

→ I am having a process whose outcomes are binary and we are repeating the process for n no. of times
→ This kind of dist. we say as Binomial dist.

① Binomial Distribution

→ In probability theory and statistics, the binomial distribution with parameters n and p is the discrete probability distribution of the number of successes in a sequence of n independent experiments, each asking a yes-no question, and each with its own Boolean-valued outcome: success (with prob. p) or failure (with prob. $q=1-p$). A single success/failure experiment is also called a Bernoulli trial or Bernoulli experiment, and a sequence of outcomes is called a Bernoulli process; for a single trial, if $n=1$, the binomial distribution is a Bernoulli distribution.

Ex: Tossing a coin {Bernoulli Distribution}

$$P(H) = 0.5 = p$$

↓
0

$$P(T) = 0.5 = p$$

↓
1

$\boxed{n=10}$

Ex: Tossing a coin for 10 times

| 1 st time | 2 nd | 3 rd | 4 th |
|----------------------|-----------------|-----------------|-----------------|
| $P(T) = p$ | p | p | p |
| $P(H) = 1-p$ | $1-p$ | $1-p$ | $1-p$ |

When we combine this

we get: Binomial Distribution.

→ Binomial dist. can combine various Bernoulli dist. in p.



Parameters n = no. of trials, p = prob. of success or failure.

p.e. Long → success or prob. of each trial.

$$q = 1-p$$

$$\text{PMF} \quad P(X) = {}^n C_k p^k (1-p)^{n-k}$$

$X \in \{0, 1, 2, \dots, n\} \rightarrow$ no. of success

④ Mean of binomial distribution

Mean = $n p \rightarrow$ no. of experiments

⑤ Variance And std

$$\text{Variance} = npq$$

$$\text{std} = \sqrt{npq}$$

Poisson Distribution &

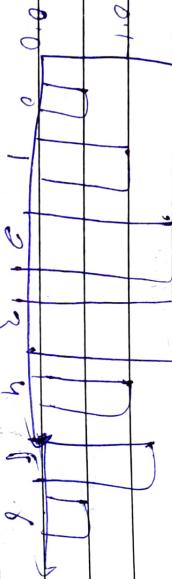
- (i) Discrete Distribution (PMF)
 Describes the number of events occurring in a fixed time interval.

e.g no. of people visiting hospital every hour
 no. of people visiting banks every hour
 no. of people visiting airport every hour
 gobarometer

$\lambda = 3 \rightarrow$ Expected no. of

events to occur at

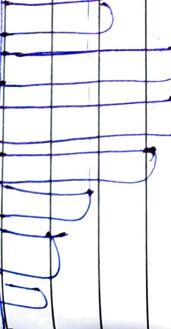
every time interval



$$\lambda = 3$$

PMF

$$P(X=5) = \frac{e^{-\lambda} \lambda^x}{x!}$$



$$\Rightarrow e^{-3} 3^5 = 0.101$$

$\approx 10.1\%$ for happening event
 to occur at least
 once at this time

$$P(X=4) + P(X=5) =$$



Mean And variance of

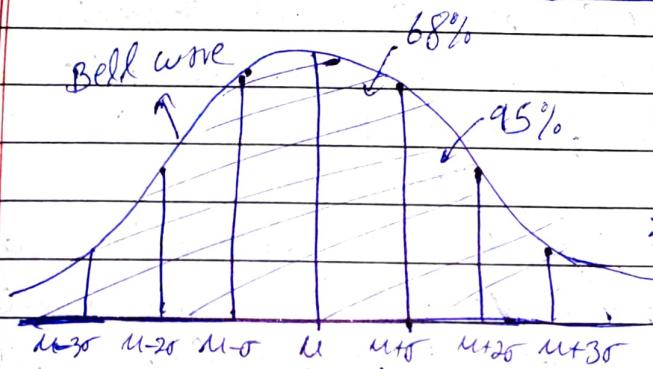
$$\text{mean} \Rightarrow E(X) = \mu = \lambda t.$$

λ = expected no. of events to occur at every time interval.

$$\begin{aligned}\text{variance} &\Rightarrow E(X^2) - \mu^2 \\ &= \lambda t.\end{aligned}$$

t = time interval.

