

Statistics for Data Science

2023



By

Santosh Lohar

IBM Data Science Professional Certified

[LinkedIn](#)



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date:

statistics for

Data scientist / Data Analyst / Business Intelligence tools /

What is Statistics

→ Statistics is the science of collecting, organizing and analyzing data for better decision making

Data → facts or piece of information that can be measured in a quantitative way.

Eg: The IQ of the class:

$$\{92, 91, 90, 76, 75, 60\}$$

Eg: Ages of students of class.

$$\{30, 21, 21, 25, 29, 16, 15\}$$

Types of Statistics

Descriptive Statistics

- It consists of organizing and summarizing data.

Inferential Statistics

Technique where we used the data that we have measured

to form conclusions.

Eg: Class room of maths students marks of the 1st Sem

$$\{84, 72, 94, 59, 89, 90, 91\}$$



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

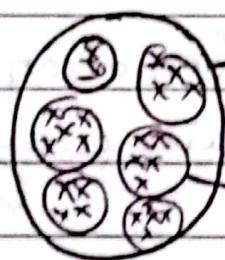
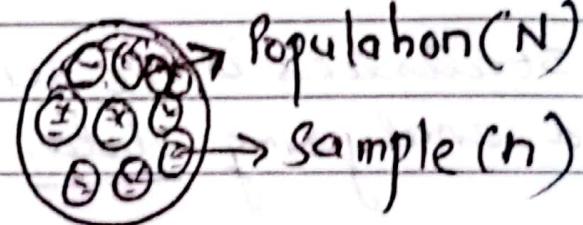
Date:

Descriptive Statistics

- What is the average marks of the students in the class? (from previous example data.)

Inferential Statistics

Are the marks of the students of this classroom similar to the age of the maths classroom in the College?



→ Total population (N)

→ Sample (taken) to make decision

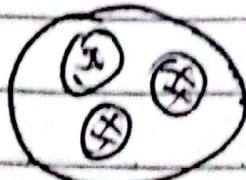
like election, exit Poll.

Population (N)

Sample(n)

Sampling Techniques

- ① → Simple Random Sampling:- Every member of the population (N) has an equal chance of being selected for your sample (n)





Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date:

② Stratified Sampling :- Where the Population (N) is split into non-overlapping groups (strata)

eg:

Gender: male

: female

Age-group: (0-10) (10-20) (20-40) (40-100)

Programming: Dot net ← PHP → Python

Not-overlapped

③ Systematic Sampling:
 $(N) \rightarrow n^{\text{th}}$ individual

Thanos

eg: Mall \rightarrow Survey (covid)
 \hookrightarrow 8th person \rightarrow Survey.

④ Convenience Sampling (only those people)

\hookrightarrow Survey \rightarrow Data Science

RBI (House-hold Survey)

\hookrightarrow Women

eg: Drug Tested



Mo Tu We Th Fr Sa Su

Date : _____

Variables:

A variable is a property that can take on any value: {182, 175, 172, 170, 165} Height
eg: Weight: {50, 55, 60, 70, 80, 100 - }

types of variables

<u>Quantitative Variable</u>	<u>Qualitative Variable</u> (categorical variable)
✓	

measured numerically
{add, sub, mul, div} → Based on some characteristics we can derive categorical variables

Eq: \overline{PQ}

0-10

10-50

50-100

Blood Group

less: 10

medium, 10

Find 10

A⁺ve (O⁻)

P-Shū-size

L

XL

XXL

m

5

Quantitative

Discrete Variable

Eg: Whole number
No. of Bank A/c

Eg: {2, 3, 4, 5, 6, 7, ...}

Total number of children

Total no. of family.

Continuous Variables

Eg: Height : 172.5, 170.5, 168.2

Weight : 100kg, 77.5kg, 98.25

Rainfall : 1.1, 2.3, 1.25, 1.35

Eg: What kind of variable Gender is? Categorical
" " " " marital status? (categorical)

River length: Continuous

Pin Code: { discrete or categorical }

Variable Measurement Scales:

4 types of measured variable.

- ① Nominal data = { Categorical data : colors }
- ② Ordinal data : order of the data matters value no.
- ③ Interval : order matters value also matters (zero no.)
- ④ Ratio :

Eg: students (marks)

100

98

57

85

Rank

1. } ordinal

2 } data

4

3



Mo Tu We Th Fr Sa Su

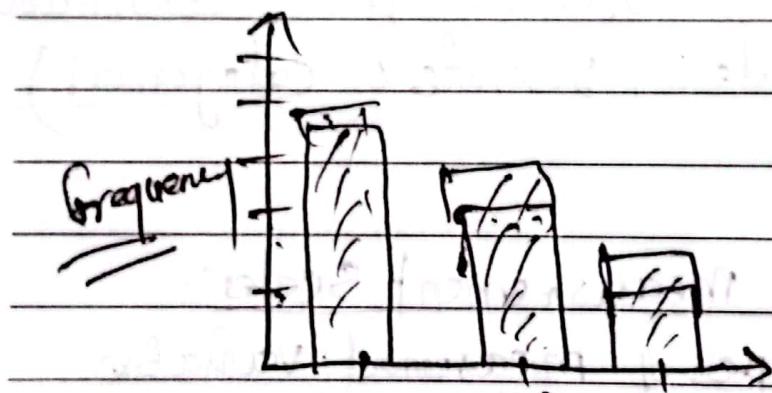
Date:

frequency Distribution:

Sample dataset: Rose, lily, Sunflower, Rose, lily, Sunflower, Rose, lily, Rose.

flowers	frequency	Cumulative freq
Rose	4	4
lily	3	7
Sunflower	2	9

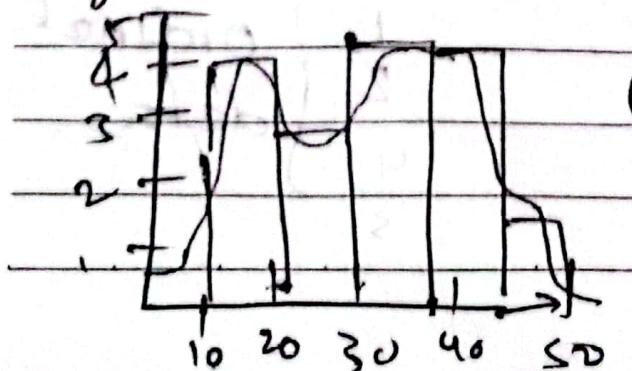
① BAR GRAPH:



Rose lily Sunflower

② HISTOGRAM - Continuous

Ages: {10, 12, 15, 16, 18, 24, 26, 35, 36, 37, 40, ...}



Bins

default bins = 10



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date:

11/
Lincolin

Measure of Central Tendency:

* Arithmetic mean for Population & Sample.

$$\text{Average} \quad \frac{\sum_{i=1}^N (x_i)}{N} \quad \frac{\sum_{i=1}^n (x_i)}{n}$$

eg: $\frac{1+1+2+2+3+3+4+5+5+6}{10}$

$$M = \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$M = \frac{\sum_{i=1}^N (x_i)}{N}$$

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

$$\bar{x} = \frac{32}{10}$$

$$\bar{x} = 3.2$$

$$M = 32/10 = 3.2$$

Central Tendency:

① Mean ② median ③ mode

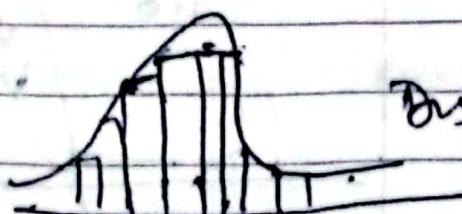
↳ Refers to the measures used to determine the center of the distribution of data.

eg: { 1, 1, 2, 2, 3, 3, 4, 5, 5, 6 } \rightarrow outlier.

