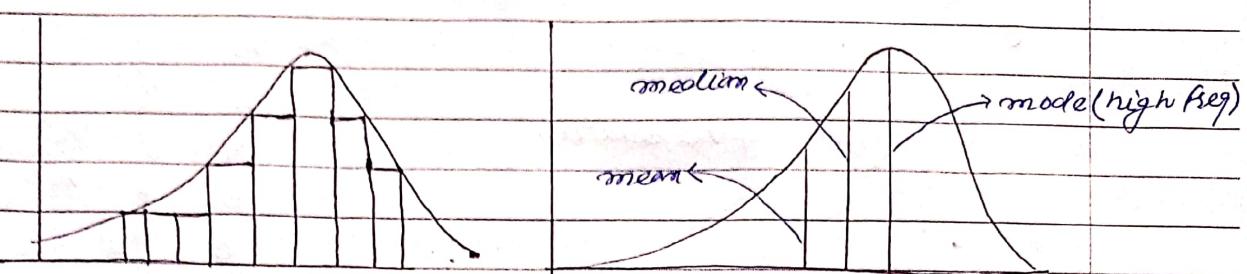


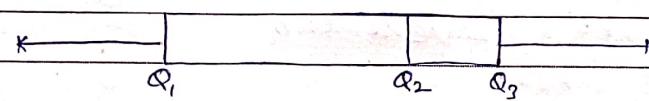
③ Left Skewed Distribution



- Left skewed also called negative skewed
- Relation between mean, median, mode

$$\boxed{\text{mean} \leq \text{median} \leq \text{mode}}$$

Box Plot



$$Q_2 - Q_1 > Q_3 - Q_2$$

Covariance and Correlation

Both are used to solve "what is the relationship between x and y?"

Correlation (co-relation between x and y)

x	y
2	3
4	5
6	7
8	9

$x \uparrow y \uparrow$
 $x \downarrow y \uparrow$ } find relation between
 $x \uparrow y \downarrow$ } $x \in$ and $y \in$
 $x \downarrow y \downarrow$

There are two types of relations are possible

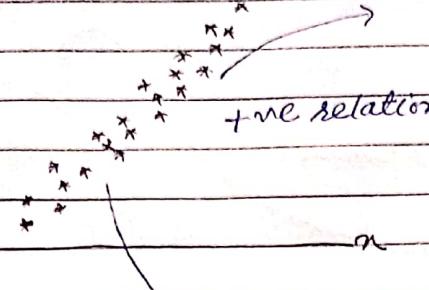
(i)	$x \uparrow y \uparrow$	or (ii)	$x \downarrow y \uparrow$
	$x \downarrow y \uparrow$		$x \uparrow y \downarrow$

Important
Example
Visit
Points
e.g
⇒Title:
Sub Title:
Highlighted:

① If

 $x \uparrow y \uparrow, x \downarrow y \downarrow$

y



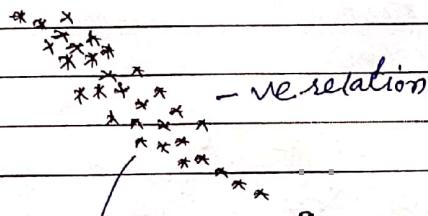
x	y
2	3
4	5
6	7
8	9

$x \uparrow y \uparrow$ } In this relation scatter plot
 $x \downarrow y \downarrow$ looks like this shape

② If

 $x \uparrow y \downarrow, x \downarrow y \uparrow$

y



x	y
5	4
6	5
7	9
9	11

$x \uparrow y \downarrow$ } In this correlation, scatter
 $x \downarrow y \uparrow$ plot looks like this

* Scatter plot can be used to find correlation between two random continuous variable.

↓ SIZE OF HOUSE ↑ PRICE OF HOUSE ↓

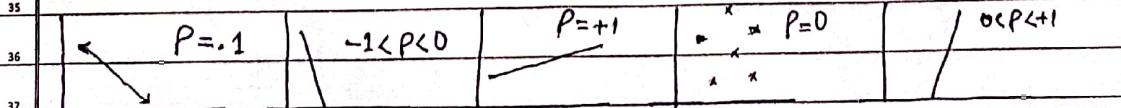
This type of relation is a high relation and strong fetecher

* Correlation value range in between -1 to +1.

The value more towards +1 they are more positively corr.

The value more towards -1 the more negatively correlated

$$P_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$$



Note ():
Note ():

Covariance

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

$x_i \rightarrow$ Data points of x

$\bar{x} \rightarrow$ Sample mean of x

$y_i \rightarrow$ Data points of y

$\bar{y} \rightarrow$ Sample mean of y

$n \rightarrow$ x or y data size

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

(i) $\text{Cov}(x, y)$ $x \uparrow \begin{matrix} y \uparrow \\ y \downarrow \end{matrix} \rightarrow +ve \text{ correlation}$

(ii) $\text{Cov}(x, y)$ $x \downarrow \begin{matrix} y \uparrow \\ x \uparrow \end{matrix} \rightarrow -ve \text{ correlation}$

Eg. x y

2 3

4 5

6 7

$\bar{x} = 4$ $\bar{y} = 5$

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

$$= \frac{[(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)]}{n-1}$$

$$= \frac{4+0+4}{2} = 8/2 = 4 \text{ +ve covariance}$$

x and y having a +ve covariance

Advantages

Disadvantages

Relationship between x and y has a +ve or -ve value

Covariance does not have a specific limit value

Covariance can be vary between

$[-\infty \text{ and } +\infty]$

$[-\infty \text{ and } +\infty]$

You can fix this Problem by using

- ① Pearson Covariance Coefficient
- ② Spearman's Covariance Coefficient

Important Example Visit Points

Title: e.g. Sub Title: Highlighted:

St-16

Page No. / /
Date / /

Person Covariance Coefficient

$$\rho_{xy} = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \quad \text{Range } [-1 \text{ to } 1]$$

* The more the value toward +ve the more correlated it is (x, y)

* The more the value toward -ve the more -ve correlated it is (x, y)

Spearman Rank Correlation Coefficient:

$$\rho_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \sigma(R(y))} \quad \text{Range } [-1 \text{ to } +1]$$

x	y	R(x)	R(y)
1	2	5	5
3	4	4	4
5	6	3	3
7	8	2	1
8	7	1	2

→ x values more high R(x) more less from 8 to 1

Feature Selection

+ve	+ve	+ve	-ve
Size of House ↑	No of Rooms ↑	Location ↑	Hunted ↑ Price ↑

① Get only important features

Correlation and Covariance by using Python

```
df = sns.load_dataset('health-exp')
```

```
df.head(); df.cov()
```

```
np.cov(df['columnname'])
```

Note ():

Note ():

Important Example Visit Points

Title:
Sub Title:
Highlighted:

8/12

Page No.....
Date/..../....

Probability Distribution Function.

→ It's denotes distribution of a data

Probability Density Function
pdf

Probability Mass Function
pmf

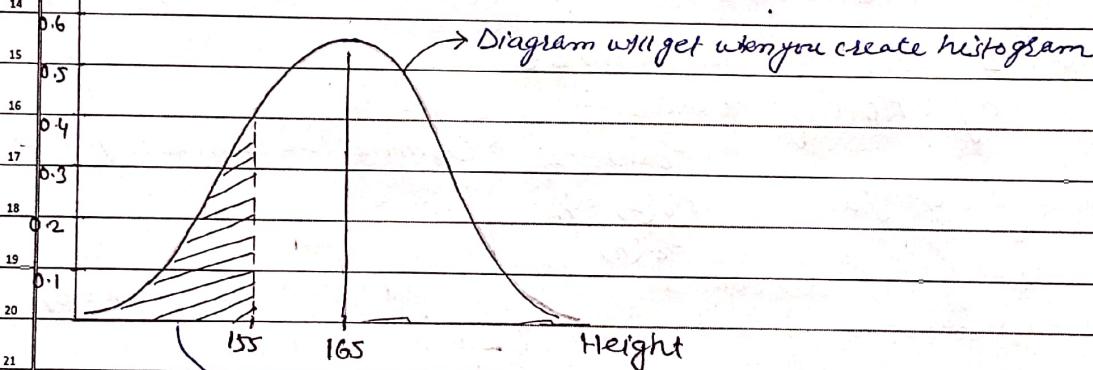
Continuous Random Variable
Range $[0, 1]$

Discrete Random Variable
 $[-\infty, +\infty]$

Probability Density Functions: It is also called
Continuous Probability Distributions