Normal / Gaussian distribution



$X = \{ \dots \}$
if X follows a normal dist then we get this curve. It's bell-shaped

\Rightarrow Symmetric Distribution
[50% each side]

so if we have a random variable that follows normal/gaussian dist then 68% of the entire data falls

$\frac{1}{2}$ std. to right
 $\frac{1}{2}$ std. to left

68 95% 99.7

$X = \{ \dots \}$

Empirical Rule [3-Sigma Rule]

With one standard deviation from the mean, Around 68% of the data falls here [i.e. below $\mu - \sigma$ to above $\mu + \sigma$]

68 - 95 - 99.7%

1st std 2nd std 3rd std

O-O plot & whether a Distribution is Gaussian/Normal

now do we check whether the data is Gaussian or not?

Distribution

[Quantile-Quantile]

Probability

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

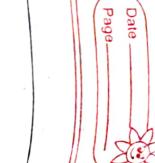
$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%$$

$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99.7\%$$

for weight, flight, IRIS DATASET

[Domain Experts]

Uniform Distribution



- ① continuous uniform distribution (PDF)
② discrete uniform distribution (PMF)

- ① continuous uniform distribution of continuous random variables

In probability theory and statistics, the continuous uniform distribution or rectangular distribution is a family of symmetric probability distributions.

The distribution describes an experiment whose tree is an arbitrary outcome that lies between certain bounds.

Re bounds are defined by the parameters a and b , which are the minimum and maximum values.

The interval can either be closed (e.g. $[a, b]$) or open ((a, b)).

Uniform

prob density fun

Notation : $f(x; a, b)$

parameters : $a < b & f > 0$

$\frac{1}{b-a}$

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } a < x < b \\ 0 & \text{otherwise.} \end{cases}$$

using maximum condition

using cumulative dist. fun



$$CDF = \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } a \leq x < b \\ 1 & \text{for } x \geq b \end{cases}$$

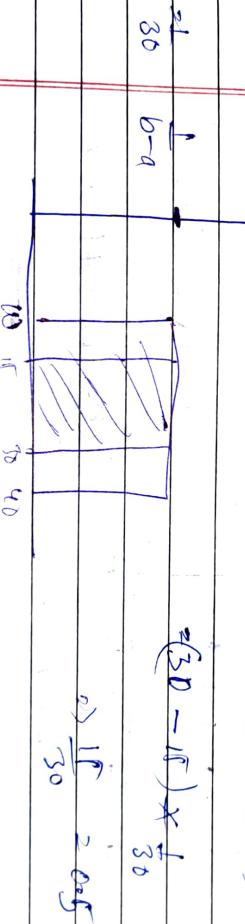


$$\text{mean} = \frac{a+b}{2} = f(x)$$

$$\text{Median} = \frac{1}{2} (\text{a} + \text{b})$$

$$\text{Variance} = \frac{1}{12} (b-a)^2 = \text{Var}(y)$$

The no. of candies sold daily at a shop is uniformly distributed with minimum of 10 and maximum of 40.



$$f(x \geq 20) = \frac{1}{30} \int_{20}^{\infty} (10 - 2x)^{-\frac{1}{2}} dx$$

$$\Rightarrow \frac{20}{30} = 0.66 \Rightarrow 66.66\%$$

$$\frac{b-a}{T} \propto (m - c_w)$$

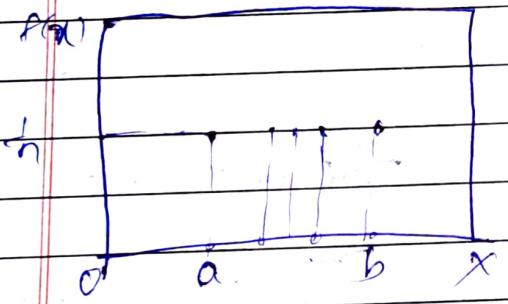
② Discrete uniform Distribution of Discrete Random variable }

In probability theory and statistics, the discrete uniform distribution is a symmetric probability distribution where in a finite no. of values are equally likely to be observed; every one of n value has equal probability $1/n$.