Mean after outlier

$$\begin{aligned} \text{Mean} &= \frac{32 + 100}{11} \\ &= 12 \end{aligned}$$

$$\begin{aligned} M &= 3.2 \quad \text{Previous} \\ &\quad \text{After outlier 100} \\ \text{Mean} &= 12 \end{aligned}$$



Distribution



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date :

Median :

$$\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 112\}$$

Step -1 : Sort the numbers .

if odd number mid number will result

if even number take average between two numbers .

Measure of Dispersion

↓
(Spread)

↳ Variance (σ^2)↳ Standard deviation (σ)* Variance :

Population Variance

Sample Variance

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

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

eg:

$$\begin{array}{cccc} x & \mu & x - \mu & (x - \mu)^2 \\ 1 & 2.83 & -1.87 & 3.34 \end{array}$$

$$\begin{array}{cccc} 2 & 2.83 & -0.83 & 0.6869 \end{array}$$

$$\begin{array}{cccc} 2 & 2.83 & 0.83 & 0.6869 \end{array}$$

$$\begin{array}{cccc} 3 & 2.83 & 0.17 & 0.03 \end{array}$$

$$\begin{array}{cccc} 4 & 2.83 & 4.17 & 1.37 \end{array}$$

$$\begin{array}{cccc} 5 & 2.83 & 2.17 & 4.71 \end{array}$$

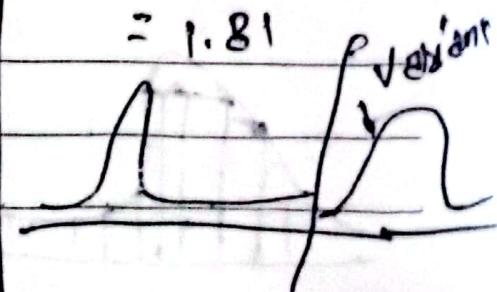
$$\boxed{\mu = 2.83}$$

$$\frac{10.84}{6}$$

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

$$= 10.84$$

$$= 1.81$$



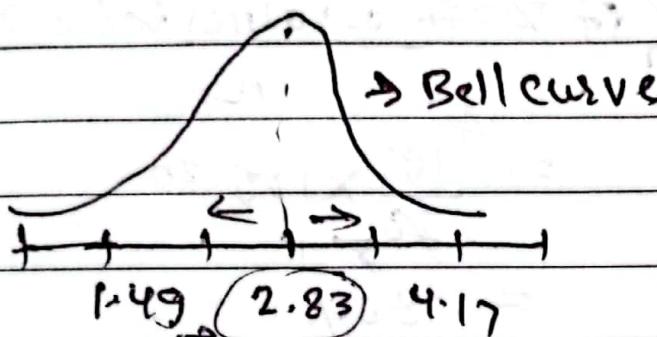


Mo Tu We Th Fr Sa Su

Date :

$$\text{Standard deviation}(\sigma) = \sqrt{\text{Variance}} = \sqrt{\sigma^2}$$

$$= \sqrt{1.81} = 1.345$$



$$\begin{array}{r} 2.83 \\ 1.34 \\ \hline 4.17 \end{array} \quad \begin{array}{r} 2.83 \\ -1.34 \\ \hline 1.49 \end{array}$$

\Rightarrow 1.5 σ from the mean

Percentiles and Quartiles (find median)

Percentage: 1, 2, 3, 4, 5

\therefore % of the numbers that are odd?

\therefore # of numbers that are odd
Total numbers,

$$= 3/5 = 0.6 = 60\%$$

Percentile: A Percentile is a value below which a certain percentage of observations lie.

eg:-

Defacto: 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 12

What is the Percentile ranking of 10?

here $n = 20$



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Date:

from the given Problem:

$$x = 10$$

Now, Percentage rank of $x = \frac{\# \text{ of values below } x}{n} \times 100$

$$= \frac{16}{20} \times 100 \\ = 80\%$$

Q)

80% of the entire distribution is less than 10

find the percentile rank of 11?

$$\text{here, } x = 11$$

$$n = 20$$

Now! Percentile Rank of $x_i = \frac{\# \text{ of values below } x_i}{n} \times 100$

$$= \frac{17}{20} \times 100 \\ = 85\%$$

Hence, 85% of the entire distribution is less than 11

* What value exists at percentile ranking of 25%?

$$\text{Value} = \frac{25}{100} \times (n+1)$$

$$= \frac{25}{100} \times (21) = 5.25 \rightarrow \text{Index Position.}$$

Locate value at the distribution $\rightarrow \textcircled{5} \text{ value}$



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Date:

❖ Five Number Summary ❖

- ↳ minimum
- ↳ first quartile (Q_1)
- ↳ Median
- ↳ Third quartile (Q_3)
- ↳ maximum.

Removing Outlier

e.g.: $\{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 27\}$ ↗ outlier

Lower fence \longleftrightarrow higher fence

$$\text{lower fence} = Q_1 - 1.5(\text{IOR})$$

$$\text{Higher fence} = Q_3 + 1.5(\text{IOR})$$

$$\left. \begin{array}{l} Q_3 = 75\% \\ Q_1 = 25\% \end{array} \right\}$$

$$\boxed{\text{Interquartile Range (IOR)} = Q_3 - Q_1}$$

$$\text{Value of } 25\% = \frac{25}{100}(n+1) = \frac{25}{100}(19+1) = 5 \text{ (index)}$$

$$= 3$$

$$\text{Value of } 75\% = \frac{75}{100}(19+1) = \frac{75}{100} \times 20 = 15 \text{ (index)}$$

$$= 7$$



Now,

$$\begin{aligned}\text{Interquartile Range (IQR)} &= Q_3 - Q_1 \\ &= 7 - 3 \\ &= 4\end{aligned}$$

$$\begin{aligned}\text{Again, lower fence} &= Q_1 - 1.5 \times (\text{IQR}) \\ &= 3 - 1.5(4) \\ &= 3 - 6 \\ &= -3\end{aligned}$$

$$\begin{aligned}\text{Higher fence} &= Q_3 + 1.5(\text{IQR}) \\ &= 7 + 1.5(4) \\ &= 13 \\ &\equiv\end{aligned}$$

$$\begin{bmatrix} \text{Lower fence} & \longleftrightarrow & \text{Higher fence} \end{bmatrix}$$

$$[-3 \quad \longrightarrow \quad 13]$$

* From the above distribution remove the elements as per lower & higher fence calculation value. So, distribution elements lies which is between -3 & 13 fence.