Eg: Height of students in classroom
lets assume $\bar{x} = 165$

Probability Density



$P(x \leq 155) = \text{Area under the curve}$

Kernel Density Estimator Smooth histogram
to draw a line plot

Probability Mass Function

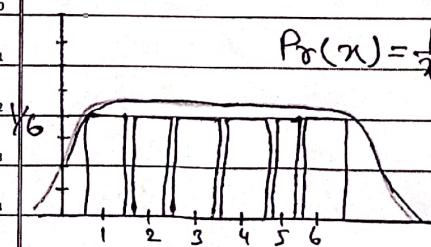
Show Distribution for Discrete Random Variable

Eg: Rolling a Dice $\{1, 2, 3, 4, 5, 6\}$

$$P(x) = \frac{1}{6} \text{ (according to dice)}$$

$$P(x=1) = \frac{1}{6}$$

$$P(x=5) = \frac{1}{6}$$



Problem

$$\begin{aligned} P(x \leq 4) &= P(x=1) + P(x=2) + \\ &\quad P(x=3) + P(x=4) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{2}{3} \end{aligned}$$

3 Cumulative Distribution Function (CDF)

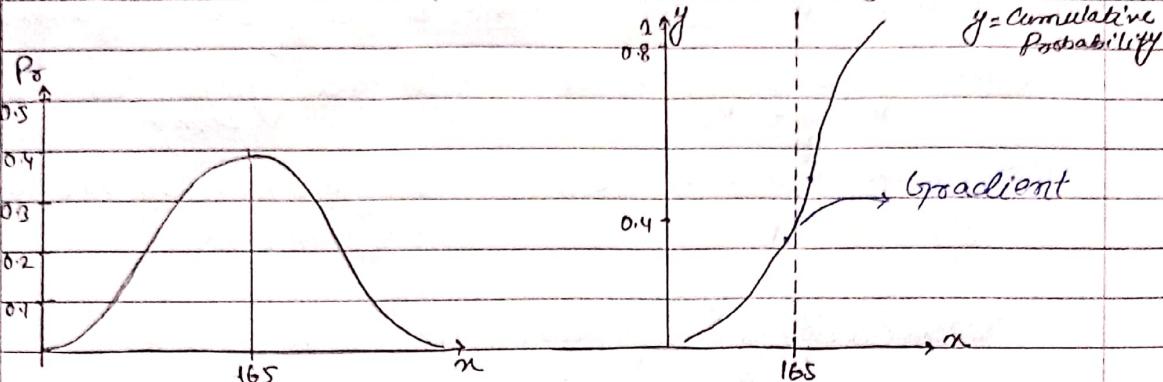
Note ():
Note ():

Important Example Visit Points
e.g. ⇒

Title:
Sub Title:
Highlighted:

Page No.
Date / /

1 Cumulative Distribution Function (cdf)

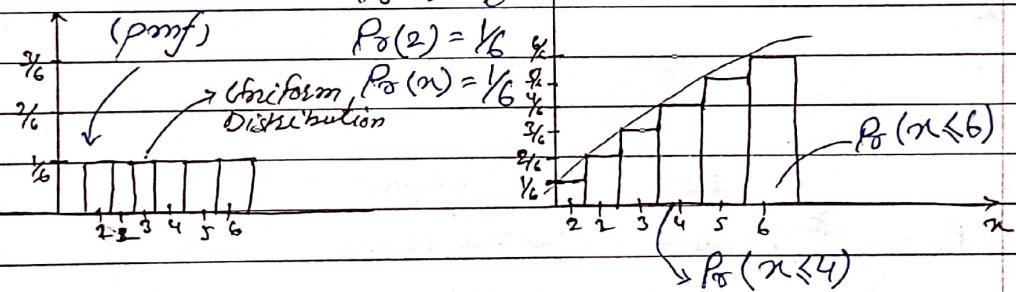


13 Probability Density Function & Probability Mass Function

15 PMF: Discrete Random Variable

16 Eg : Rolling a dice $\{1, 2, 3, 4, 5, 6\}$

$$P_{\sigma}(1) = \frac{1}{6} \uparrow \text{Cumulative Probability}$$

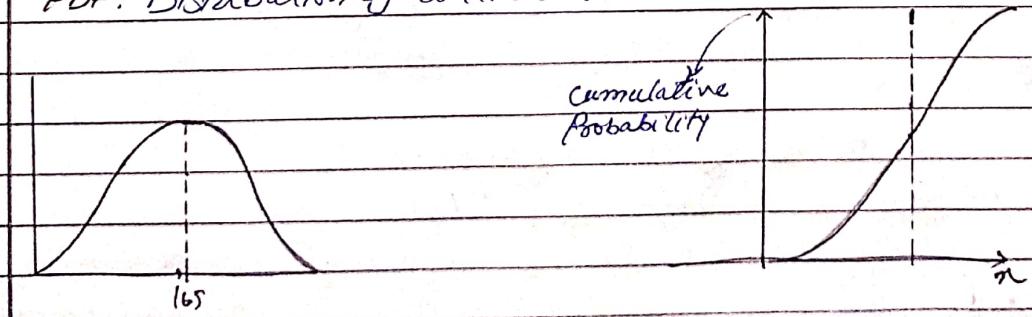


$$P_{\sigma}(x \leq 2) = P_{\sigma}(x=1) + P_{\sigma}(x=2)$$

$$= \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

$$P_{\sigma}(x \leq 6) = 1$$

27 PDF: Distribution of continuous random variable



35 Probability density \rightarrow Gradient of cumulative curve

37 Note ():

Note ():

Types of Probability Distribution

- 1 ① Normal / Gaussian Distribution
- 2 ② Bernoulli Distribution
- 3 ③ Uniform Distribution
- 4 ④ Poisson Distribution
- 5 ⑤ Log Normal Distribution
- 6 ⑥ Binomial Distribution

1. Bernoulli Distribution

In probability theory and statistics, the Bernoulli distribution, named after Swiss mathematician Jacob Bernoulli;

(i) is the discrete probability distribution of a random variable which has two outcomes Y and N also known as James or Jacques.

(ii) or outcomes Binary

Eg: Tossing a coin {H, T}

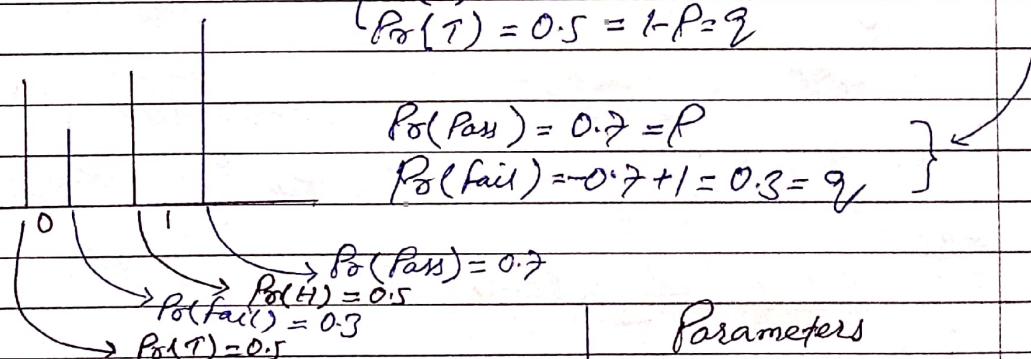
Eg: Whether the person will pass or fail,

$$\rightarrow P(H) = 0.5 = P$$

$$P(T) = 0.5 = 1 - P = q$$

$$P(\text{Pass}) = 0.7 = P$$

$$P(\text{fail}) = 0.3 = 1 - P = q$$



Parameters

$$0 \leq P \leq 1$$

$$q = 1 - P$$

$$K = \{0, 1\}$$

0 = fail, 1 = success

Pf Pmf

P

$$P(x=0) = 0.5 \text{ and } P(x=1) = 0.5$$

$$P(x=0) = 0.3 \text{ and } P(x=1) = 0.7$$

Note (): Discrete has finite value from 0, 1, 2 ...