Another way of saying "discrete uniform distribution" would be "a known, finite no. of outcomes equally likely to happen".

discrete uniform
prob. needs n^{th}

e.g., rolling a dice {1, 2, 3, 4, 5, 6}



$$P(X=1) = 1/n \quad P(X=2) = 1/n \quad P(X=3) = 1/n$$

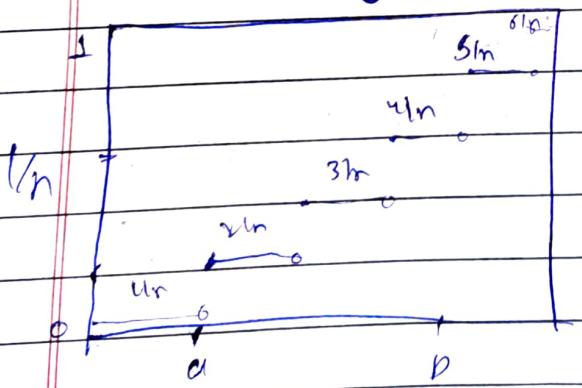
$$a=1 \quad b=6$$

$$n = b - a + 1 = 6$$

$$n = 5 \text{ where } n = b - a + 1$$

cumulative dist. funcⁿ

notation $U(a, b)$



parameters a and b $b \geq a$

$$\text{MF} \quad \frac{1}{n}$$

$$\begin{cases} \text{Mean} \\ \text{Median} \end{cases} = \frac{a+b}{2}$$



A

In uniform Distribution the prob. of getting the output is equal.

Standard Normal Distribution & Z-score

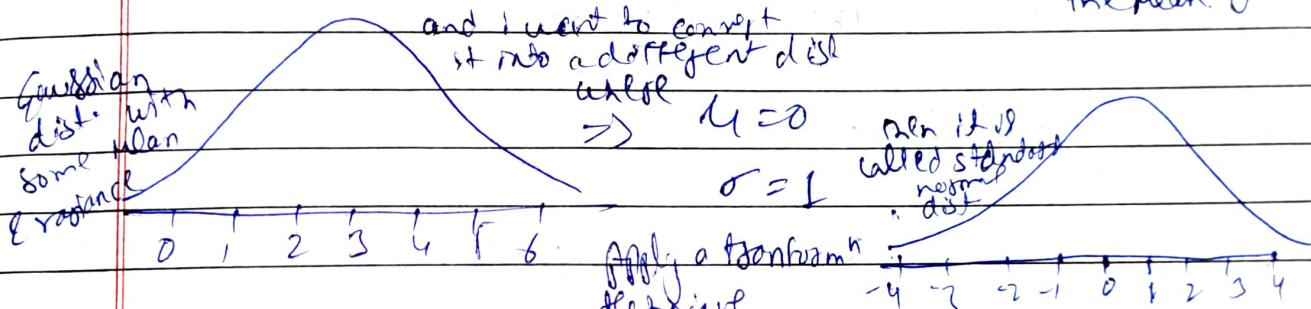
$$X = \{1, 2, 3, 4, 5\}$$

$$\mu = 3$$

$$\sigma = 1.414 \quad z = \frac{x - \mu}{\sigma}$$

It helps to identify that any value that how much std is away from the mean.

Gaussian dist. with some mean & variance



$$X = \{1, 2, 3, 4, 5\} \quad z\text{-score} = \frac{x_i - \mu}{\sigma} \quad Y = \{-2, -1, 0, 1, 2\}$$

$$\textcircled{1} \frac{1-3}{1} = -2$$

$$\textcircled{2} \frac{3-3}{1} = 0$$

$$\textcircled{3} \frac{2-3}{1} = -1$$

$$\textcircled{4} \frac{4-3}{1} = 1$$

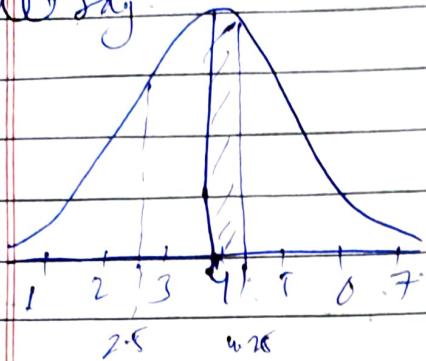
$$X \sim \text{SND}(\mu=0, \sigma=1)$$

z-score helps to know how many standard deviations a value is from the mean

Date _____

Page _____

g) Let say:



$$\mu = 4$$

$$\sigma = 2$$

z-score instance

Q How many standard deviation 4.25 is away from the mean?