New distribution

Remaining { 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 27 }
- data

{ removed }

$$\text{Median} = \frac{5+5}{2} = 10 / 2 = 5$$

IN : SANTOSH - LOHAR



Mo Tu We Th Fr Sa Su

Date :

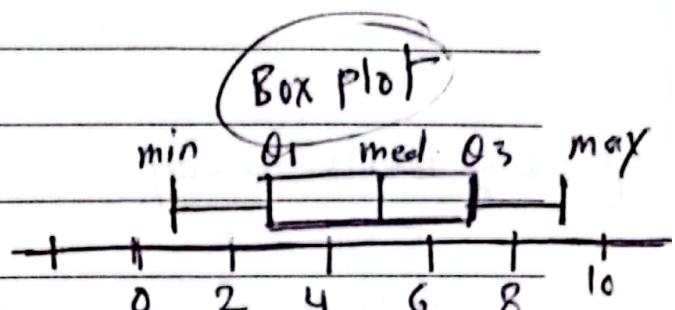
Minimum = 1

Q₁ = 3

median = 5

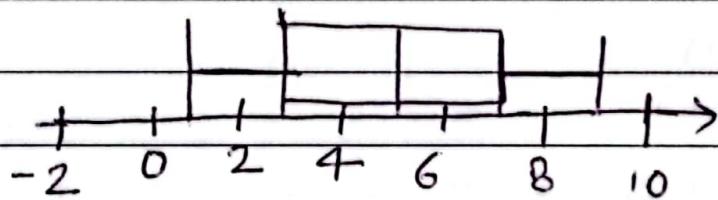
Q₃ = 7

maximum = 9



Box Plot

Box plot is used to determine Outlier



Application : Data Visualization . (Boxplot)

Why Sample Variance is $n-1$?

{



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date:

Advance statistics

① Distribution

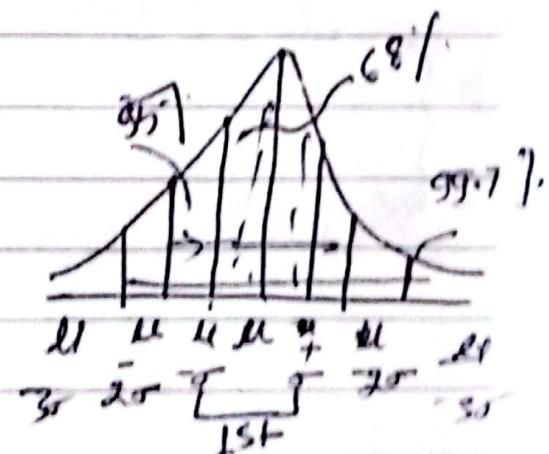
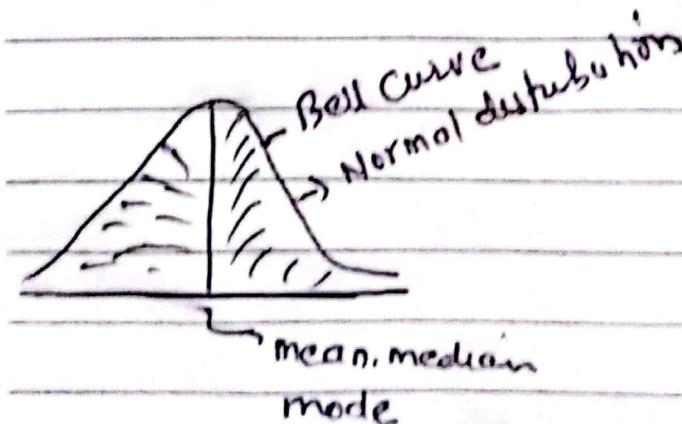
- Normal distribution
- Standard Normal distribution.
- Z score
- Log normal distribution.
- Bernoulli's distribution.
- Binomial distribution.

② Distribution

e.g.: Ages = { 24, 26, 28, 29, 20, 32, ... }



* Gaussian / Normal distribution



68% dataset distribution lies in 1st standard deviation.

Empirical formula.

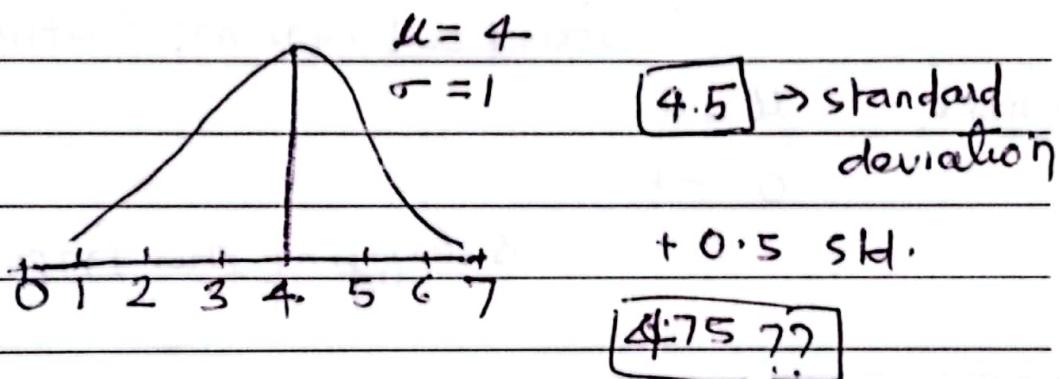
(68% - 95% - 99.7% Rule)

1st standard
deviation
Values

2nd standard
deviation
Values

3rd standard
deviation
Values

eg:



$$Z\text{-score} = \frac{x_i - \mu}{\sigma} = \frac{4.75 - 4}{1} = 0.75 \text{ SD}$$

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

$$\text{Find } z\text{-score}, z = \frac{x_i - \mu}{\sigma}$$

Here,

$$\mu = 4, \sigma = 1$$

$$\text{Now, } z(1) = \frac{1-4}{1} = -3 \quad | \quad z(6) = \frac{6-4}{1} = 2$$

$$z(2) = \frac{2-4}{1} = -2 \quad | \quad z(7) = \frac{7-4}{1} = 3$$

$$z(3) = \frac{3-4}{1} = -1$$

$$z(4) = \frac{4-4}{1} = 0$$

$$z(5) = \frac{5-4}{1} = 1$$

Now, after z -score

new dataset will be

$$z = \{-3, -2, -1, 0, 1, 2, 3\}$$



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

Where, $\mu = 4$ $\sigma = 1$ Normal Distribution

After z-score calculation

dataset - z-score = $\{-3, -2, -1, 0, 1, 2, 3\}$

Standard Normal Distribution

Where, $\mu = 0$

$\sigma = 1$

Satisfying this Property.

Practical Application:

DATASET: unit (standardization) z-score

Age (yr) Salary (Rs.) Weight (kg)

24 40K 60

25 50K 70

26 60K 49

27 45K 59

Normalization

$\{\mu=0, \sigma=1\}$

<

(-1 to 1)

(0 to 1)



Practical eg:
=

IND vs SA {ODI Series}

ODI Series 2021 (CRICKET)

Scores Average 2021 = 250

Standard deviation = 10

Team final score = 240

Q1

Compare the both the series in which year Rishabh final score

was better?