Note ():

Pmf for Bernoulli Distribution

$$\text{Pmf} \rightarrow P(X=x) = p(x) = \begin{cases} P^n (1-P)^{1-n} & n=0, 1 (x \text{ belongs to } 0 \text{ and } 1) \\ 0 & \text{otherwise} \end{cases}$$

where $P = \text{Probability of success}$

$1-P = q = \text{Probability of failure}$

number of trial = 1 (trial \rightarrow Experiments)

number of trial = one

Success (P)

failure ($1-P$)

$$\text{If } x=1 \Rightarrow P(n) = P^n (1-P)^{1-n}$$

$$= P(1) = P^1 (1-P)^0$$

$$= P(1) = P \quad \# \text{Probability of Success}$$

$$\text{If } x=0 \Rightarrow P(n) = P^n (1-P)^{1-n}$$

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

$$P(0) = 1 (1-P)^1$$

$$P(0) = 1 - P = q$$

* If you know
 $(n)^0 = (n)^1 = 1$

Verify is this Pmf or not, if proof

$$\sum_{n=0}^1 P(n) = 1$$

Step $\sum_{n=0}^1 p(n) =$

* $\sum_{n=0}^1 p(0) = P^0 (1-P)^{1-0} = 1 - P = q$

Pmf $\sum_{n=0}^1 p(1) = P^1 (1-P)^{1-1} = P$

$k=n$

$$\therefore \sum_{n=0}^1 P(n) = 1 - P + P = 1$$

Moment Generating Function.

$$M_X(t) = \sum_{n=0}^1 e^{tn} P(n) \Rightarrow M_X(t) = \sum_{n=0}^1 e^{tn} P^n (1-P)^{1-n}$$

If $x=0$

$$M_X(t) = \sum_{n=0}^1 e^{tn} P^n (1-P)^{1-n} = e^0 P^0 (1-P)^{1-0} = 1 - P$$

If $x=1$

$$M_X(t) = \sum_{n=0}^1 e^{tn} P^n (1-P)^{1-n} = e^t P (1-P)^0 = Pe^t$$

Now $1 - P + Pe^t = q + e^t P$

Note ():

Note ():

Mean of Bernoulli Distribution:

$$E(K) = \sum_{i=0}^K k \cdot P(k)$$

$$= [0 \cdot 0.4 + 1 \cdot 0.6] \\ = 0.6 = P$$

$$k = x$$

mean → E

$$\text{lets } P_{\sigma}(K=1) = 0.6$$

$$P_{\sigma}(K=0) = 0.4$$

It means the mean of bernoulli distribution be always P , $\neq 0$ does not has any value.

Median of Bernoulli Distribution

$$\text{median} \begin{cases} 0 & \text{if } P < \frac{1}{2} \\ [0,1] & \text{if } P = \frac{1}{2} \\ 1 & \text{if } P > \frac{1}{2} \end{cases}$$

Variance

$$\text{var} = P * (1-P)$$

Standard deviation

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

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 question "yes → success (with probability b) or no → failure $q = 1-p$ ".

A single success/failure experiment is also called a Bernoulli trial or Bernoulli Experiment.

A sequence of outcomes is called a Bernoulli process for a single trial; i.e. $n=1$, the binomial distribution is a Bernoulli Distribution (in case of single experiment).

① Every experiment outcomes "yes"

② Then experiment is performed "no" n trials

Notation : $B(n, p)$ B for binomial distributionParameters: $n \in \{0, 1, 2, 3, 4, \dots\}$ → Number of trial / experiment $p \in [0, 1]$ → Success probability for each trial

$$q = 1 - p$$

It is a fixed number of experiments.

Eg: Tossing a coin 10 times

 P = Probability of success / failure n = numbers of experiments

It is fixed or countable

Support: $n \neq k$ $k \in \{0, 1, 2, 3, \dots, n\}$ → Number of success

Pmf $P_x(k, n, p) = {}^n C_k p^k (1-p)^{n-k}$

Binomial Distribution

for $k = 0, 1, 2, 3, 4, \dots, n$ where $x = 1, 2, 3, 4, \dots, n$

$${}^n C_k = \frac{n!}{k!(n-k)!}$$

no of trials, n = finiteP. of success, p = constant

Eg: 35% of all voters support Proposition A. If

a random sample of 10 voters is polled.

What is the probability that exactly three of them support the proposition

Find $P(n=3)$ if $n=10$ and $p=0.35$

$$P(n=3) = \frac{n!}{x!(n-x)!} p^x q^{n-x} = \frac{10!}{3!7!} (0.35)^3 (0.65)^7 = 0.2522$$

Ans: 25.22%.

Mean :

mean = np

variance

var = npq

std

std = \sqrt{npq}

Proof mean of binomial distribution

Important
Example
Visit
Points

Title:
Sub Title:
Highlighted:

31-22-a

Page No.....
Date / /

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

Note ():
Note ():

Important
Example
Visit
Points
⇒

Title:
Sub Title:
Highlighted:

Properties of Binomial Distribution

1. The number of observation n is fixed
2. Each observation is independent
3. Each observation represents one of two outcomes ("success" or "failure")
4. The probability of "success" p is the same for each outcomes.
5. No previous event impact on next event.

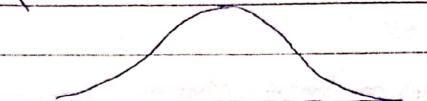
Note { }:
Note { }:

Random Variables.

Probability Distribution of a Discrete Random Variable



Probability Distribution of a Continuous Random Variable



Poisson Distribution

In statistics, Poisson distribution is a probability distribution used to show how many times an event is likely to occur over a specified period.

In a poisson distribution number of experiment

is very large (∞) but number of success $\{ X = 1, 2, 3, \dots, \infty \}$
is very low. $\{ \text{number of trial, } n \rightarrow (\text{very large})$
 $\text{Prob. of success, } P \rightarrow 0 \text{ (very low)}$

Eg: Lic holder

number of policy holder is very high = ∞

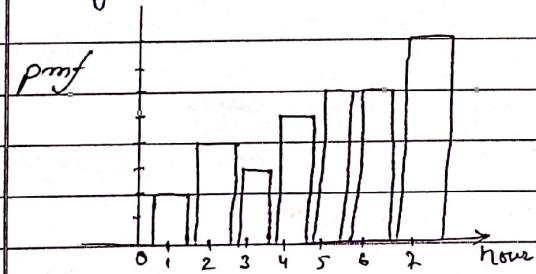
but number of policy holder death is very low

Eg2: Numbers of word on book page is very high = ∞

but number of miss print word in book very low

*Describe number of event occurring in a fix time interval

Eg1: Number of people visiting hospital every hour



$$\lambda = 3$$

Lambda is a -expected number of events occur at every time interval.

Question based on graph.

What is the probability of a person to come (in bank)
at the fifth hour?

Pmf

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

(x = # of time interval)

lets $\lambda = 3$

$$P(x=5) = \frac{e^{-3} 3^5}{5!} = 0.101 \approx 10\%$$

Question 2: What is the probability for a person to comes less than or equal to fifth hours

some

$$P_0(x \leq 5)$$

$$P_0(n=1) + P_0(n=2) + P_0(n=3) + P_0(n=4) + P_0(n=5)$$

Mean:

Variance:

$$\text{mean} = E(n) = \lambda t = \lambda * t$$

Expected of n

 λ = is a expected number of events

at every time interval

t = time interval

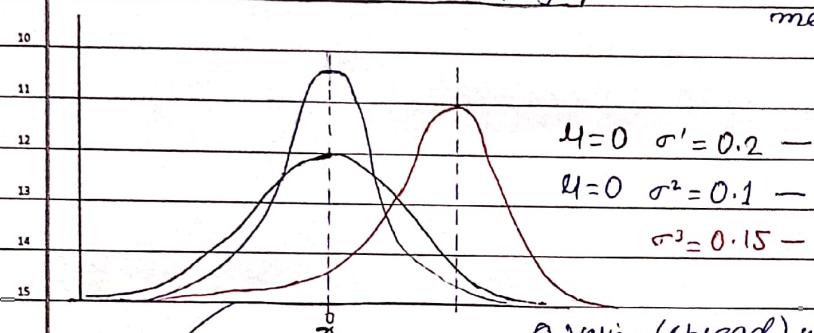
Normal/Gaussian Distribution

In statistics, a normal distribution or Gaussian distribution is a

- type of continuous probability distribution
- for a real-valued random variable
- and a bell-shaped curve
- with an equal number of measurements above and below the mean value

Continuous random variable = pdf

mean = median = mode



① variance (spread) will increase if std increase

Probability density function

If variance = n and std = y

then $n \uparrow \Rightarrow y \uparrow$] Positive
 $n \downarrow \Rightarrow y \downarrow$] Relation

