mu, \frac{\sigma}{\sqrt{n}}) \Rightarrow \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim \text{Normal}(0, 1)$$

→ Generally we will be finding T statistics in our problems since we will be having sample data.

$$t_{\text{value}} = \pm |T.\text{INV}(\alpha/2, n-1)|$$

↳ excel formula

NOTE :-

null hypothesis cannot have

\neq $>$ $<$

It can only have:-

$=$	\geq	\leq
↓	↓	↓
two tailed test	one tail LHS	one tail, RHS

we wish to test a claim that average age of men MBA students across various MBA programs in US is greater than 28 years. For this we collect data on average age of men MBA students across sample of 40 MBA programs in US.

Step 1 formulate hypothesis

null hypothesis $H_0: \mu > 28$ X as H_0 cannot have $<$
 $H_A: \mu \leq 28$

hence

$$H_0: \mu \leq 28 \quad \checkmark$$

$$H_A: \mu > 28$$

Step 2

calculate t statistics

$$t_{\text{stat}} = \frac{\bar{x} - \mu}{s/\sqrt{n}} = 0.22$$

Step 3 :- cutoff value for t stat

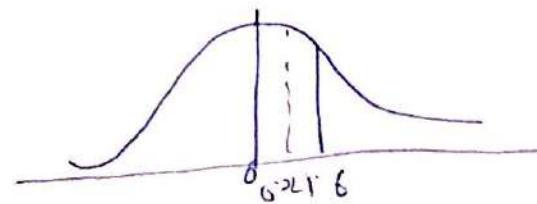
(single tail test on RHS)

$$= |T. \text{ INV } (2, 40-1) | = 1.6$$

$\downarrow 0.5$

Step 4 :-

do not reject null hypothesis.



in 2

Hypothesis test for population proportion

> where we interest in population proportion rather than population mean.

e.g. University introduces a new lunch facility on campus on a trial basis. University operates lunch facility for few months and then decides to survey student body. Based on survey, university would make this facility a permanent fixture or do away with it. Specifically, if more than 70% students approve then facility would be made permanent else it would shut down.

→ university conducts a survey with 750 random selected students on campus and finds that 510 of these students or (68% of sampled student) approves new facility & remaining 240 students or 32% do not approve.

→ Based on criteria should university should make this permanent or not?

Ans: Step 1 formulate hypothesis

$$H_0: p \geq 0.70$$

$$H_a: p_a < 0.70$$

Step 2 :-

Calculate Z statistics

$$Z \text{ statistic} = \frac{\overbrace{\bar{p} - p_{\text{hypothesis}}^{\uparrow}}^{0.68}}{\sqrt{\frac{p_{\text{hyp}}(1-p_{\text{hyp}})}{n}}} \rightarrow 0.7$$

$\sim 7.50 \text{ studies}$

- * here in this case we always use Z statistic not t statistic in population proportion case

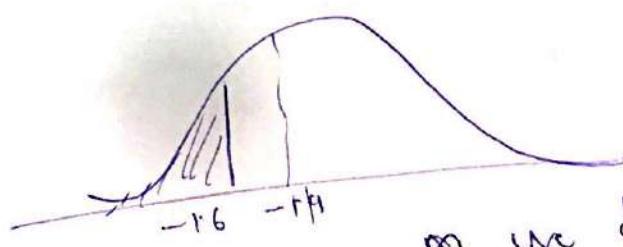
$$Z_{\text{stat}} = \frac{\overbrace{\bar{p} - p^{\uparrow}}^{0.68}}{\sqrt{\frac{p(1-p)}{n}}} = \cancel{0.68} - 1.1952$$

Step 3 :- cut off for Z statistic.

$$z_{\text{cut off}} = -\left\{ \text{NORM. INV}\left(d, \underbrace{0}_{\text{mean}}, \underbrace{1}_{\sigma^2}\right) \right\}^{0.05}$$

$$= -1.6$$

rejection region should be on LHS.



→ hence we do not reject Null hypothesis

NOTE

Z test for comparing two proportions :

$$Z = \frac{(p_1 - p_2) - \bullet \cdot p_{\text{hyp}}}{\sqrt{\frac{p(1-p)}{n_1 + n_2}}}$$

Difference In means hypothesis test :-

> when we are comparing population means across two different populations.