ODI Series 2020 (CRICKET)

Series Avg 2020 = 260

standard deviation = 12

Team final score = 245

Now, 2021

$$z\text{-score} = \frac{\mu_1 - \mu}{\sigma} = \frac{240 - 250}{10} = -1$$

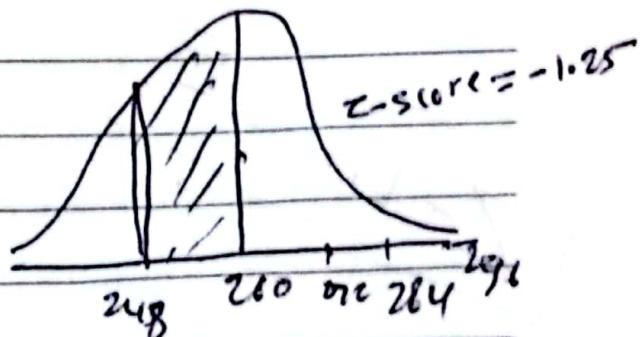
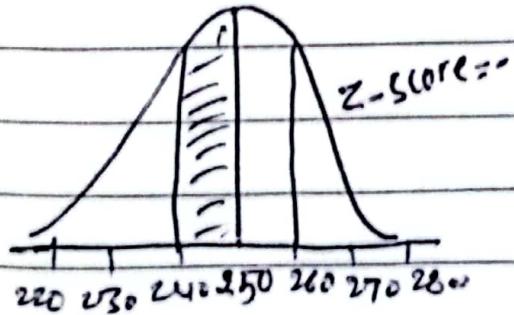
2020

$$z\text{-score} = \frac{\mu_1 - \mu}{\sigma} = \frac{245 - 260}{12} = -1.25$$

In 2021:

$$\mu = 250, \mu_1 = 240, \sigma = 10$$

$$\mu = 260, \mu_1 = 245, \sigma = 12$$



IN

Santosh - Lohar

Mo Tu We Th Fr Sa Su

Date _____

Assignment

What Percentage of scores falls above 4.25 ?

$$\mu = 4$$

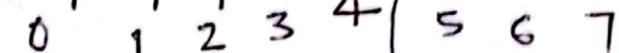
$$\sigma = 1$$

Area of Body

0.59871

Tail

1 - Left Area



4.25

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

Locate z-table

Now,

$$\begin{aligned}
 \text{Right Area (Tail)} &= 1 - \text{left area} \\
 &= 1 - 0.59871 \\
 &= 0.4013 \\
 &= \underline{\underline{40\%}}
 \end{aligned}$$



→ PRACTICAL



Mo Tu We Th Fr Sa Su

Date:

PRACTICAL

①

```
import seaborn as sns
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
%matplotlib inline
```

} import
} Python
} Library

②

```
import statistics
```

③

```
df = sns.load_dataset('tips')
```

④

```
df.head()
```

	total_bill	tip	sex	smoker	day	time	size
--	------------	-----	-----	--------	-----	------	------

0	16.99	1.01	female	No	Sun	Dinner	2
---	-------	------	--------	----	-----	--------	---

1	10.34	1.66	male	No	Sun	Dinner	3
---	-------	------	------	----	-----	--------	---

2	21.01	3.50	male	No	Sun	Dinner	3
---	-------	------	------	----	-----	--------	---

3	23.68	3.31	male	No	Sun	Dinner	2
---	-------	------	------	----	-----	--------	---

4	24.59	3.61	female	No	Sun	Dinner	4
---	-------	------	--------	----	-----	--------	---

5

⑤

```
np.mean(df['total_bill'])
```

⇒ 19.7959...

⑥

```
np.median(df['total_bill'])
```

⇒ 17.795

⑦

```
statistics.mode(df['total_bill'])
```

⇒ 13.42

IN

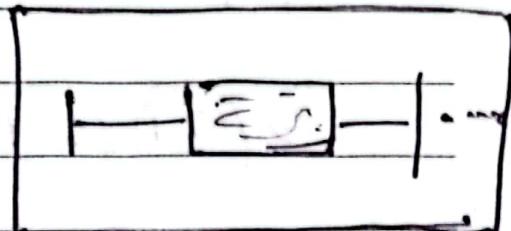
Mo Tu We Th Fr Sa Su

SANTOSH LOHAR

Date:

Box plot (import seaborn as sns)

sns.boxplot(df['total_bill'])



sns.histplot(df['total_bill'])

sns.histplot(df['sepal_length'], kde=True)

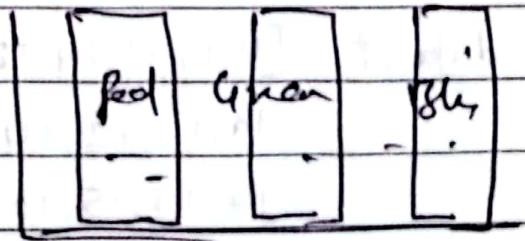
↓
sepal_width

Before 'sepal_length' load data.

df = load_dataset('iris')

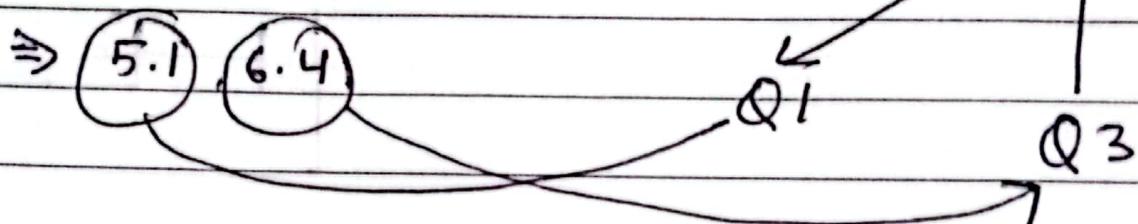
BarGraph → countplot

sns.countplot(df['species'])



Percentile

↳ np.percentile(df1['sepal-length'], [25, 75])



Note : we can change parameter as we need.

IQR

* Outliers

- ① import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline.] import python library.
- ② dataset = [11, 10, 12, 14, 13, 15, 14, 13, 15, 102, 12, 14, 17, 19, 107, 10, 13, 12, 14, 12, 108, 12, 11, 14, 13, 15, 10, 15, 12, 10, 14, 13, 15, 10]] → Define our dataset



Mo Tu We Th Fr Sa Su

Date:

Find outlier | z-score

③

Outliers = []

def outlier_find(data):

threshold = 3

mean = np.mean(data)

std = np.std(data)

for i in data:

z-score = (i - mean) / std

if np.abs(z-score) > threshold:

outliers.append(i)

return outliers

④

outlier_find(dataset)

↳ [102, 107, 108]

Python Code

IQR → Interquartile Range

- ① ↳ sort the particular dataset
- ② ↳ calculate Q_1 & Q_3
- ③ ↳ $IQR = Q_3 - Q_1$
- ④ ↳ find lower fence ($Q_1 - 1.5(IQR)$)
- ⑤ ↳ find higher fence ($Q_3 + 1.5(IQR)$)

dataset = sorted(dataset)

dataset

[Display Dataset]

- # $Q_1, Q_3 = np.percentile(dataset, [25, 75])$
Print(Q_1, Q_3)

⇒ 12.0, 15.0

$\overbrace{Q_1}$ $\overbrace{Q_3}$

find lower & higher fence

$$\# \text{lower-fence} = Q_1 - 1.5(IQR) = 7.5$$

$$\# \text{higher-fence} = Q_3 + 1.5(IQR) = 19.5$$

Where, $IQR = Q_3 - Q_1 \Rightarrow 3$



Mo Tu We Th Fr Sa Su

Date:

import seaborn as sns

sns.boxplot(dataset)

[Box Plot -]

Probability

Probability is a measure of the likelihood of an event.

e.g.