More variance more gradient

② Relation between variance and gradient

If variance = n and gradient = y

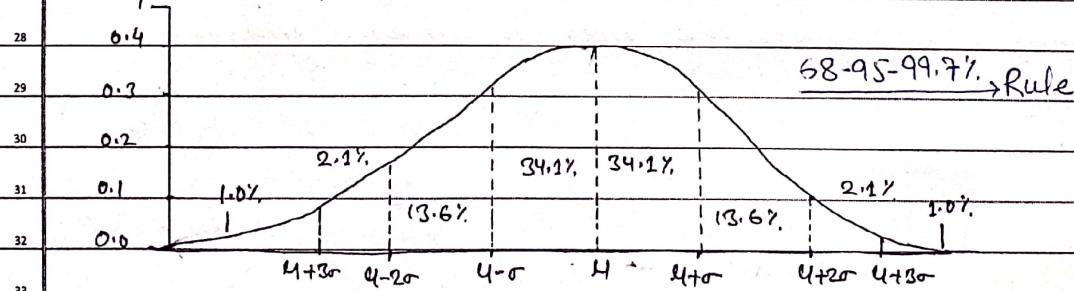
then $n \uparrow \Rightarrow y \uparrow$] Positive Relation
 $n \downarrow \Rightarrow y \downarrow$

low variance
low gradient

*

Area under the curve is always = 1

Empirical Rule of Normal Distribution:



$P(\mu-\sigma < x \leq \mu+\sigma) \approx 68\%$. 1st std

$P(\mu-2\sigma < x \leq \mu+2\sigma) \approx 95\%$. 2nd std

$P(\mu-3\sigma < x \leq \mu+3\sigma) \approx 99.7\%$. 3rd std

} Data Distribution.

Note ():
Note ():

pdf

Notation: $N(\mu, \sigma^2)$

Parameters:

$\mu \in R$ = mean

$\sigma^2 \in R > 0$ = variance

$x \in R$

Eg: In which types of data set follow Normal Distribution

(1) Height of the people

(2) Height of the student in the class

(3) Height of student in the class

(4) IRIS data set

(5) Q-Q plot {Quantile-Quantile Plot}

continuous
discrete

Probability in Uniform Distribution

In probability theory and statistics, the

continuous uniform distribution or rectangular distribution is a family of symmetric probability distributions.

The distribution describes an experiment where there is an arbitrary outcome that lies between certain bounds.

Where a is minimum and maximum limit

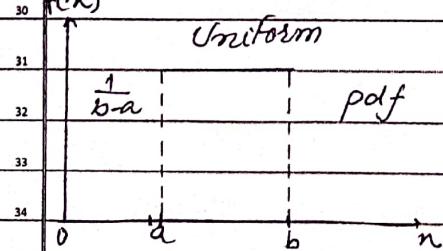
The bounds are defined by the parameters, a and b , which are the minimum and maximum values

$f(x)$

Notation: $U(a, b)$

Uniform a, b

Parameters: $-\infty < a < b < \infty$



$$\text{pdf} \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

$$\text{cdf} \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a, b] \\ 1 & \text{for } x > b \end{cases}$$

Note ():
Note ():

Important Example Visit Points

Title:
e.g.
Sub Title:
Highlighted:

St 23

Page No.....
Date/...../.....

Mean!

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

Median

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

Variance

$$\frac{1}{12}(b-a)^2$$

Eg 1: The number of canopies sold daily at a shop is uniformly Distribution with a maximum of 40 and minimum of 10.

i) Probability of daily sales to fall between 15 and 30

$$(where n_1 = 15 \text{ and } n_2 = 30)$$

$$P(n_1 \leq n \leq n_2) = (n_2 - n_1) \frac{1}{b-a}$$
$$= (30-15) \frac{1}{40-10} = \frac{1}{2} = 0.5$$

$$\text{Solve 2 } P(n \geq 20) = (40-20) \frac{1}{30}$$
$$= 0.66 = 60\%$$

Discrete Uniform Distribution (pmf)

In probability theory and statistics, the discrete uniform distribution is a symmetric probability distribution wherein a finite number of values are equally likely to be observed; every one of n values has equal probability $\frac{1}{n}$. An other way of saying "discrete uniform distribution" would be a known finite number of outcomes equally likely to happen

$$f(x)$$

Eg: Rolling a dice

$$\{1, 2, 3, 4, 5, 6\} \quad P(1) = \frac{1}{6}$$

$$a=1 \text{ and } b=6 \quad P(2) = \frac{1}{6}$$

$$\frac{1}{n} \Rightarrow n=b-a+1 \quad P(3) = \frac{1}{6}$$

$$\frac{1}{n}$$

$$p_{mf}$$

$$a \quad b \quad x$$

cdf

Notation: $U(a, b) \rightarrow$ uniform

Parameters: a, b with $b > a$

Pmf $\frac{1}{n}$

Mean $\frac{a+b}{2}$

Median $\frac{a+b}{2}$

Note ():

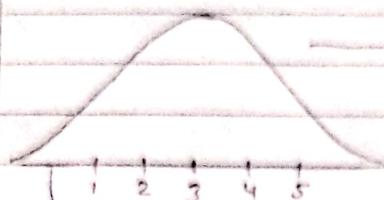
Note ():

Z-Score and Z-test

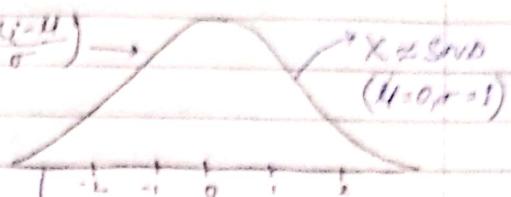
Standard Normal Distribution And Z-score

→ where mean is zero

$X = \{1, 2, 3, 4, 5\}$ where mean = 3 and $\sigma = 1.414 \approx 1$



$$\rightarrow Z\text{-Score} = \frac{(x_i - \mu)}{\sigma}$$

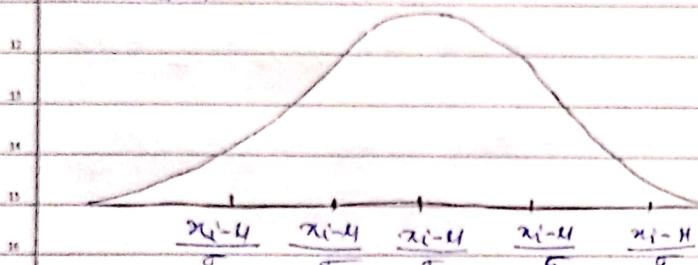


$$X \sim N(0, 1)$$

Normal Distribution

Standard Normal Distribution

Solution



$$\frac{x_1 - \mu}{\sigma} = \frac{1 - 3}{1} = -2$$

$$\frac{x_2 - \mu}{\sigma} = \frac{2 - 3}{1} = -1$$

Eg1: In $X = \{1, 2, 3, 4, 5\}$, how much away from Standard deviation

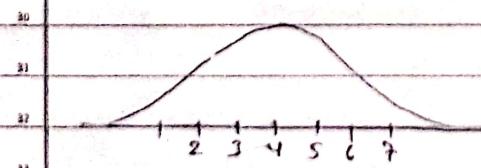
some Z-score = $\frac{x_i - \mu}{\sigma}$, where $x_i = 4$ and $\mu = 3$

$$= \frac{4 - 3}{1} = 1 \quad (+ve to the right)$$

So 4 is 1 Standard deviation away from mean

Z-score define how much away from the mean a single random variable

Eg2: How many Standard deviation 4.5 is away from mean?
where $\mu = 4$ and $\sigma = 1$



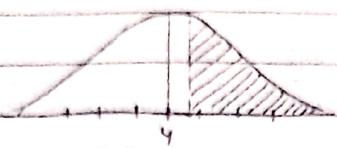
$$Z\text{-Score} = \frac{x_i - \mu}{\sigma} \quad x_i = 4.5$$

$$\mu = 4$$

$$\sigma = 1$$

$$= \frac{4.5 - 4}{1} = 0.5$$

Eg: How much % of data is greater than 4.5



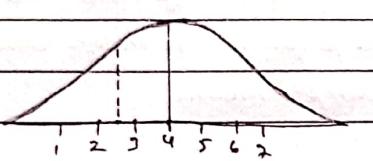
given $\bar{x}_i = 4.5$
 $\mu = 4$
 $\sigma = 1$

Some Z-Score = $4.5 - 4 = 0.5$

Area under the curve (> 4.5) = $1 - 0.69146$ +ve
table value
(white area under the curve) $= 0.30854$

In Percentage = 30.85%.