$$x_i = 4.25$$

$$z\text{-score} = \frac{4.25 - 4}{2} = 0.25$$

$$x_i = 2.5$$

$$z\text{-score} = \frac{2.5 - 4}{2} = -1.5$$

| g) | Dataset | units age diff | | | weight | x x |
|---|---------|----------------|-----|------------------|--------|-----|
| | | year | kg | cm's + INR or \$ | | |
| when we solve a ML problem we have all the data | age | 70 | 175 | 40K | x | x |
| we have all the data | 25 | 60 | 160 | 50K | x | x |
| same don't scale | 26 | 58 | 170 | 60K | x | x |
| scale | 27 | 40 | 130 | 30K | x | x |
| Don't forget to bring it into same unit | 30 | 30 | 175 | 20K | x | x |
| scale | 31 | 28 | 180 | 70K | x | x |
| use z-score | 4 | 4 | 4 | 4 | x | x |
| use z-score of standardization | 4 | 4 | 4 | 4 | x | x |

$$z\text{-score} = \frac{x_i - \mu_{\text{age}}}{\sigma_{\text{age}}} \quad \frac{x_i - \mu_{\text{weight}}}{\sigma_{\text{weight}}} \quad \frac{x_i - \mu_{\text{height}}}{\sigma_{\text{height}}} \quad \frac{x_i - \mu_{\text{INR}}}{\sigma_{\text{INR}}} \rightarrow 4 \text{ models}$$

- (1) clustering algorithm
- (2) linear regression
- (3) logistic regression

[to bring estimate in a model]

Z-table [z-table & -ve z-table]

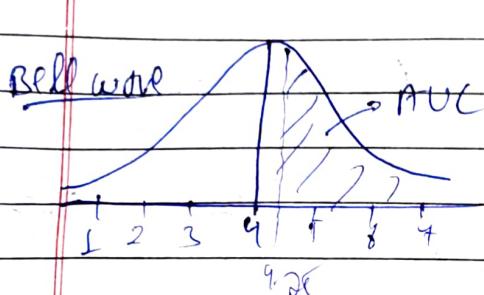
Date _____
Page _____

problem statement on Z-score [Z-table]

$$X = \{1, 2, 3, 4, 5, 6, 7\}$$

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

let, $\mu = 4$ $\sigma = 1$ Using these we can construct the curve



Area under curve

Total area = 1 \Rightarrow Symmetric Distribution

Z-table

Q) what percentage of scores fall above 4.25?

$$\bar{x}_i = 4.25$$

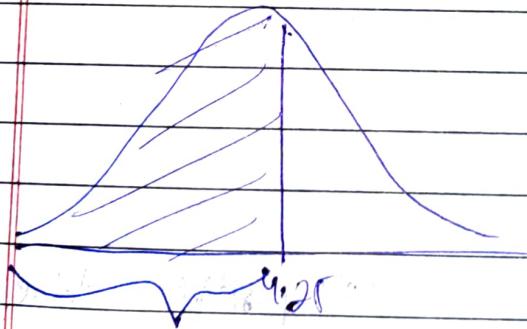
$$Z\text{-score} = \frac{\bar{x}_i - \mu}{\sigma}$$

$$\mu = 4$$

$$\sigma = 1$$

$$Z = \frac{4.25 - 4}{1} = 0.25 \Rightarrow \text{std}$$

Q) Look z-table



Q) What % of score falls above 4.25?
using z-table

$$1 - 0.59871$$

$$\Rightarrow 0.4013$$

$$\Rightarrow 40.13\%$$

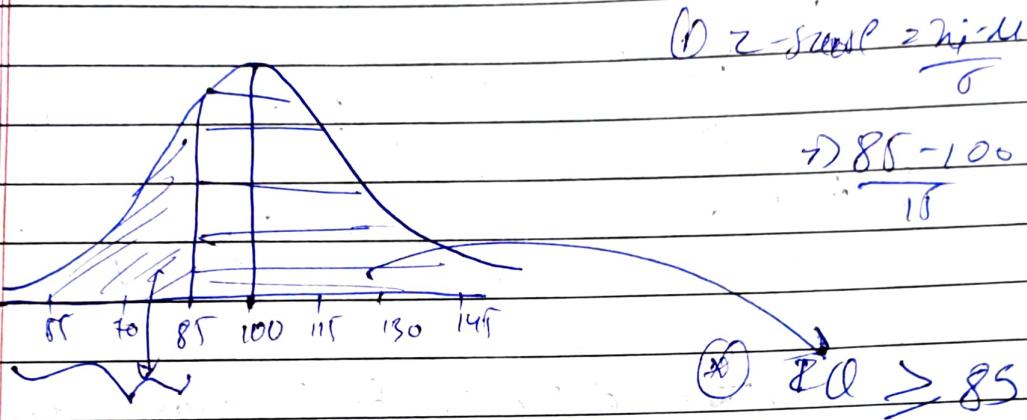
but we want Right area

[or above curve]

[In Z - ve z-table we'll get -0.25]