Roll a dice $\{1, 2, 3, 4, 5, 6\}$

$$\text{Pr}(6) = \frac{1}{6} = \frac{\# \text{ of ways an event can occur}}{\# \text{ of possible outcome}}$$

e.g. Toss a coin $\{H, T\}$

$$\text{Pr}(H) = \frac{1}{2}$$

■ Addition Rule (Probability, "or")

* Mutual Exclusive Event *

→ Two events are mutual exclusive if they can not occur at the same time

e.g. Rolling a dice $\{1, 2, 3, 4, 5, 6\}$

e.g. Tossing a coin $\{H, T\}$

* Non-Mutual Exclusive Event *

→ Multiple events can occur at the same time

e.g. Deck of cards $\{Q, Q\}$

↳ red color

IN

SANTOSH LOHAR

Mo Tu We Th Fr Sa Su Date

eg. If I toss a coin, what is the probability of the coin landing on heads or tails?

→ Mutual exclusive (Addition Rule we called)

$$\begin{aligned} \Pr(A \text{ or } B) &= \Pr(A) + \Pr(B) \\ &= \frac{1}{2} + \frac{1}{2} \\ &= 1 \end{aligned}$$

$\therefore \Pr(A \text{ or } B) = 1$

eg. Roll a dice {1, 2, 3, 4, 5, 6}

$$\begin{aligned} \Pr(1 \text{ or } 3 \text{ or } 6) &= \Pr(1) + \Pr(3) + \Pr(6) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ &\stackrel{+}{=} \frac{3}{6} = \frac{1}{2} = 0.5 \end{aligned}$$

eg. Non-mutual exclusive

↳ you are picking a card randomly from a deck
what is the probability of choosing a card that is queen or a heart?

⇒ Non-mutual exclusive

$$\Pr(Q) = \frac{4}{52}, \Pr(\heartsuit) = \frac{13}{52}, \Pr(Q \text{ and } \heartsuit) = \frac{1}{52}$$

\hookrightarrow multiplication

Addition Rule for non-mutual exclusive events.

$$\begin{aligned} \Pr(A \text{ or } B) &= \Pr(A) + \Pr(B) - \Pr(A \text{ and } B) \\ &= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \\ &= \frac{16}{52} = \frac{4}{13} \end{aligned}$$



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date:

Multiplication Rule

↳ Independent Events

e.g. Rolling a dice $\{1, 2, 3, 4, 5, 6\}$

Each & every are independent

It may be $\{1, 1, 2, \dots\}$

↳ Dependent Events

e.g.

<input type="checkbox"/>				
<input type="checkbox"/>				

$Pr(\bullet) = 4/9$

$Pr(0) = 5/9$

Now, After pickup/taken ^{out} a ball (\bullet) from bag

<input type="checkbox"/>				
<input type="checkbox"/>				

$Pr(\bullet) = 4/8$

$Pr(0) = 4/8$

Independent Events

e.g.

What is the probability of rolling a "5" and then a "4" in a dice?

Ans \Rightarrow Independent event

Multiplication Rule

$$\begin{aligned} Pr(A \text{ and } B) &= P(A) \times P(B) \\ &= 1/6 \times 1/6 \\ &= 1/36 \end{aligned}$$



eg. What is the probability of drawing a Queen and then a Ace from a deck of cards?

Ans: Dependent

$$P(A \text{ and } B) = P(A) \times P(B/A)$$

conditional probability.

$$= \frac{4}{52} \times \frac{4}{51}$$

Permutation and Combination



eg: School trip {chocolate factory} \rightarrow {Dairy milk, Jerm, Silig, dairy, milk, Juicy} \rightarrow $6 \times 5 \times 4 = 120$

formula:

$$n_{Pr} = \frac{n!}{(n-r)!} = \frac{6!}{(6-3)!} = \frac{6 \times 5 \times 4 \times 3!}{3!} = 120$$

Dairy Germ Eclair

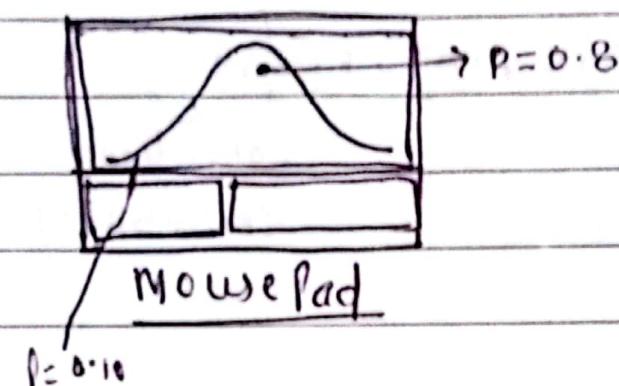
Combination: $n_{Cr} = \frac{n!}{r!(n-r)!} = \frac{6!}{3!(6-3)!}$

formula:

$$\frac{6 \times 5 \times 4 \times 3!}{3! \times 2 \times 1 \times 3!}$$

$$= \underline{\underline{20}}$$

P value



Every 100 times I touch
the mousepad 80 times I
touch this specific region

Hypothesis testing

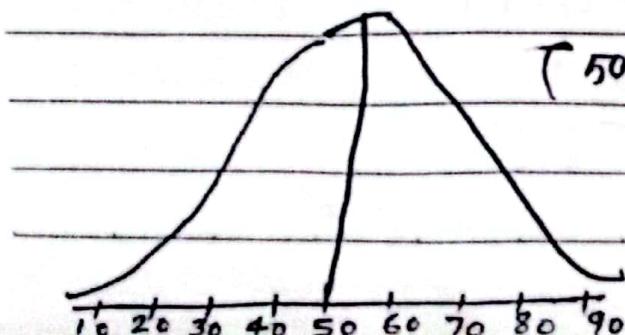
Coin → Test whether this coin is a fair coin or not
by perform 100 tosses

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

50 times head (The coin is fair)

Hypothesis Testing

- ① → Null Hypothesis Testing : Coin is fair
- ② → Alternative Hypothesis : - Coin is unfair
- ③ → Experiment
- ④ → Reject or Accept the Null Hypothesis .