Eg: What % of data is falling below 2.5?



given $\bar{x}_i = 2.5$
 $\bar{x}_i - \mu = 2.5 - 4$
 σ $\mu = 4$
 $= -1.5$ $\sigma = 1$

Area under the curve (< 2.5) = 0.06681
= 6.6%.

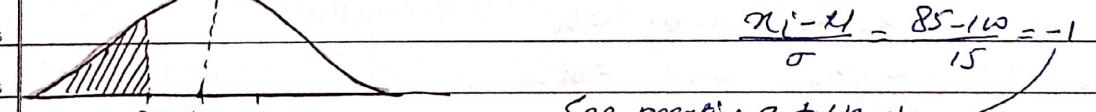
Eg: In India the average IQ is 100, with a standard deviation of 15. With a percentage of the population would you expect to have an IQ lower than 85?

Solve

given

$$\bar{x}_i = 85, \sigma = 15, \mu = 100$$

$$\frac{\bar{x}_i - \mu}{\sigma} = \frac{85 - 100}{15} = -1$$



See negative Z-table -1

Area under the curve = 0.15866.
= 15.866

Area under the curve (> 85) = $1 - 0.15866$
 $\approx 84\%$

Central Tendency Vs Central limit Theorem

Central Tendency

The central tendency refers to a measure that describes the typical or central values of a data set. The most common measures are the mean, median, and mode.

Central limit Theorem

The CLT states that if you have a large sample size from a population with a finite mean & variance, then the distribution of the sample mean will approximate a normal distribution, regardless of the distribution of the original population. This theorem is important because it allows us to make inference about the population mean based on the sample mean.

Central limit Theorem (CLT)

The central limit theorem says that

① $\bar{x} \sim N(\mu, \sigma)$ (a) the sampling distribution of the mean will always be normally distributed, as long as the sample size is large enough.

(lets $n=20$ and $S=\text{sample}$

$$S_1 = \{x_1, x_2, x_3, \dots, x_{20}\} = \bar{x}_1,$$

$$S_2 = \{x_1, x_3, x_4, \dots, x_{20}\} = \bar{x}_2,$$

$$S_3 = \{x_5, x_6, x_7, \dots, x_{20}\} = \bar{x}_3,$$

$$S_4 = \{x_1, x_9, x_3, \dots, x_{20}\} = \bar{x}_m$$

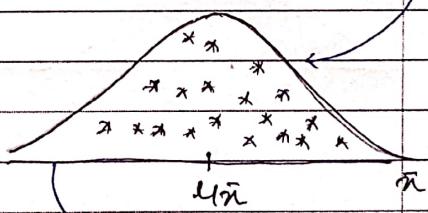
$$S_m = \{ \dots \} = \bar{x}_m$$

Sample data can be repeated

Gaussian Distribution

$n = \text{any size from population in GP}$

$$\bar{x} = \{\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_m\}$$



Normal Distribution

If you have population that is normally distributed.

And if you get different sample from population

$S_1, S_2, S_3, \dots, S_m$ (Some data can be repeated) and number of sample

is n , then find all samples mean where $S = 1/n$ and $n = \bar{x}$

Now you get Normal Distribution.

② (b) Regardless of whether the population has a normal, Poisson, binomial, or any other distribution.

Note (): Relate to (a) and (b)

Note ():

$$\{ X \sim N(\mu, \sigma^2) \rightarrow n \geq 30 \}$$

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

Non Gaussian

$$S_1 = \{x_1, x_2, x_3, \dots, x_{30}\} = \bar{x}_1$$

$$S_2 = \{x_2, x_5, x_7, \dots, x_{30}\} = \bar{x}_2$$

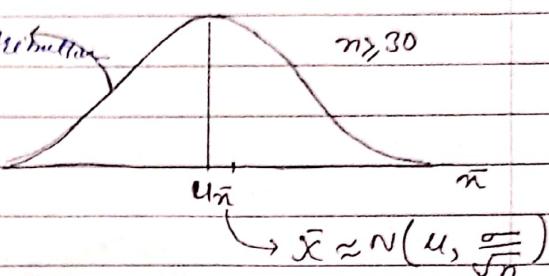
$$S_3 = \{x_9, x_{12}, x_{13}, \dots, x_{30}\} = \bar{x}_3$$

$$S_m = \{x_9, x_{12}, x_{13}, \dots, x_{30}\} = \bar{x}_m$$

Gaussian Distribution

$n \geq 30$

CLT is important for Interview



Estimate

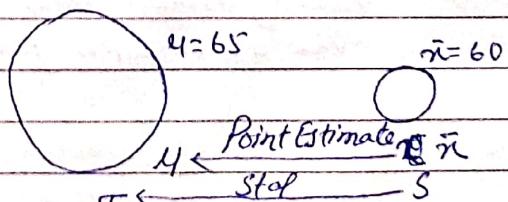
It is an observed numerical value used to estimate an unknown population parameter

- ① Point Estimate
- ② Interval Estimate

Point Estimate : Single numerical value used to estimate the unknown population parameter

* Sample mean is a point estimate of a population mean

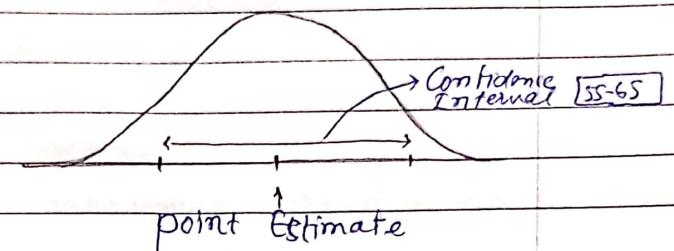
→ Some time we are not able to find population mean due to data loss (null value, NaN) etc
Then, we use estimate to solve



Problem: It may be large different between population mean and sample mean or σ and S . Then we go with Interval Estimate.

1 Internal Estimate: Range of value used to
2 estimate the unknown population parameters

3 * Internal estimate of population parameters are
4 called confidence intervals.

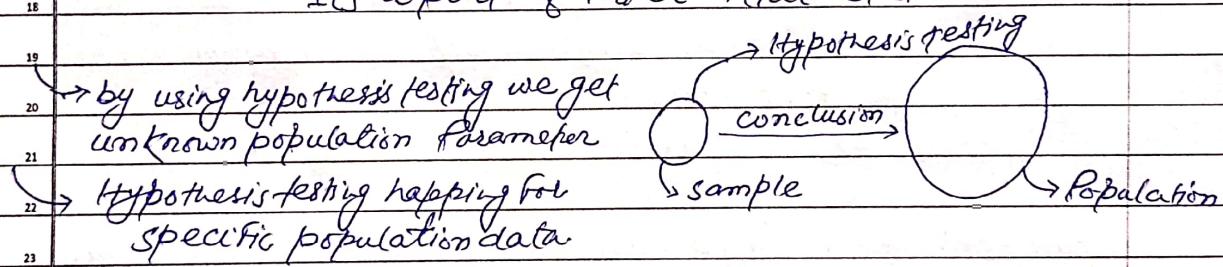


13 Hypothesis.

14 Hypothesis and hypothesis testing.

15 A hypothesis is an assumption about a
16 population parameter that is being tested using data
17 from a sample.

18 It's a part of inferential stats.



25 Hypothesis Testing Mechanism

26 Eg: Person Crime → Court

27 Step 1: Null Hypothesis (H_0) = The person is not guilty
28 - The assumption, you are beginning with.

29 Step 2: Alternate Hypothesis (H_1) = The person is guilty
30 - Opposite of null hypothesis

33 Step 3: Experiments → Proof collect (DNA Test, Finger Print),
34 Statistics Analysis using ↘

35 Step 4: Accept the null hypothesis or reject hypothesis

Eg: Colleges at District A stats it is average passed % of student are 85%. A new collage opened in the district and it was found that a sample of student 100 have a pass % of 90% with a standard deviation of 4%.

Does this have a different passed percentage.

Solution:

Null Hypothesis $H_0 = \mu = 85\%$.

Alternative Hypothesis $H_1 = \mu \neq 85\%$.

P Value (Hypothesis 3rd step)

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.

Pvalue are used in hypothesis testing to help decide whether to reject the null hypothesis

Eg: in case of your Keyboard space key



Out of 100 touches in this key $\rightarrow P=0.8$

the probability of touching in this region 80%.