→ we have athletes data , and claim was that average height of men aged 18 years to 45 years across the world was 173cm.

Hence hypothesis test can be used to test this claim for population of men athletes at olympics.

→ Step 1 :-

Formulate hypothesis

$$H_0 : \mu_{\text{height}} > 173 \text{ cm}$$

$$H_a : \mu_{\text{height}} \leq 173 \text{ cm} \quad \times$$

but we cannot use that inequality here so we flip it

$$H_0 : \mu_{\text{height}} \leq 173 \text{ cm}$$

$$H_a : \mu_{\text{height}} > 173 \text{ cm}$$

Free rejection region is on RHS

Step 2 :- calculate + statistic

$$t_{\text{star}} = \frac{\bar{x} - \mu_{\text{height}}}{S/\sqrt{n}}$$

Step 3 :- cutoff for t statistic

$$t_{cutoff} = + |T.INV(1, n-1)| \\ = 3.03$$

Step 4 :- reject NULL Hypothesis.

→ final conclusion based on our data evidence we cannot claim reject the claim that male athletes at olympics are taller than 173 cm.

→ same study also claimed that the difference in average heights across men & women from various walks of life around world was 12.5 cm.

Now, to find
 $(\bar{M}_{men} - \bar{M}_{women})$

Hypothesis test in difference in means

Claim to be tested

Difference between the population mean height of men and women and athletes is 12.5 cm.

Step 1 :- Formulate hypothesis

$$H_0: \bar{M}_{men} - \bar{M}_{women} = 12.5 \text{ cm}$$

$$H_a: \bar{M}_{men} - \bar{M}_{women} \neq 12.5 \text{ cm}$$

two tailed test

Step 2 :- calculate t-statistic

* For calculation of t-statistic we need to see the nature of variation or variance in two populations.

* Is variance more similar or more dissimilar

* Need to assume either equal variance or unequal variance across two populations.

→ In most applications conclusions do not change based on 'equal' or 'unequal' variance assumption.

Equal variance assumption

$$t_{std} = \frac{\bar{x}_1 - \bar{x}_2 - \mu}{\sqrt{\left(\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \right) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

we can do this using excel in excel we can use t-test when we are doing diff. in mean. this formula is not applicable for single mean.

Get data → data analysis →

T-test assuming equal variance = OK

Variable 1 range = men data range B1:B16

Variable 2 range = C1:C16
Mean difference = 12.5

Unequal variance assumption

$$t_{std} = \frac{\bar{x}_1 - \bar{x}_2 - \mu}{\sqrt{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)}}$$

we will get all results ✓

Paired T test :-

→ Used to test claims around difference in two population means
↳ * what is difference in paired test & normal difference in means hypothesis test.
we have conduct the training and we have to check that training is effective or not.
we have taken 50 students and we need to check that atleast on an average after score should be 10 points greater than before score?

Claim to be tested :-

→ Training was effective. That is on an average after scores are atleast 10 points greater than before score.

Sample Data :-

<u>Name</u>	<u>Before Score</u>	<u>After Score</u>
-	-	-
-	-	-
-	-	-
-	-	-

we called it paired because before and after data is of the similar student. They are paired with each other. last week example men or women ~~at~~ was not paired example.

Step 1 :-

$$H_0: \mu_{\text{before}} - \mu_{\text{after}} \geq 10$$

$$H_a: \mu_{\text{before}} - \mu_{\text{after}} < 10$$

$$n_{\text{before}} = 50 \quad n_{\text{after}} = 50$$

Step 2

* use data analysis tool box of excel and select

paired t-test.

NOTE :- you may have false conclusions
generated in using paired
either equal or unequal

→ we file our information
of paired T test and you

T test VS
variance T test.

in data analysis
will get result.

$$t_{\text{stat}} = -1.82$$

$$t_{\text{critical}} = -1.6$$

Conclusion :- I Reject Null Hypothesis.

Applications :-

Difference in Means T test :-

Examples :-

- ① Testing whether customer satisfaction across service companies in US is greater than Europe.
- ② food and drug company want to compare difference between two disease.
- ③ Testing whether average productivity of woman employees is greater than that for men employees.

Three types of difference in means T test :-

- ① Paired T test :- used when there is pairing in the data.
- ② ^{T test} difference in means 'assuming equal variance'
- ③ T test difference in means 'assuming unequal variance'.