Mo Tu We Th Fr Sa Su

Date:

linkedin: SANTOSH.CO.HR

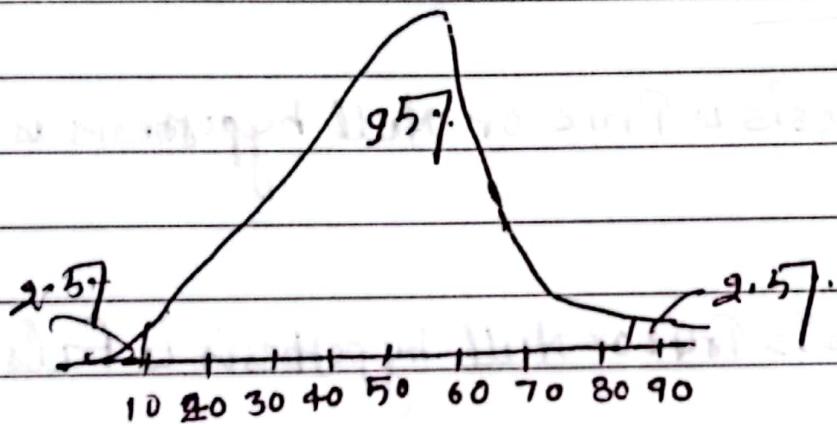
Significance Value: α (Alpha) α

$\alpha = 0.05$ (This is defined by Domain Expert)

Here, $1 - 0.05 = 0.95 = 95\%$.

Which means,

Confidence Interval (CI) = 95%.



Note! The experiment is falls between 95%, then it should be assumed fair (coin)

Note! Outside 2.5%, it is unfair (coin)



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date:

Type 1 and Type 2 Error

Null Hypothesis (H_0) = Coin is fair

Alternative Hypothesis (H_1) = Coin is not fair

Reality Check

→ Null hypothesis is True or Null hypothesis is false

Decision

Null hypothesis is True or Null hypothesis is False.

Outcome 1: We reject the Null Hypothesis, When in reality it is False \rightarrow yes

Output 2: We reject the Null Hypothesis, When in reality it is true \rightarrow Type 1 Error

Output 3: We accept the Null Hypothesis When in reality it is False \rightarrow Type 2 Error

Output 4: We accept the Null Hypothesis When in reality it is true \rightarrow Good



Mo Tu We Th Fr Sa Su

Date:

Confusion matrix

	P	N	
T	TP	TN	→ Typ 2
F	FP	FN	

↓
Typ1

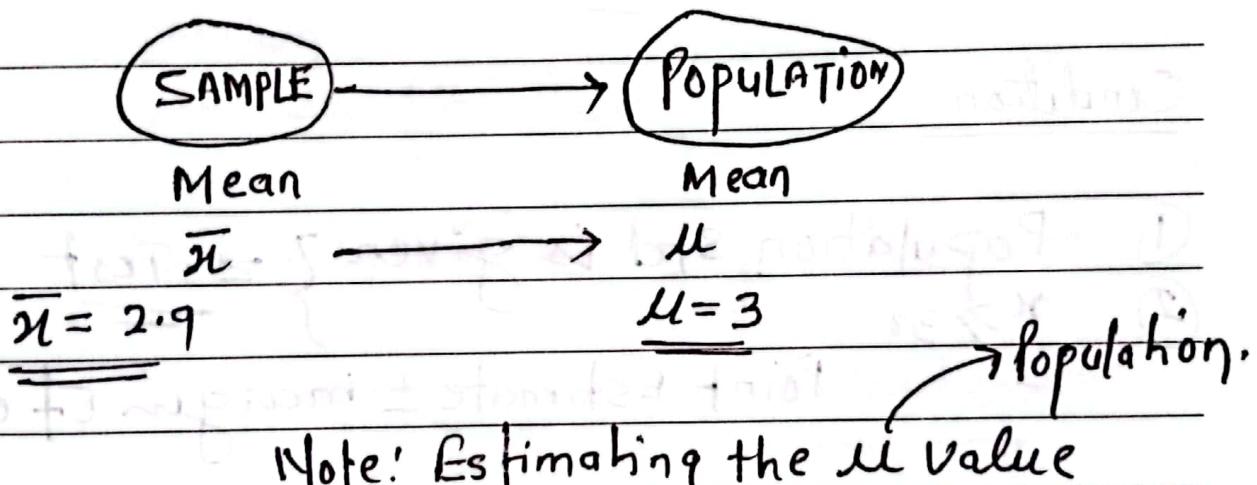
<input type="checkbox"/>						
--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------

Date :

Confidence Intervale

Point estimate: The value of any statistic that estimates the value of a parameter

In Inferential Stats.



Note! Estimating the μ value

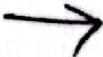
Confidence Intervale

Point estimate.

Point Estimate \pm margin of Error

e.g' On the quant test of CMAT Exam, the standard deviation is known to be 100 . A sample of 25 test takers has a mean of 520 score. Construct a 95% CI about the mean?

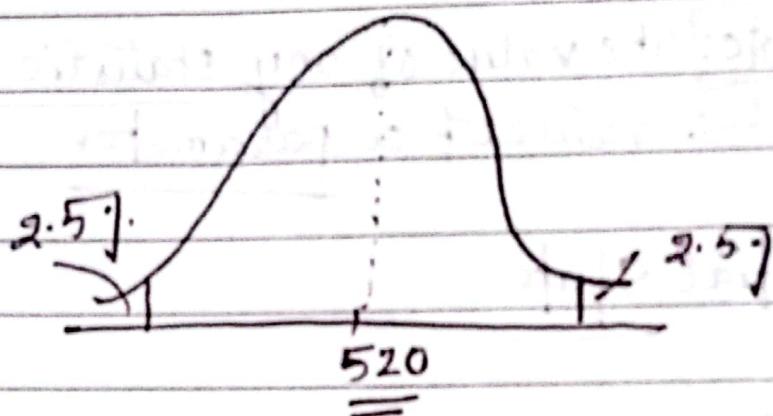
$$\text{Ans} \quad \sigma = 100, n = 25, \alpha = 0.05, \bar{x} = 520$$



80/19

here,

$$\alpha = 1 - 0.95 = \underline{\underline{0.05}}$$



Condition

- ① Population std is given } \Rightarrow z-test
 ② $n \geq 30$ } =

- Point Estimate \pm margin of error

$$\bar{x} \pm z_{\frac{\alpha}{2}} \left[\frac{\sigma}{\sqrt{n}} \right] \rightarrow \underline{\underline{\text{standard error}}}$$

$$\begin{aligned} \text{Upper bound} &= \bar{x} + z_{0.05} \frac{100}{\sqrt{25}} \\ &= 559.2 \text{ / } \end{aligned}$$

$$\begin{aligned} \text{Lower bound} &= \bar{x} - z_{0.05} \frac{100}{\sqrt{25}} \\ &= 480.8 \end{aligned}$$



eg(2)

On the quant test of CAT exam, a sample of 25 test takers has a mean of 520 with a standard deviation of 80. Construct 95% CI about the mean?