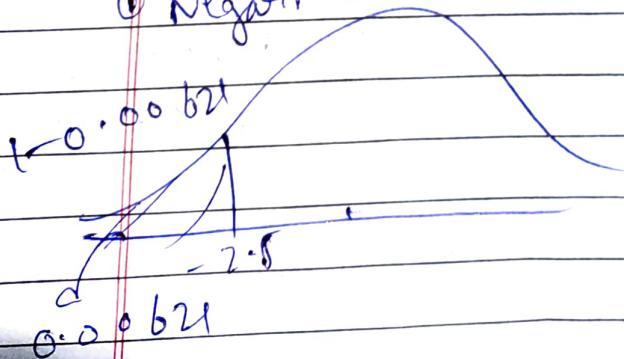
Q In India the average IQ is 100, with a standard deviation of 15. what is the percentage of the population would you expect to have an IQ lower than 85?

$\mu = 100 \quad \sigma = 15$

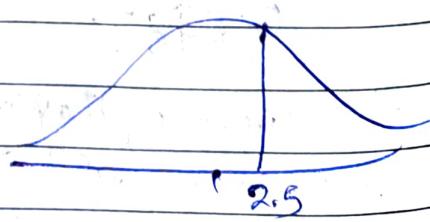


$TS \geq IQ \leq 100 \rightarrow$ Internal Assignment

① Negative z-table



② positive z-table



- It is imp. to understand some concepts like Sampling distribution or the properties of sampling dist.
- To calculate prob.

Date _____
Page _____

Central limit theorem

The central limit theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large no. of samples taken from population.

The central limit theorem says that the sampling distribution of the mean will always be normally distributed.

as long as the ~~size~~ sample size is large enough. regardless of whether the population has a normal, poisson, binomial or any other distribution, the sampling distribution of the mean will be normal.

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

$n = \text{sample size} \Rightarrow \text{any real}$

sample data

$$S_1 = \{n_1, n_2, n_3, \dots, n_n\} = \bar{n}_1$$

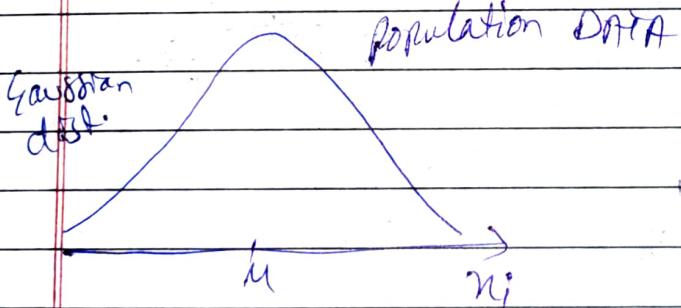
$$S_2 = \{n_2, n_3, \dots, n_n\} = \bar{n}_2$$