To calculate pvalue,
we have formula

$$pvalue = 2 * T.DIST \left(- |t\text{-statistic}|, \text{residual df, } \text{TRUE} \right)$$

↙ formula to find
pval.

Note,

Hypothesis Testing in regression context.

→ t-cutoff approach { already seen}

→ pvalue approach.

→ confidence interval approach

confidence interval approach :-

whenever we find the regression output we will see that the upper end and lower end of that coefficient.

Ex:-

$$\begin{aligned} \text{Step 1} :& H_0: \beta = 500 \\ & H_a: \beta \neq 500 \end{aligned}$$

$$\begin{aligned} \text{Step 2} :& \text{ consider 95% confidence interval} \\ & \text{for } \beta \\ & = [212.6, 1084.6] \end{aligned}$$

Conclusion :-

since 500 falls in the confidence interval hence we cannot reject NULL hypothesis.

When we do regression test,
we get p values ^{automatically} what are they?

→ Excel is doing hypothesis test for every
coefficients and giving us p values
These hypothesis are of kind as

Model :-

$$\rightarrow \text{Sales} = \beta_0 + \beta_1 \text{Adver} + \beta_2 \text{Bram}$$

Automatic hypothesis testing by excel.

$$H_0 : \beta_0 = 0$$

$$H_a : \beta_0 \neq 0$$

$$H_0 : \beta_1 = 0$$

$$H_a : \beta_1 \neq 0$$

$$H_0 : \beta_2 = 0$$

$$H_a : \beta_2 \neq 0$$

→ Automatic p values for every coefficient
comes by doing the test as whether
null hypothesis is whether coefficient is zero or
not.

→ If p value comes out to be low
we reject null hypothesis and that
is good, that variable can help to
explain variation in y variable.

Low p values are good. → as we
reject null hypothesis

→ In other words that variable is
important in regression.

R^2 : (goodness of fit)

* overall variation in Y variable : Total sum of squares

Explained variation in Y variable : Regression sum of squares

unexplained variation in Y variable : Residual sum of squares

Total sum of squares = Regression SS + Residual SS

$$R\text{-square} = \frac{\text{Regression SS}}{\text{Total SS}}$$

Adjusted R square decrease if we add unnecessary explanatory variable. , by value of R^2 is never decreased.

Dummy Variables :-

→ we change categorical variable to numeric variable by applying statistics. When using regression analysis.

→ We use If and else statement for that.

$$\text{Ex:- } \text{If}(B2 = "A", 1, 0)$$

NOTE:-

If focus is more on predictions then, low ~~R^2~~ is problem

If focus is on understanding relation between X and Y variable then low R^2 may not be problem.

Mean centering variables in Regression model

$$\text{weight} = \beta_0 + \beta_1 \text{Male} + \beta_2 \text{Height}$$

(this is categorical so we make it dummy)

$$\begin{aligned}\text{Male} &= 1 \\ \text{Female} &= 0\end{aligned}$$

<u>Gender</u>	<u>Height</u>	<u>weight</u>	<u>Gender</u> (male)
M	:	:	= $\Sigma (\text{P}_i = "m", 1, 0)$
F	:	:	:
:	:	:	:
:	:	:	:

make regression model.

	<u>coefficient</u>	<u>p value</u>
Intercept	-101	6E-100
Male	5.5	2.83E-20
Height	0.966	2.83E-40

→ Now let us interpret each and every coefficients.

Height: every 1cm increase in height of olympian weight increases by 0.97kg

Male: Male olympian than female olympian have 5.5kg more.