Eg: Hypothesis Testing

$\text{Exp} = \text{Coin is fair or not with } 100 \text{ tossing}$

Confidence Interval

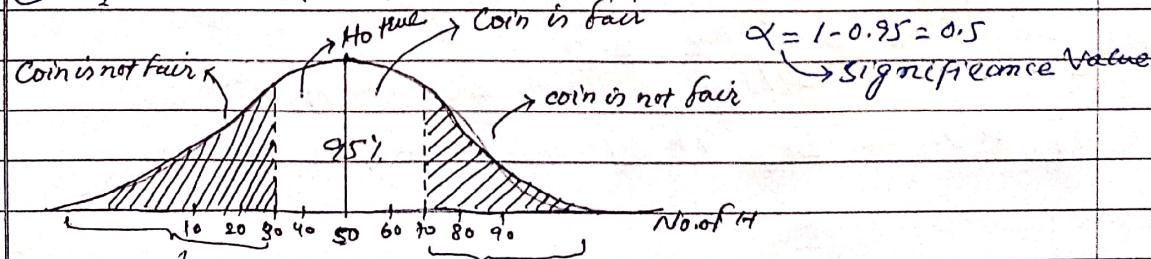
$$\chi = \{H, T\}$$

① $H_0 = \text{Coin is fair } 30 \leq P_n(H) \leq 70$

② $H_1 = \text{Coin is not fair}$

$$C.I = 0.95$$

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



Note (): \rightarrow Rejection Region \rightarrow Confidence interval
Note (): \rightarrow Rejection Region \rightarrow Confidence interval

Hypothesis Testing And Statistical Analysis

Z-table for Z-Test

Average { Z-Test }

t-table for t-Test

t-Test

Categorical Data \leftarrow Chi-Test

Variance \leftarrow ANOVA

Use Z-test { given population std }
given sample size $n \geq 30$

Eg: The average height of all height residents in a city is 168 cm with $\sigma = 3.9$. A doctor believe the mean to be different. He measured the height of 36 individual and found the average height to be 169.5 cm.

(a) State null and Alternate Hypothesis

(b) At a 95% confidence level,

Solution:

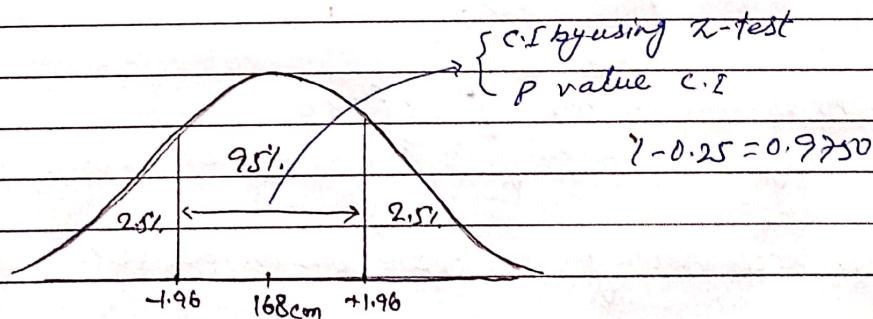
① Given $\mu = 168 \text{ cm}$, $\sigma = 3.9$, $n = 36$, $\bar{x} = 169.5 \text{ cm}$

a) Null Hypothesis $H_0 = \mu = 168 \text{ cm}$

b) Alternate Hypothesis $H_1 = \mu \neq 168 \text{ cm}$ { 2 tail Test }
if less than or greater than true H_0 rejected

c) C.I = 95% so $\alpha = 1 - 0.95 = 0.05$

Right and left $\alpha = 0.05/2$



\Rightarrow Statistical Analysis for sample data

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

$$= \frac{169.5 - 168}{3.9/\sqrt{36}} = 2.31$$

$2.31 > +1.96$ { we reject the H_0 }

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

for population data

② P-value: Says if significance value $> H_0$ reject H_0

Eg2: A factory manufactures bulbs with a average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will manufacture in less than 5 years. He test a sample of 40 bulbs and find the average time to be 4.8 year.

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

Solution: $\mu = 5$, $\sigma = 0.50$, $n = 40$, $\bar{x} = 4.8$

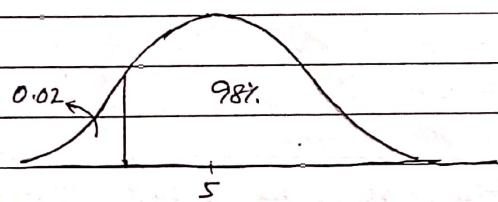
$$\textcircled{1} \quad H_0 = \mu = 5$$

$$\textcircled{2} \quad H_1 = \mu < 5 \quad \{ \text{1 tail test} \}$$

$\textcircled{3}$ Decision Boundary

Pvalue:

$$\begin{aligned} Z\text{-test} &= \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \\ &= \frac{4.8 - 5}{0.50/\sqrt{40}} \\ &= -2.53 \end{aligned}$$



IF Pvalue < significant = False Area under the curve = -2.53

$0.0570 < 0.02 \Rightarrow \text{False}$ Z-table value is = 0.0570

{ we accept H_0 }

Pvalue = 0.0570

Conclusion: The warranty needs to be revised.

Student t distribution

In statistical analysis using z-score we need (σ):
but "How do we perform an analysis when we don't know the population standard deviation?"

Solution: By using "Student t distribution"

$$Z\text{-test} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \Rightarrow t = \frac{\bar{x} - \mu}{S/\sqrt{n}}, s = \text{sample std}$$

Degree of freedom

$df = n - 1$, $n = \text{sample size}$

Note ():

Note ():

Important
Example
Visit
Points

e.g.
Sub Title:
Highlighted:

T-stats t-test (one sample t-test)

- In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A simple of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication effect intelligence? C.I = 95%.

Solution

Ans → given $\mu = 100$, $n = 30$, $\bar{x} = 140$, $\sigma = 20$

C.I = 95%, $\alpha = 5\%$.

① Null Hypothesis $H_0 : \mu = 100$

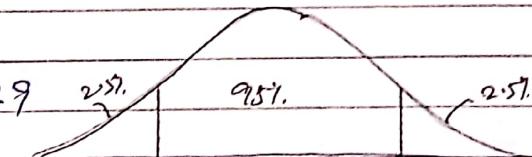
Alternate " $H_1 : \mu \neq 100$ {2 tail test}

② $\alpha = 5\% = 0.05$

③ degree of freedom

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

Decision Rule



→ If the test less than -2.045 and greater than 2.045, Reject the H_0

④ Calculate t-test statistics

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65} = 10.96$$

⑤ Decision Rules: If t is less than -2.045 and greater than 2.045, Reject the null hypothesis

$$t = 10.96 > 2.0452 \rightarrow \text{Reject the null hypothesis}$$

Answer:

Conclusion: Medication has increased the intelligence

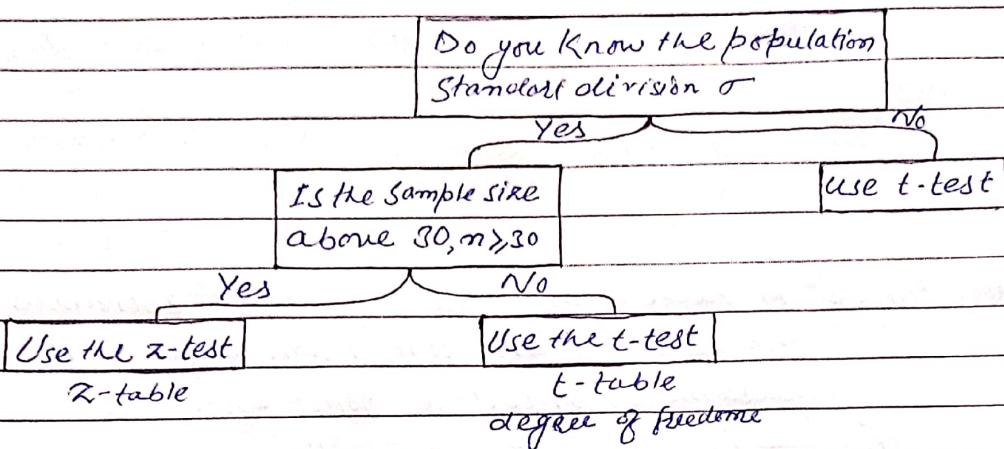
Note ():
Note ():

Important Example Visit Points
e.g.