→

Condition:

$$n = 25, \bar{x} = 520, s/d = 80 \\ \alpha = 0.05$$

Here, Population sd is
not given. → (t-test)

Point Estimate \pm margin of error

$$\bar{x} \pm t_{\frac{\alpha}{2}} \left(\frac{s}{\sqrt{n}} \right) \rightarrow \text{standard error}$$

$$\text{Upper bound} = \bar{x} + t_{\frac{\alpha}{2}} \left(\frac{s}{\sqrt{n}} \right)$$

$$\text{Degree of freedom} = n - 1 = 25 - 1 = 24$$

$$= 520 + 2.064 \left(\frac{80}{\sqrt{5}} \right)$$

$$= 553.024, 11$$

$$\text{Lower bound} = \bar{x} - t_{0.05/2} \left(\frac{s}{\sqrt{n}} \right)$$

$$= 520 - 2.064 \left(\frac{80}{\sqrt{5}} \right)$$

$$= \underline{\underline{486.97}}$$

$$[486.97 \rightarrow 553.024]$$

Mo Tu We Th Fr Sa Su

① One Sample Z-Test

↳ Population sd is given

↳ Sample size $n > 30$

* In the Population, the average IQ is 100 with a sd of 15. wants to test a new medication to see if there is position or negative effect on intelligence, or no effect at all. A sample of 30 Participants who have taken the medication has a mean of 140. Did the medication affect the intelligence?

$$\alpha = 0.05, (C.I = 95\%)$$

Ans/

i) Define Null Hypothesis

$$H_0: \mu = 100$$

2) Alternative Hypothesis $H_1: \mu \neq 100$

3) State Alpha

$$\alpha = 0.05$$

4) State Decision Rule



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Date:

(2) tail test

95%

 $\alpha = 0.05$

Z table

2.57

 $\frac{+}{-}$

-1.96

85

2.57

1 - 0.025 = 0.975

 $\frac{-}{-}$ $\frac{-}{-}$ $P \leq 0.05$

5) Calculate Test statistics

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

$$\left(\frac{\sigma}{\sqrt{n}} \right)$$

standard error

$$= \frac{140 - 100}{\frac{15}{\sqrt{30}}} = \frac{40}{15} \times \sqrt{30} = 14.60 //$$

State Our Decision:

$$14.60 > 1.96$$

$$\{ Z = 0.2$$

$$\{ Z = 14.60 \}$$

If Z is less than -1.96 or greater than 1.96 , reject the Null Hypothesis //

Meditation improves the intelligence.

Or, decrease?

\Rightarrow Improves Intelligence //

IN

SANTOSH-LOHAR

Date:



Mo Tu We Th Fr Sa Su

One-Sample t-test

Z -Test \Rightarrow Population std.

t-test \Rightarrow Unknown population std.

Q) Population the average IQ = 100

$$n = 30, \bar{x} = 140, s = 20$$

Did the meditation affect intelligence?

$$\alpha = 0.05$$

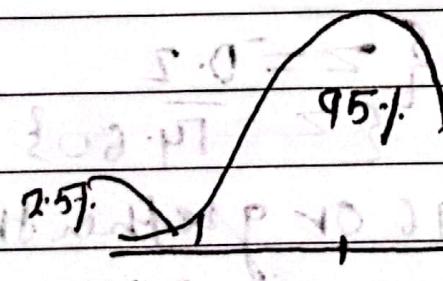
Ans) 1) $H_0 : \mu = 100$

2) $H_1 : \mu \neq 100$

3) Calculate the degree of freedom

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

4) State Decision Rule:



$$\alpha = 0.05, t = 10.96 > 2.045$$

Reject Null Hypothesis

$P \leq$ significance value

⑤ T-Test

$$\bar{x} = 140$$

{ Increase the }
Intelligence

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

$$\mu = 100$$

$$s = 20$$

$$n = 30$$

$$= 10.96$$

Reject the Null Hypothesis //

CHI SQUARE TEST

1) Chi-square test ~~clears~~ about population proportional claims proportions

It is a non parametric test that is performed on categorical (normal or ordinal) data.

Q.) In the 2000 Indian census, the age of the individuals in a small town were found to be the following

Less than 18	18-35	> 35
20%	30%	50%

In 2010, age of $n = 500$ individuals were sampled. Below are the results.

48	< 18	18-35	> 35
121	288	91	

Now, using $\alpha = 0.05$, would you conclude the population distribution of ages has changed in the last 10yrs.

80%	< 18	18-35	> 35	
20%	30%	50%		→ Expected.

< 18	18-35	> 35	$n = 500$
121	288	91	→ Observed
$100 = 500 \times 0.20$	$500 \times 30\%$	$500 \times 50\%$	→ Expected

$$(121 - 100)^2 / 100 = 150$$

$$(288 - 250)^2 / 250 = 250$$

IN

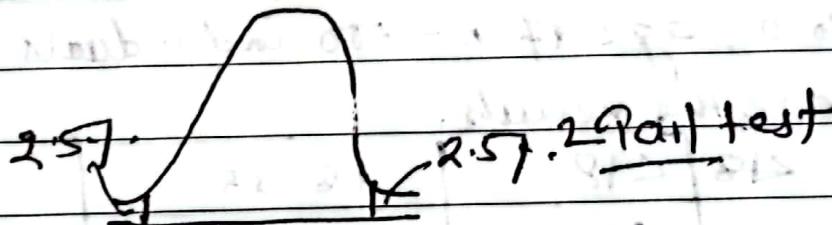
SANTOSH LOHAR

distribution = 3

Mo	Tu	We	Th	Fr	Sa	Su	Date:
< 18	18-35						735
121	288				91		Observation
100	150				250		- Expected

Follow: Chi square Table

- Now
- ① H_0 = The data meets the distribution 2011 Census
 - H_1 = The data does not meet
 - ② $\alpha = 0.05$ (95% C.I.)
 - ③ Degree of freedom = $(n-1) = (3-1) = 2$
 - ④ Decision Boundary:



If χ^2 is greater than 5.991 reject H_0

See Chi square table.

With the help of α & df

- ⑤ Calculate test statistics -

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$\begin{aligned}
 &= (121-100)^2 + (288-150)^2 + (91-250)^2 \\
 &\quad 100 \qquad 150 \qquad 250 \\
 &= 232.494
 \end{aligned}$$

IN

SANTOSH - LOHAR



<input type="checkbox"/>						
--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------

Date :

Here,

$$\chi^2 = 232.494 > 5.99$$

\Rightarrow Reject the Null Hypothesis

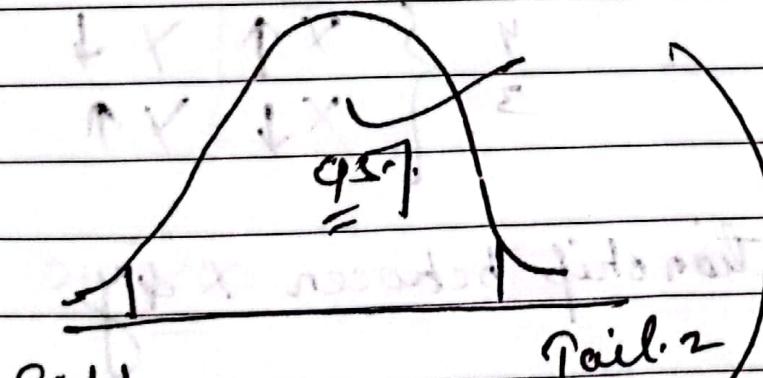
$$0.11 > 0.05$$

\hookrightarrow Accept the Null.

P-value $<$ Significance

\hookrightarrow Reject the Null Hypothesis
OR

Accept the Null Hypothesis



Tail-1

$$P = 0.002 < 0.05$$

$$0.11 > 0.05$$

\checkmark
 \downarrow
Accept

\downarrow
Reject the Null Hypothesis

Covariance

X	Y	
weight	height	{ we can say that
50	160	$X \uparrow \rightarrow Y \uparrow$
60	170	$X \downarrow \rightarrow Y \uparrow$
70	180	
75	181	

No. of hours study play.