Intercept:- Female with zero height -101 is weight of the person, but it does not make sense since height cannot be zero. (does not have managerial interpretation)

→ we can make it meaningful by
mean centering height.

$$\text{male} \quad \begin{array}{l} \text{Height (mean)} \\ = \bar{x} - \text{average (fit)}_{100} \end{array}$$

$$\begin{array}{c} \text{Height} \\ ; \\ ; \\ ; \\ ; \end{array} \quad \begin{array}{c} \text{Weight} \\ ; \\ ; \\ ; \\ ; \end{array}$$

Now, again run regression (male, weight
and height)

$$\begin{array}{ll} \text{new,} & \text{Coeff} \\ \text{Intercept} & 69.66 \\ \text{male} & 5.5 \\ \text{height mean} & 0.96 \end{array} \quad \begin{array}{l} \text{rvalue} \\ 0 \\ 2.83e^{-10} \\ 2.2e^{-20} \end{array} \quad \begin{array}{l} \text{it is same} \\ \} \end{array}$$

→ now, we can interpret intercept as
managerial interpretation :-

→ weight of the female with whose
height is equal to average height in the
entire data is 69.66 kg .

$$\text{and male Olympiad weight} \\ = 69.66 + 5.5$$

Interpretation

→ on an average male Olympian has
weight that is 5.5 kg more as compared to female
Olympian, all other variables kept at same
level.

Interaction effects in Regression :-

NOTE

$$\text{weight} = \beta_0 + \beta_1 \text{Male} + \beta_2 \text{Height}$$

* Effect of height on weight is same across gender. As we interpret as increase in 1cm of height will increase the weight by 0.966 kg, irrespective of gender, since it was held constant.

* No gender difference allowed based on the height.

Q How do we allow impact of 'height' based on gender?

This is where interaction effects come into play.

$$\text{weight} = \beta_0 + \beta_1 \text{Male} + \beta_2 \overline{\text{Height}} + \beta_3 \overline{\text{Height} \times \text{Male}}$$

↳ interaction variable

Now, let us interpret the coefficient of

height.

→ The impact of one cm increase in height on the weight, all other variables held constant.

→ To keep the interaction variable 'at some level', we need to consider female olympians because 1cm increase in height will also increase interaction variable and we have to keep that constant.

here, β_2 is the impact of 1cm increase in height on weight of female olympians.

Now, what about impact of height on weight of male olympians?

Ans: $\beta_3 \rightarrow$ it gives us additional impact on weight for male olympians

Thus, total impact of height on weight for male olympians is $(\beta_2 + \beta_3)$

Note, Transformation of variables in Regression :-
(Improving linearity)

Q Why we need transformation?
→ If relation between y and x variable is not linear.

Variable Transformations

- * Invert one or more variable
- * Take a square of variable
- * Take a square root
- * Log transformation



Log Log Model

~~DDDDDD~~

$$\ln(Y) = \beta_0 + \beta_1 \ln(X_1) + \beta_2 \ln(X_2) + \dots$$

Interpretation,
when natural log of X_1 increases by
1 unit, then natural log of Y increases by
 β_1 units, all other variables are kept at
same level. OR

we read it like this way one percentage
increase in X_1 corresponds to β_1 percentage
increase in Y , all other variables at same
level.

Semi-log model :-

$$\ln(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

↓

Interpretation

one unit increase in X_1 , corresponds to $100\beta_1$ percentage increase in Y , all other variables at same level.

Now,

NOTE :-

Eg:- Annual demand for cocoa in millions pounds over a period of time.

$$\ln(\text{demand}) = \beta_0 + \beta_1 \ln(\text{price}) + \beta_2 \ln(\text{pcinc}) + \beta_3 \text{year}$$

→ Here we have Log-Log as well as

semilog interpretation.

→ β_1 and β_2 is log-log interpretation also known as elasticity interpretation.

→ β_3 is semi log interpretation also known as growth rate interpretation.

as growth rate interpretation

→ It is used for analyzing relationship between one categorical and one numerical variable like total there are 60 people divided equally into 4 diets.