$$S_3$$

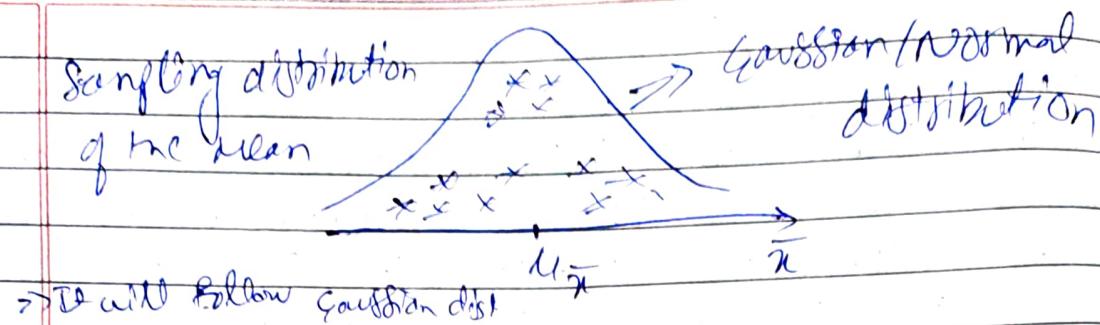
$$S_4$$

$$S_m$$

$$\bar{n}_m$$

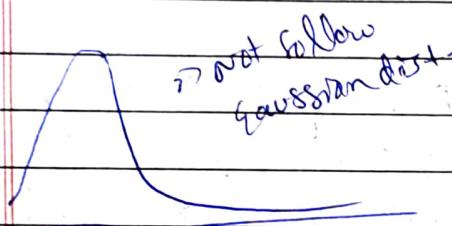


CLT

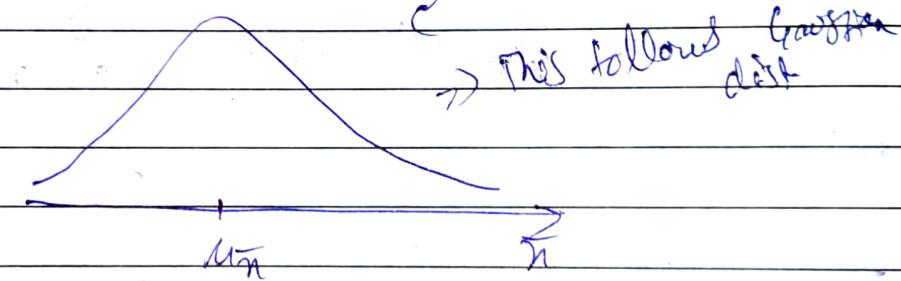


$$\textcircled{a} \quad X \not\sim N(\mu, \sigma)$$

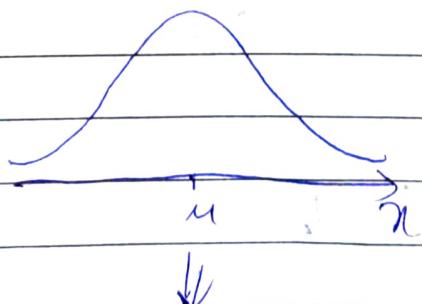
$\rightarrow (n \geq 30) \rightarrow$ Sample size

 s_1 \bar{x}_1 s_2 \bar{x}_2 \vdots
 \bar{x}_m

b
CLT



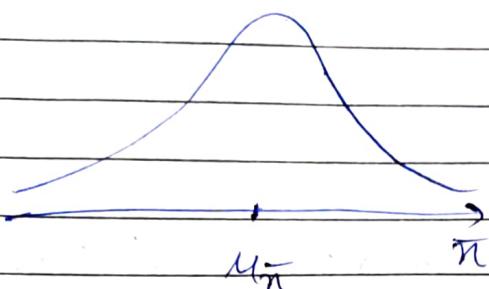
• Normal Distribution



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

normal
distribution

sampling distribution of mean



σ = population std
 μ = population mean
 n = sample size

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

\bar{x}_1

\bar{x}_2

\vdots

\bar{x}_m

effea

(c)
Design

now what's →

Inferential Statistics

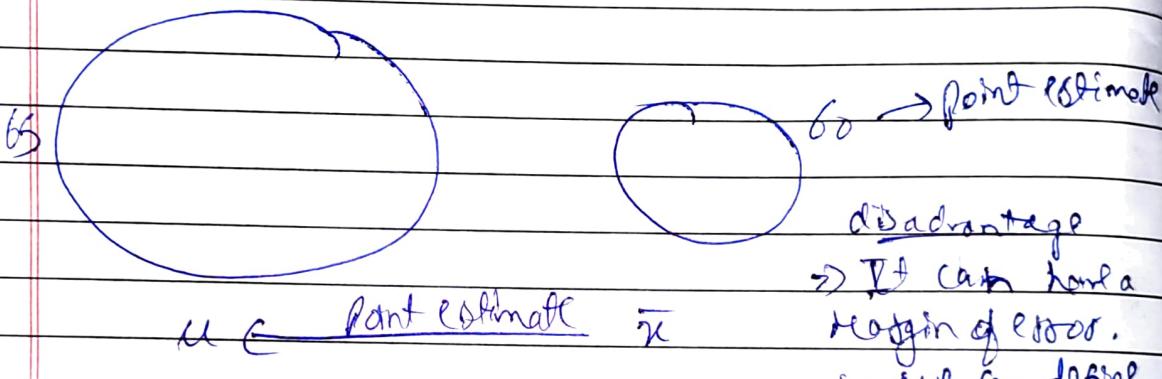
ISSE PECHHE DESCRIPTIVE
STATS THIS

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

Type(s)

① Point estimate: single numerical value used to estimate an unknown population parameter.

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



Disadvantage

→ It can have a margin of error.

→ So we can define interval.