Title: Sub Title: Highlighted:

Page No.....
Date / /

When to use t-test Vs z-test



Type 1 & Type 2 errors

		Reality	
Decision		H_0 True	H_0 False
Hypothesis	Fail to Reject H_0	✓	Type II Error
	Reject H_0	Type I Error (α)	✓

Reality: Null hypothesis is true, or null hypothesis is false

Decision: Null hypothesis is true, or null hypothesis is false

Conclusion

Outcomes 1: We reject the null hypothesis when in reality it is false → Good

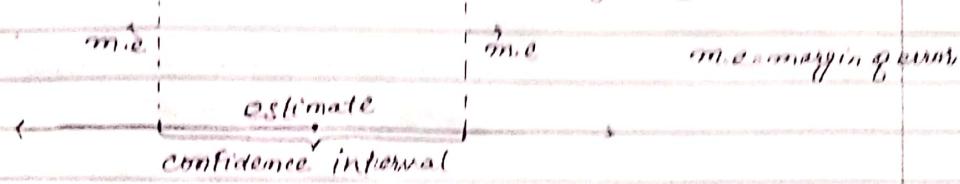
Outcomes 2: We reject the null hypothesis when in reality it is true → Type I error

Outcomes 3: We retain the null hypothesis, when in reality it is false → Type II error

Outcomes 4: We retain the null hypothesis, when in reality it is true → Good

Note ():
Note ():

Confidence Interval and Margin of Error



Point Estimate: A value of any statistics that estimate the value of an unknown population parameter is called Point estimate

$$\bar{x} \rightarrow \mu \text{ or } S \rightarrow \sigma \text{ etc.}$$

Confidence Interval: We construct a confidence interval to help estimate that the actual value of the unknown population mean is.

$$\text{Point Estimate} \pm \text{Margin of error}$$

$$Z\text{-test C.I} = \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Margin of Error

$$E = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

α = is a significant value

\bar{x} = Sample mean

Margin of Error

n = Sample size

Margin of error has based on

$Z_{\alpha/2}$ = Two tailed test

Z -test and t -test table has different values

Get value from Z -table

Eg: On the verbal section of cat exam, the σ is to be 100. A sample of 25 test taken has a mean of 520. construct 95% C.I about the mean?

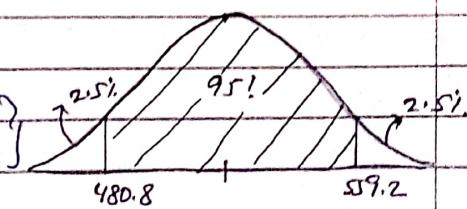
Solve,

Given $\sigma = 100$, $n = 25$, $\bar{x} = 520$, C.I = 95%, $\alpha = 5\%$.

$$Z\text{-test C.I} = \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$\text{Lower C.I} = 520 - (1.96) * \frac{100}{\sqrt{25}} = 480.8$$

$$\text{Higher C.I} = 520 + (1.96) * \frac{100}{\sqrt{25}} = 559.2$$



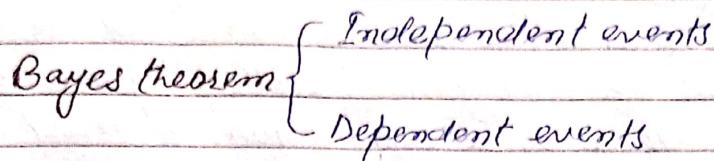
I am 95% confident that the mean CAT score lies b/w 480.8 and 559.2 confidence interval

Note ():

Note ():

Bayes theorem / Bayes Statistics

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



Independent Events

The probability of each event remains the same value whether or not the other event occurs.

Eg 1: Rolling a dice {1,2,3,4,5,6}

$$P_{\theta}(1) = \frac{1}{6}$$

$$P_{\theta}(2) = \frac{1}{6} \dots$$

Eg 2: Tossing a coin

$$P_{\theta}(H) = 0.5$$

$$P_{\theta}(T) = 0.5$$

Dependent Events

The probability of the second event is influenced by the outcomes of the first event

$$\begin{array}{c} \text{cup} \\ \text{000} \end{array} \rightarrow P_{\theta}(R) = \frac{2}{5} \quad (1)$$

after occurring this event

$$\begin{array}{c} \text{cup} \\ \text{001} \end{array} \rightarrow P_{\theta}(B) = \frac{3}{4}$$

So first event effect next event

$$* P_{\theta}(R \text{ and } B) = P_{\theta}(R) * P_{\theta}(B|R)$$

$$P_{\theta}(A \text{ and } B) = P_{\theta}(B \text{ and } A)$$

$$P_{\theta}(A) * P_{\theta}(B|A) = P_{\theta}(B) * P_{\theta}(A|B)$$

considering B event happened after A event

Bayes theorem

$$P(B/A) = \frac{P_{\theta}(B) * P_{\theta}(A|B)}{P_{\theta}(A)}$$

Conditional Probability.

Parameters:

↓ it is also written

$$P_{\theta}(A|B) = \text{Probability of } A \text{ given } B$$

A given B is True (happened)

$$P_{\theta}(B)$$

A event happened when B event

$$P_{\theta}(B|A) = \text{Probability of } B \text{ is given } A \text{ is true (happened)}$$

$P_{\theta}(A), P_{\theta}(B)$ = They are independent probability of A and B

In machine learning:

Dataset

Size of House	No. of rooms	Location	Price
---------------	--------------	----------	-------

x_1

x_2

x_3

y

↳ independent

↳ dependent

$$P_{\theta}(y|x_1, x_2, x_3) = P_{\theta}(y) * P_{\theta}(x_1, x_2, x_3|y)$$

$$P_{\theta}(x_1, x_2, x_3)$$

Bayes theorem

Note ():

Note ():

Important
Example
Visit
Points
e.g.
⇒

Title:
Sub Title:
Highlighted:

Page No.....
Date / /

CHI SQUARE TEST

The chi square test for goodness of fit test claims about population proportions. (categorical variance).

It is a non parametric test that is performed on categorical (ordinal, nominal) data.

Eg: There is a population of men who likes different colors of bike

	Theory	Sample
Yellow bike	$\frac{1}{3}$	22
Orange bike	$\frac{1}{3}$	17
Red bike	$\frac{1}{3}$	59

Theory about population → Observed categorical distribution
 Theory categorical distribution → distribution

Observed word used for sample data set

Goodness of fit test

(i) In a student class of 100 students, 30 are Right handed. Does this class fit the theory 12% of people are right handed.

	Observation	Expected
Right handed	30	12
Left handed	70	88
Sample	100	100

Theory Categorical Distribution

CHI SQUARE For Goodness of Fit

In 2010 census of the city, the weight of the individuals in a small city were found to be the following

Theory about population	<50 kg	50-75	>75
	20%	30%	50%

In 2020, weight of $n=500$ individuals were sampled. Below are the results

Observed / Sample data	<50	50-75	>75
	140	160	200

Note ():
Note ():

Using $\alpha = 0.05$, what you conclude the population difference of weight has changed in the last 10 years?

Solution:

2010	<50kg	50-75	>75
Expected	20%	30%	50%

2020	<50	50-75	>75
Observed	140	160	260
<i>m=200</i>			

	<50	50-75	>75
Expected	0.2×200	0.3×200	0.5×200
	=100	=150	=250

- ① Null Hypothesis H_0 : The data meets the expectations
 Alternative Hypothesis H_1 : The data does not meet the expectation

② $\alpha = 0.05$ and C.I = 0.95

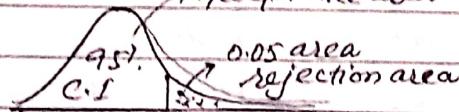
③ degree of freedom

$$df = K - 1 = 3 - 1 = 2$$

* K = numbers of categories

④ Decision Boundary

→ Acceptance area



CV = Critical value

χ^2 = Chi square notation

If χ^2 is greater than 5.99, Reject H_0 .
 Else:

We fail to reject the null hypothesis

⑤ Calculate Chi-square test Statistics

$$\chi^2 = \sum \frac{(O_{\text{observed}} - E_{\text{expected}})^2}{E_{\text{expected}}}$$

Check Chi-square table to find value

by using degree of freedom ($df = 2$) and

$\alpha = 0.05$ value

$$\begin{aligned} &= \frac{(140 - 100)^2}{100} + \frac{(160 - 150)^2}{150} + \frac{(260 - 250)^2}{250} \\ &= \frac{1600}{100} + \frac{100}{150} + \frac{100}{250} = 16 + 0.66 + 1 \\ &= 26.66 \end{aligned}$$

$$\chi^2 = 26.66 \quad (\text{Chi square})$$

Now $\chi^2 = 26.66 > 5.99$, Reject H_0 .