ANOVA

	Water	Coffee
29		28
30		29
31		27
31		30
29		29



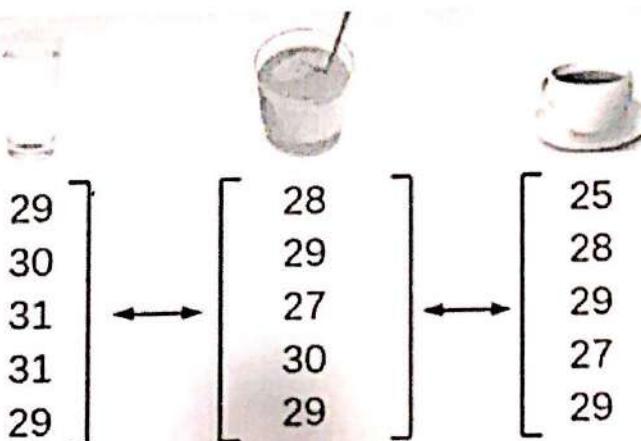
Variation within each group

Diet	weight loss (in kg)
A	16, 17, 17, 17
B	- - -
C	- - -
D	- - -

H_0 = mean weight loss for all diets are same

H_a = At least one is different

First example



Variation between the groups

and the variation between the groups,

Assumption :-

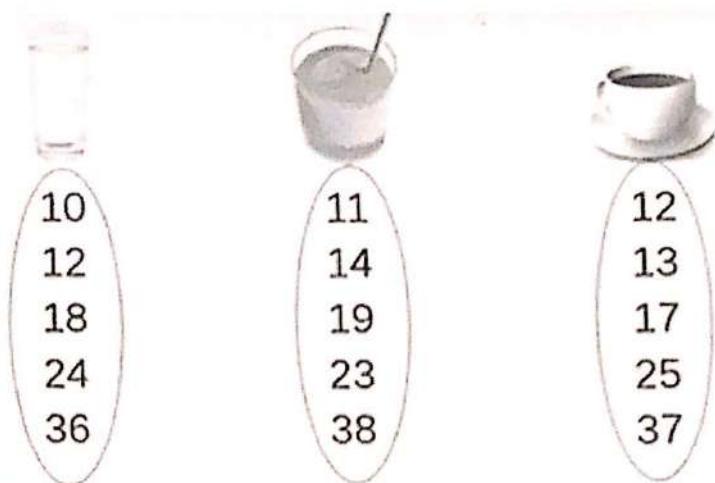
① Independent samples, people on diet A is independent of diet B

② Standard deviation of each group is roughly same

③ each group has large sample size. as groups are approx

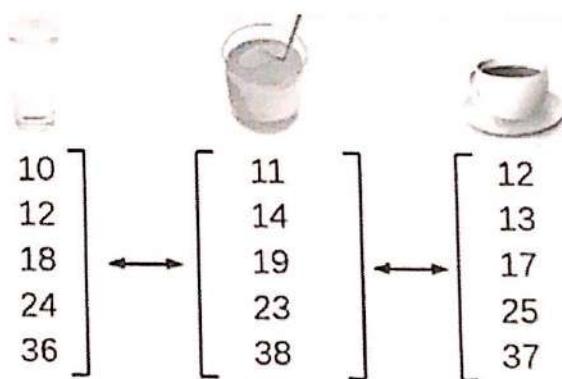
normal distribution.

If n is small the we can use Kruskal test or bootstrap approach.



There's a lot of variation in each group...

to make it easier to see the patterns.



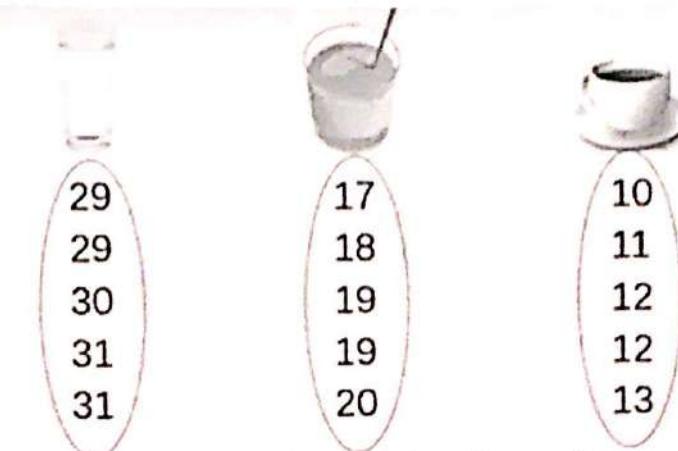
...but each group looks pretty much the same.