② Interval estimate: range of values used to estimate the unknown population parameter.

$$[55 - 65] \Rightarrow \text{Sample Mean}$$

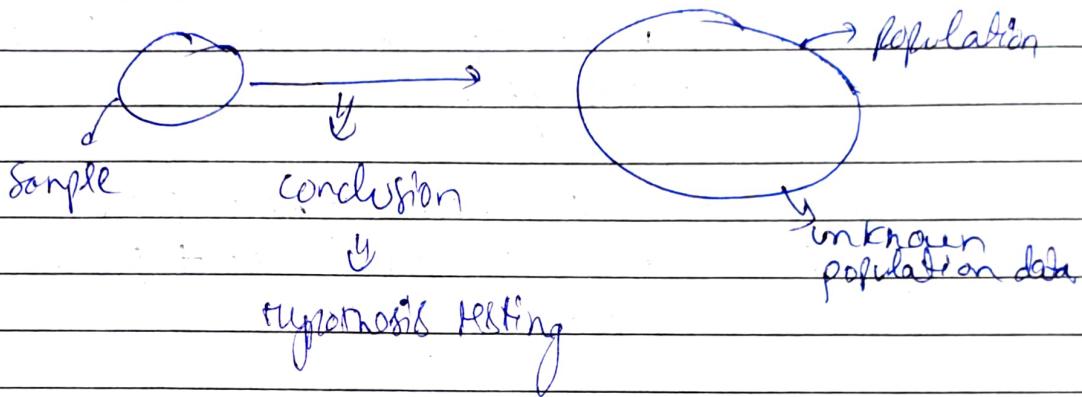
55 Point estimate 65

confidence interval.



A hypothesis And hypothesis testing mechanism

Inferrential stats → conclusion or inference



Hypothesis testing mechanism

e.g. person → crime
court
not guilty until proven

- ① Null Hypothesis (H_0) - person is not guilty.
- The assumption you are beginning with.
- ② Alternate hypothesis (H_1) - The person is guilty.
- opposite of Null hypothesis
- ③ Experiments → Statistical Analysis
→ collect proof. (DNA finger test)
(RT-PCR)
- ④ Accept the null hypothesis or reject the null hypothesis

{ p value,
significance value}

If colleges at District A states its average passed percentage of students are 85%. A new college was opened in the district and it was found that a sample of student 100 have a pass percentage of 90% with a standard deviation of 4%. Does this college have a different passed percentage.

Null hypothesis (H_0) = $\mu = 85\%$

Alternate hypothesis (H_1) = $\mu \neq 85\%$

Ex helps us to make some conclusions about no particular
data, so it is used in inferential statistics.

Date _____
Page _____

p values

The p value is a number, calculated from a statistical test, that describes how likely you are to have found a particular set of observation if the null hypothesis were true. If values are low in hypothesis testing to help decide whether to reject the null hypothesis.

$$p = 0.02 \rightarrow p < 0.05$$

out of 100 studies,

get



we found around 20 times in that region

hypothesis testing i.e. coin is fair or not (see Higgs)

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

$$P(H_0) = 0.6 \quad P(T) = 0.4$$

$$P(H_1) = 0.4 \quad P(T) = 0.3$$

① Null hypothesis:

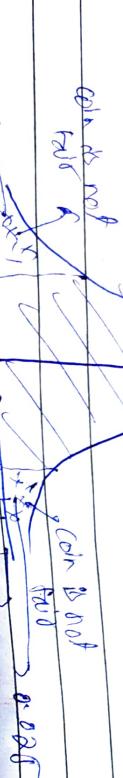
H_0 → coin is fair

② Alternative hypothesis:

H_1 → coin is not fair.

③ Experiment 100 times \rightarrow 45% tail (it is a concern)

coin is not fair \rightarrow coin is not fair \rightarrow may not be fair



(4) Significance value : $\alpha = 0.05$

$$CI = 1 - 0.05 = 0.95$$

(5) Conclusion : $P < \text{significance}$

Reject the null hypothesis
~~else~~

fail to reject the null hypothesis

Output is the f value t .

$P \leq \text{Significance}$
then reject the null type.

$P > \text{sig.}$
then accept or fail to reject the null type.

hypothesis testing and statistical analysis

- (1) 2 - test → ^{dependent} _{independent} → ^{ratio} _{nominal} → 2 t-test → 2 score and ^t value
 (2) ^t test → average → ^{interval} _{nominal} → t-table ^p value
 (3) CHI-SQUARE → categorical data
 (4) ANOVA → variance of data.