Answer: The weight of 2020 population are different than those expected in the 2010 population.

Note { }:
 Note { }:

F-Test (Variance Ratio Test)

① Following data shows the number of bulbs produced daily for same days by 2 workers A and B

A B

40 39

30 38

38 41

41 33

38 32

35 39

40

34

Can we consider based on
the data worker B is more
stable and efficient

$$\alpha = 0.05$$

Solution

① Null Hypothesis : $H_0 \Rightarrow \sigma_1^2 = \sigma_2^2$

Alternate Hypothesis : $H_1 \Rightarrow \sigma_1^2 \neq \sigma_2^2$

② calculate the variance

A

x_i \bar{x} $(x_i - \bar{x})^2$

40 37 9

30 37 49

38 37 1

41 37 16

38 37 1

35 37 4

B

x_i \bar{x}_2 $(x_i - \bar{x}_2)^2$

39 37 4

38 37 1

41 37 36

33 37 16

32 37 25

39 37 4

40 37 9

34 37 9

$$S_1^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

$$S_2^2 = \sum_{i=2}^n (x_i - \bar{x}_2)^2 = \frac{84}{7} = 12$$

$$S_1^2 = \frac{80}{5} = 16$$

→ Calculating of variance Ratio F-test

$$F = \frac{S_1^2}{S_2^2} =$$

③ Decision Rule

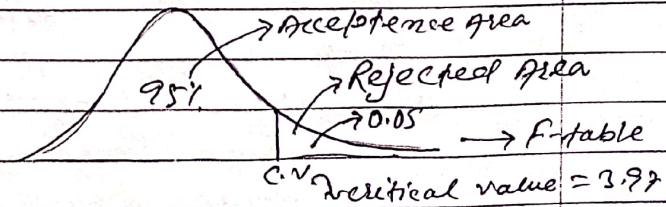
$$df_1 = 6-1 = 5$$

$$df_2 = 8-1 = 7$$

$$\alpha = 0.05$$

Note () :

Note () :



If F-test is greater than 3.77, Rejected the null hypothesis H_0 .
 $1.33 < 3.77$ we fail to reject the null hypothesis.
 Conclusion: Worker B is not efficient when worked to worker A.

ANOVA (Analysis of Variance)

Definition: ANOVA is a statistical method used to compare the means of 2 or more groups.

Anova two important things

- (i) Factors (variable)
- (ii) levels

Eg1: Factor = Medicines

Level = 5mg 15mg 20mg (Dosage)

Eg2:

Mode of Payment

Level → GPay PhonePE IMPS NEFT

ASSUMPTIONS IN ANOVA

(i) Normality of Sampling Distribution of means.
 The distribution of sample mean is normally distributed.

(ii) Absence / Removal of outliers

Outlying score need to be removed from dataset.

(iii) Homogeneity

Each and every one of the population that (Sample/column) used in level has same variance

$$(\sigma_1^2 = \sigma_2^2 = \sigma_3^2)$$

Note ():

Note ():

Important Example Visit Points
e.g. Sub Title: Highlighted:

Page No.....
Date / /

population variance in different levels of each independent variable are equal.

(iv) Samples must be independent and random

Types of ANOVA

(i) One Way ANOVA: One factor with at least 2 levels. These levels are independent.

Eg: Doctors want to test a new medication to decrease headache.

They split the participant in 3 conditions.
(5mg, 10mg, 15mg).

Doctors ask the participant to rate the headache (1-10),

levels	Medication	factor (single)
5mg	10mg	15mg
5	7	2
9	8	7
-	-	-
-	-	-

all levels are independent

(ii) Repeated Measures ANOVA: One factor with at least 2 levels, levels are dependent

levels	Day 1	Day 2	Day 3
Day 2 dependent	8	5	6
Day 1 etc	7	4	3
	9	8	7

(iii) Factorial ANOVA: Two or more factors, each of which at least 2 levels, levels can be either independent and dependent.

Independent level	Day 1	Day 2	Day 3	levels dependent
Gender	Man 8	5	6	
	7	4	3	
	woman 6	5	4	

Note ():

Note ():

Hypothesis Testing for ANOVA

Null Hypothesis $H_0 : \mu_1 = \mu_2 = \mu_3$

Alternate Hypothesis $H_1 : \text{At least one of the mean is not equal}$

Test Statistics

$$F = \frac{\text{Variance between Sample}}{\text{Variance within Sample}}$$

X_1	X_2	\bar{X}_3	Variance between Sample
1	6	5	$H_0 : \bar{X}_1 = \bar{X}_2 = \bar{X}_3$
2	7	6	$H_1 : \text{At least one sample mean is not equal}$
4	3	3	
5	2	2	
3	1	4	

$$\sum X_1 = 15 \quad \sum X_2 = 19 \quad \sum X_3 = 20$$

$$\bar{X}_1 = 3 \quad \bar{X}_2 = 4 \quad \bar{X}_3 = 4$$

One way ANOVA

(1) Doctors want to test a new medication with reduce headache. They split the participant into 3 condition (15mg, 30mg, 45mg). After on the doctor ask the patient to refer the headache between (1-10). Are there any difference between 3 conditions using $\alpha = 0.05$?

Solution:

	15mg	30mg	45mg
9	7	4	
8	6	3	
7	6	2	
8	7	3	
8	8	4	
9	7	3	
8	6	2	

(1) Define H_0 and H_1 ,

$H_0 : \mu_{15} = \mu_{30} = \mu_{45}$ and $H_1 : \text{not all } \mu_i \text{ are equal}$

Note ():
Note ():

Important Example Visit Points
e.g. $\oplus \rightarrow$
Title:
Sub Title:
Highlighted:

② State significant value.

$$\alpha = 0.05 \text{ and } C.I = 0.95$$

③ calculate Degree of Freedom

$$N = 21 \quad a = 3 \quad n = 7$$

$$df_{\text{between}} = a - 1 = 3 - 1 = 2$$

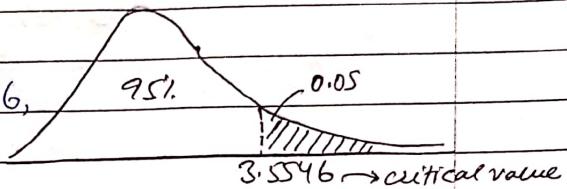
$$df_{\text{within}} = N - a = 21 - 3 = 18 \quad (2, 18)$$

$$df_{\text{total}} = N - 1 = 21 - 1 = 20 \quad F \text{ table}$$

④ State Decision Rule

If F is greater than 3.5546,

reject the H_0 .



⑤ Calculate Test statistics

SS_{between} SS_{within} SS_{total}

$$\textcircled{1} \quad SS_{\text{between}} = \frac{\sum (\sum a_i)^2 - T^2}{n} - \frac{N}{n}$$

$$SS_{\text{between}} = 9 + 8 + 7 + 8 + 7 + 8 + 8 = 57$$

$$SS_{\text{within}} = 7 + 6 + 6 + 7 + 8 + 7 + 6 = 47$$

$$SS_{\text{total}} = 4 + 3 + 2 + 3 + 4 + 3 + 2 = 21$$

$$= \frac{57^2 + 47^2 + 21^2}{7} - \frac{(57 + 47 + 21)}{21}$$

$$= 98.67$$

$$\textcircled{2} \quad SS_{\text{within}} = \sum y^2 - \frac{\sum (\sum a_i)^2}{n}$$

$$= \sum y^2 - \left(\frac{57^2 + 47^2 + 21^2}{7} \right)$$

$$\sum y^2 = 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + 8^2 + 7^2 + 6^2 \\ + \dots$$

$$= 853$$

$$= 853 - \left(\frac{57^2 + 47^2 + 21^2}{7} \right) = 10.29$$

$$\textcircled{3} \quad SS_{\text{total}} = \sum y^2 - \frac{T^2}{N}$$

$$= 853 - \frac{125^2}{21} = 108.95$$

Note ():
Note ():

	SS	df	MS	F
Between	98.67	2	49.34	
Within	10.29	18	0.54	
Total	108.95	20		

$$F = \frac{MS_{\text{between}}}{MS_{\text{within}}}$$

$$F = \frac{49.34}{0.54} = 86.56$$

If F is greater than 3.5546, reject the H_0 .
 $86.56 > 3.5546$, Reject the null hypothesis