But all the groups look pretty much alike;

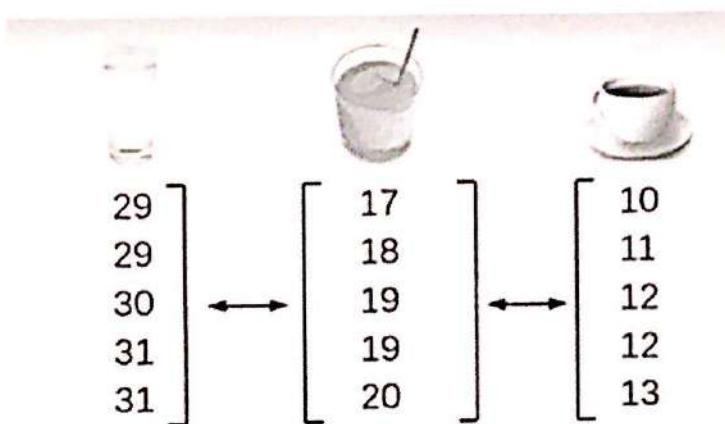
- > help(aov)
 - > ? aov
 - > boxplot (weight loss ~ Diet)
 - > aov (weight loss ~ Diet)
 - > ANOVA1 <- aov (weight loss ~ Diet)
 - > summary (ANOVA1)
 - $F = 6.113$ $p = 0.0013$
- hence p is small, we
reject Null Hypothesis

- > Tukey HSD (ANOVA) // to check if difference exists
 - > p-hat (Tukey HSD)
- जोले के mean के
तर्क, वहाँ से
group के difference
ज्ञात होते हैं।

Second example



There's not much variation in each group...



...but the groups look very different.

In this case we will reject the NULL Hypothesis.

↪ we assume over null hypothesis

$$\text{to be } H_0 = \mu_1 = \mu_2 \\ H_0:$$

Figure out how much of the total variance comes from:

The variance *between* the groups

The variance *within* the groups

Calculate the ratio:

$$F = \frac{\text{between groups}}{\text{within groups}}$$

The larger the ratio, the more likely it is that the groups have different means (reject H_0).

STRESS DATA (SACRIVI), 2014		
4 AM to MIDNIGHT	MIDNIGHT TO 8 AM	8 AM TO 4 PM
2	5	1
7	6	3
6	3	4
9	5	3
7	4	1
7	6	1
6	5	2
7	4	6
7	5	5
9	5	4
	6	3
	7	4
	6	5

Within Groups Mean Sum of Squares

$$MSS_W = \frac{\sum_{g \in G} (X - \bar{X}_g)^2}{n - k}$$

Between Groups Mean Sum of Squares

$$MSS_B = \frac{\sum_{g \in G} n_g (\bar{X}_g - \bar{X}_G)^2}{k - 1}$$

Test Statistic for Tests Concerning the Differences between the Variances of Two Populations (Normally Distributed Populations)

$$F = \frac{MSS_B}{MSS_W}$$

$$df_B = k - 1$$

$$df_W = n - k$$

STRESS DATA (SACKING, 2016)			
4 AM TO DUSK	TECHNIGHT TO 9 AM	8 AM TO 4 PM	
7	5	1	
7	6	4	
6	3	3	
9	5	1	
3	4	1	
7	6	2	
6	5	6	
7	4	5	
2	5	4	
9	5	3	
	6	4	
	7	5	
	6	5	
	10	13	
	36	13	
\bar{x}_G	7.40	6.15	3.23
\bar{x}_g		5.08	
\bar{x}_G		3	
k			
$H_0: \mu_1 = \mu_2 = \mu_3$ $H_A: \mu_1 \neq \mu_2 \neq \mu_3$			

$MSS_W = \frac{\sum_{g \in G} (x - \bar{x}_g)^2}{n - k}$
 $MSS_B = \frac{\sum_{g \in G} n_g (\bar{x}_g - \bar{x}_G)^2}{k - 1}$
 $F = \frac{MSS_B}{MSS_W}$
 $df_B = k - 1$
 $df_W = n - k$

4 AM TO MIDNIGHT	MIDNIGHT TO 7 AM	8 AM TO 4 PM
2	5	1
7	6	3
6	3	4
9	5	3
3	4	1
7	6	1
6	5	2
7	4	6
3	5	5
9	5	4
	6	3
	7	4
	6	5
n_g	10	13
\bar{x}_g	3.6	
\bar{x}_g	7.40	5.15
\bar{x}_g		3.23
k		3