Z-test ^{continuous}
 population std in 22.30

- (1) The average height of all residents in a city is 168cm with a $\sigma = 3.9$. A doctor believes the mean to be different. He measured the height of 36 individuals and found the average height to be 169.5cm.
 At a 95% confidence level, is there enough evidence to reject the null hypothesis?

$$\begin{aligned} \mu &= 168\text{cm} & \sigma &= 3.9 & n &= 36 & \bar{x} &= 169.5 \\ \text{C.I.T} &= 0.95 & \alpha &= 1 - 0.95 & & & & = 0.05 \end{aligned}$$

- ① null hypothesis $H_0: \mu = \bar{x} = 168\text{cm}$

- ② alternate hypothesis $H_1: \mu \neq 168\text{cm}$

- ③ Based on C.I. we will draw decision boundary



Probability Density Function

Date _____

Page _____

$$\Rightarrow 1 - 0.025 = 0.975 \approx 2 \text{-Score}$$

$$Area = +1.96$$

Conclusion If Z is less than -1.96 or greater than $+1.96$, reject the null hypothesis.

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

$$Z = \frac{169.5 - 168}{3.9/\sqrt{36}} = \frac{1.5}{0.65} = \underline{\underline{2.31}}$$

$$\Rightarrow \underline{\underline{169.5 - 168}}$$

$$3.9/\sqrt{36}$$

$$Z = \frac{1.5}{0.65} = \underline{\underline{2.31}}$$

Conclusion

$2.31 > 1.96$ Reject the null hypothesis

$$P(Z < 0.975) = 0.975$$

$$P(Z < 2.31) = 0.98956$$

$$0.975 \quad 0.0104 \quad 1 - 0.0104 = 0$$

$$0.98956 \quad 0.0104 \quad 1 - 0.0104 = 0$$

$$-2.31$$

2) final conclusion the average ~~is 18.8m~~

The average height seems to
increasing based on sample heights

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

$$P < 0.05$$

$0.02088 < 0.05 \Rightarrow$ Reject the

null hypothesis

- Q2) A factory manufactures bulbs with a average
warranty of 5 years with standard deviation
of 0.5. A worker believes that the
bulbs will malfunction in less than 5 years.
He takes a sample of 40 bulbs and finds
the average time to be 4.07 years
(a) State null and alternative hypothesis
(b) At a 2% significance level, is there enough
evidence to support the idea that the warranty
should be revised?

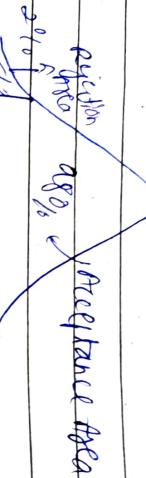
$$\mu = 5 \quad \sigma = 0.50 \quad n = 40 \quad \bar{x} = 4.07$$

Q3 Null hypothesis $H_0: \mu \leq 5$

Alternative hypothesis $H_1: \mu > 5$ { 1 tail test}

Decision Boundary

Date _____
Page _____



(Q) Z_{test}

$$Z_t = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{48 - 50}{0.50/\sqrt{50}}$$

$$= -2.53$$

Area under curve with Z -score $-2.53 = 0.0570$

$$\text{P-value} = 0.0570$$

$$\alpha = 0.02$$

compare P-value with α

$$0.0570 > 0.02 \Rightarrow \text{false}$$

we accept the null hypothesis

we fail to reject the null hypothesis.

student t distribution

In Z stats when we perform any analysis using Z score we require σ (population standard deviation) \rightarrow is already known

* How do we perform any analysis when we don't know the population standard deviation?

student's t distribution

t-stats

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

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

S = sample
standard deviation

Z table

t table \rightarrow t test

* Degree of freedom

$$df = n - 1 \Rightarrow 3 - 1 = 2$$

let consider

3 people & three chairs in a room



when the first person comes inside a room, there has 3 diff. options, then after this 2nd person comes it has 2 option.

but the third person doesn't have any option.
so degree of freedom is 2

T-test \div T test \rightarrow One Sample T-test

- ① In the population the average IQ is 100.
A team of researchers want to test a new
dissertation medication to see if it has either
a positive or negative effect on intelligence
or no effect at all.

A sample of 30 participants who have taken
the medication has a mean of 114 with a
standard deviation of 20.

Does the medication affect intelligence?
 $C_i = 95\%$ $S_i = 5$

$$\bar{x} = 114 \quad n = 30 \quad \bar{N} = 140 \quad S = 20$$

② Null hypothesis: $H_0 \div \mu = 100$

③ $\alpha = 0.05$

④ Degrees of freedom

$$df = n - 1 = 30 - 1 = 29$$

Decision Rule