4 AM TO MIDNIGHT	MIDNIGHT TO 7 AM	8 AM TO 4 PM	$(x_i - \bar{x}_g)^2$	$(x_i - \bar{x}_c)^2$	$(x_i - \bar{x}_w)^2$
2	5	1	0.16	0.02	4.92
7	6	3	0.16	0.72	0.05
6	3	4	1.96	4.64	0.59
9	5	3	2.56	0.02	0.05
3	4	1	0.36	1.33	4.98
7	6	1	0.16	0.72	4.98
6	5	2	1.96	0.02	1.51
7	4	6	0.16	1.33	7.67
3	5	5	0.36	0.02	3.13
-	9	4	2.56	0.02	0.59
	6	3		0.72	0.05
	7	4		3.41	0.59
	6	5		0.72	3.13
n_g	10	13	13	$\Sigma = 10.40$	$\Sigma = 13.69$
\bar{x}_g	3.6				$\Sigma = 32.31$
\bar{x}_g	7.40	5.15			
\bar{x}_g		5.08			
k		3			

$$MSS_w = \frac{\sum_{i=1}^n (X_i - \bar{X}_w)^2}{n-k}$$

$$MSS_c = \frac{\sum_{i=1}^n n_g (X_i - \bar{X}_c)^2}{k-1}$$

$$F = \frac{MSS_c}{MSS_w}$$

$$df_c = k-1$$

$$df_w = n-k$$

$(x_1 - \bar{x}_1)^2$	$(x_2 - \bar{x}_2)^2$	$(x_3 - \bar{x}_3)^2$
0.16	0.02	4.98
0.16	0.72	0.05
1.96	4.64	0.59
2.56	0.02	0.05
0.36	1.33	4.98
0.16	0.72	4.98
-1.96	0.02	1.51
0.16	1.33	7.67
0.36	0.02	3.13
2.56	0.02	0.59
	0.72	0.05
	3.41	0.59
	0.72	3.13
$\Sigma = 10.40$	$\Sigma = 13.69$	$\Sigma = 32.31$
$SS_W = 56.4$		
$n-k = 33$		
$MS_W = \frac{56.4}{33} = 1.71$	\rightarrow This is mean squared within the group	

$$\begin{aligned}
 n_1(\bar{x}_1 - \bar{x}_G)^2 &= 10(7.40 - 5.06)^2 = 53.82 \\
 n_2(\bar{x}_2 - \bar{x}_G)^2 &= 13(5.19 - 5.06)^2 = 0.06 \\
 n_3(\bar{x}_3 - \bar{x}_G)^2 &= 13(3.23 - 5.06)^2 = 46.49 \\
 \Sigma &= 98.37
 \end{aligned}$$

$MS_B = \frac{98.37}{2} = 49.19$

\rightarrow This is mean squared between the group

$F = 6.27$
$F = \frac{49.19}{1.71}$
$F = 28.77$

$$F_{\text{stat}} < F_{\text{critical}}$$

→ fail to reject
NULL Hypothesis

F_{critical}
$F = 49.19$
1.71
$F = 29.77$
$\text{df}_B = 3 - 1 = 2$
$22 \quad \text{df}_W = 35 - 3 = 33$
$F_{\text{stat}} = 3.32$
$F_{\text{stat}} > F_{\text{critical}}$
$\therefore \text{Reject } H_0$

$$F_{\text{stat}} > F_{\text{crit}}$$

→ reject NULL
Hypothesis

How to find F_{critical}

$$F_{\text{crit}} = \frac{\text{df}(\text{between})}{\text{df}(\text{within})}$$

$$= \frac{(K-1)}{(n-k)} = \frac{3-1}{33}$$

(2 and 33)
new set
value in graph

to find

F_{critical} is found

value to be

$$F_{\text{critical}} = 3.32$$

Note:

$$F \uparrow \Rightarrow P \downarrow$$

→ If p-value is less
then we reject
NULL Hypothesis

- ① It works with categorical target variable "success" or "failure"
 - ② It performs only binary splits.
 - ③ CART (Classification & regression tree)
 - uses Gini method to create binary splits.
 - > TAB ← take (Gender, Smoke)
 - > Chi-sq. test (TAB, correct = T)
 - $p\text{ value} = 0.1866$
 - Chi-square method :-
- $\chi^2 = \sum \left[\frac{(f_{0i} - f_{ei})^2}{f_{ei}} \right]$
- If \rightarrow observed frequency \rightarrow Chi-square test (TAB, correct = T)
 $f_{0i} \rightarrow$ expected frequency hence we fail to
 f_{ei} reject null hypothesis.

Working Rule :-

- ① Consider the Null Hypothesis : There is no difference between observed frequency & expected frequency. / X and Y are independent
- H₀
- H_a : X and Y are dependent.
- ② calculate expected frequency f_e corresponding to each cell.
- ③ calculate χ^2 distribution by formulae

$$\chi^2 = \sum \left[\frac{(f_{ei} - f_{0i})^2}{f_{ei}} \right] \text{ & calculate degree of freedom.}$$

Assumptions :-

- ① Groups are independent
 - ② All cell ≥ 1
 - ③ All expected cell counts ≥ 5
- if it is not met then we can use Fisher's test or Bootstrap test
- it means कोई ग्रूप नहीं बनाया जा सकता है।
 कोई ग्रूप नहीं बनाया जा सकता है।

④ See ~~the~~ and find the value of χ^2_{critical} from table.

⑤ If $\chi^2 < \chi^2_{\text{critical}}$ then Null hypothesis is accepted otherwise Null hypothesis is rejected.

NOTE :- Have to find expected frequency,

A/B	B ₁	B ₂	
A ₁	f_{11}	f_{12}	N ₁
A ₂	f_{21}	f_{22}	N ₂
	N ₃	N ₄	N

$f_{11} = \frac{N_1 \times N_3}{N}$
 $f_{12} = \frac{N_1 \times N_4}{N}$
 $f_{21} = \frac{N_2 \times N_3}{N}$
 $f_{22} = \frac{N_2 \times N_4}{N}$

→ Chi square statistic is used for testing relationships between categorical variables. The null hypothesis of chi square is that there is no relationship between categorical variables of the population & they are independent.

→ It is also called 'goodness of fit' because it measures how well the observed distribution of data fits with distribution that is expected if variables are independent.

4) eg freedom :-

Binomial distribution : (n-1)

Poisson distribution : (n-2)

Normal distribution : (n-3)

(MxN)table : (m-1) (n-1)

→ Pearson's chi square test can be used
to find the correlation between two categorical
variables, if sample size is greater than 50.

In sample survey of public opinion :-

Question 1: do you drink

Question 2: Are you in favour of local option on
sale of liquor

Q2:

	Yes	No	Total
Yes	56	31	87
No	18	6	24
Total	74	37	111

can you infer or not that local option on

sale of liquor is dependent on individual drink?

Given that χ^2_{critical} at dof = 5 and level of significance = 5%
 $= 3.841$?

$$f_{11} = \frac{87 \times 74}{111} = 58, f_{12} = \frac{87 \times 37}{111} = 31, f_{21} = \frac{24 \times 74}{111} = 16 \\ f_{22} = \frac{24 \times 37}{111} = 8$$

$$\chi^2 = \frac{(58-56)^2}{56} + \frac{(31-29)^2}{29} + \frac{(16-13)^2}{13} + \frac{(8-6)^2}{6}$$

$$\chi^2 = 0.95$$

Hence null hypothesis is accepted that

Dice is tossed 120 times

No. turned up :	1	2	3	4	5	6
frequency :	30	25	18	10	22	15

Test hypothesis that dice is unbiased ($\chi^2_{0.05, 5} = 11$)

Ans :-

H₀: Dice is unbiased

Calculation of expected frequency :-

$$P(\text{gt}) = \frac{1}{6}$$

$$\text{expected frequency} = 120 \times \frac{1}{6} = 20 \text{ for } g_i = 1, 2, 3, 4, 5, 6$$

$$\chi^2 = \frac{(30-20)^2}{20} + \frac{(25-20)^2}{20} + \frac{(18-20)^2}{20} + \dots \quad \checkmark$$

$$\chi^2 = 12.90$$

Hence Null hypothesis is rejected.

→ Now we can use Chi-square technique on decision trees.

→ Since, Chi-square is algorithm to find out statistical significance between the difference between sub-node and parent node.

→ Higher the Chi-square value higher the statistical significance between sub-node & parent Node.

→ It generate tree called (CHAD)

(Chi-square automatic Interaction Detector)

→ Ex यह variable को target variable को साझा contingency table करता है, जिस variable को chi-square, ANOVA, regression आदि, को variable root node करता है।