2	6
3	4 { $X \uparrow Y \downarrow$
4	3 { $X \downarrow Y \uparrow$

Now, quantify relationship between X & Y.

Covariance

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

$n-1$ → For Sample distribution.

N → For Population distribution.

 y 

-ve covariance.

+ve correlation

 $x \downarrow \rightarrow y \uparrow$ $x \uparrow \rightarrow y \downarrow$ Negative -ve Correlation,
-ve Covariance

$\begin{matrix} x \\ \downarrow \\ x \end{matrix} \rightarrow$ -ve correlation
 $\begin{matrix} \downarrow \\ x \end{matrix}$ $\begin{matrix} x \\ \uparrow \\ x \end{matrix}$ -ve covariance

Covariance
 $\begin{matrix} x \\ \uparrow \\ x \end{matrix} = 0$

Disadvantage of covariance

Pearson Correlation Coefficient

(-1 to +1)

⇒ The more towards +1 more positively Correlation

The more towards -1 more negative correlated

$$f(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} \Rightarrow (-1, +1)$$



$$\text{Spear}(x, y) = \frac{\text{Cov}(R(x), R(y))}{R_{\sigma_x} \times R_{\sigma_x}}$$

Height	weight	$(R(x))$	$(R(y))$
170	75	2	2
160	62	3	3
150	60	4	4
145	55	5	5
180	85	1	1

\rightarrow this will be
cancelled

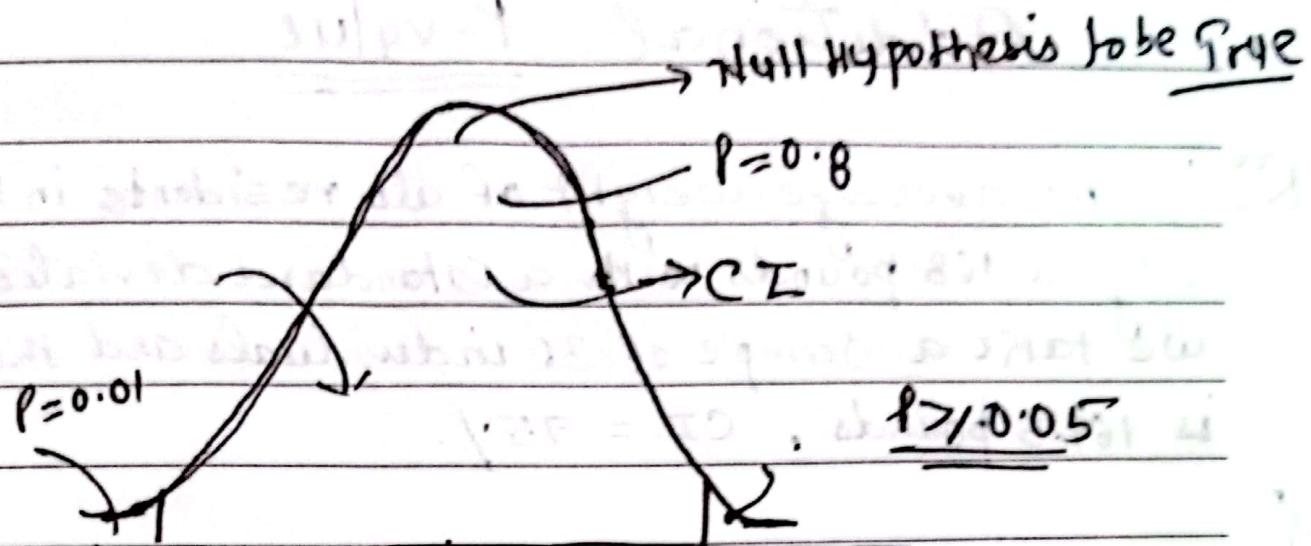
we are interested
for Computing.

Non-linear Proportions

$\beta \leq 0.05 \rightarrow$ Reject the Null Hypothesis
 $\rightarrow \alpha = 0.05$

$\rightarrow 5\%$, Probability the Null Hypothesis is
Correct

$\beta > 0.05 \rightarrow$ Accept the Null Hypothesis



Here, Null Hypothesis should be Accepted because of

$$\underline{P > 0.05}$$

P-value \leq C.I.

\hookrightarrow Reject the Null hypothesis

P-value \geq C.I.

\hookrightarrow Accepted Null H.

30.2.2021 Additional P-value

Q>:- The average weight of all residents in Bangalore city is 168 pounds with a standard deviation 3.9 we take a sample of 36 individuals and the mean is 169.5 pounds, CI. = 95%.

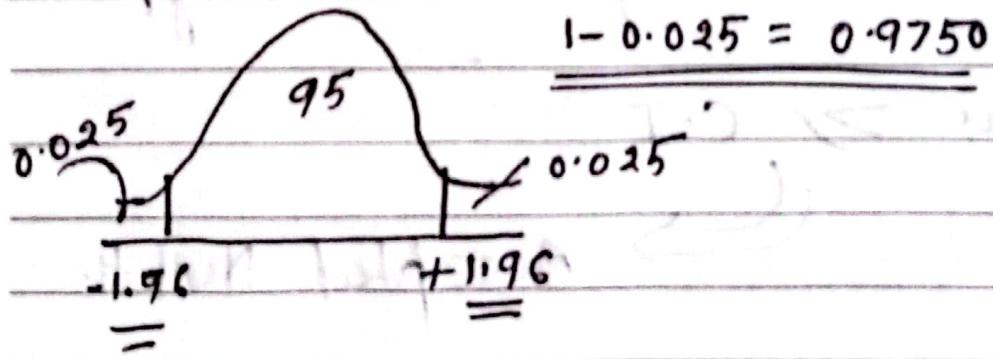
Solⁿ:

$$\Rightarrow \mu = 168, \sigma = 3.9, \bar{x} = 169.5, n = 36, \alpha = 0.05$$

① $H_0 = \mu = 168$
 $H_1 = \mu \neq 168$

② $\alpha = 0.05$

③ Decision Boundary:



④ Z-test

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = \frac{1.5 \times 6}{3.9} = 2.307$$

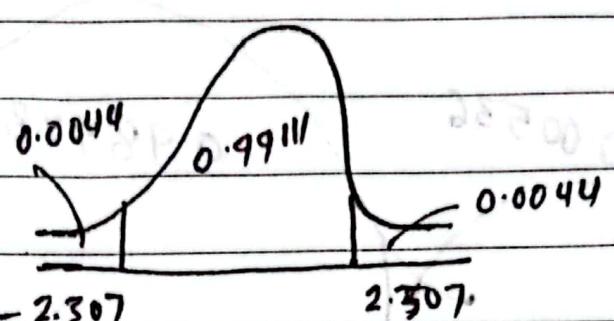
$\therefore Z = 2.307 > 1.96$

⇒ Reject Null Hypothesis



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						

Date:

P-value

$$1 - 0.9911$$

$$\begin{aligned} \text{Pvalue} &= 0.0044 + 0.0044 \\ &= 0.0088 \end{aligned}$$

Here, $\text{Pvalue} < 0.05$

$0.0088 < 0.05 \rightarrow \text{Reject the Null Hypothesis}$

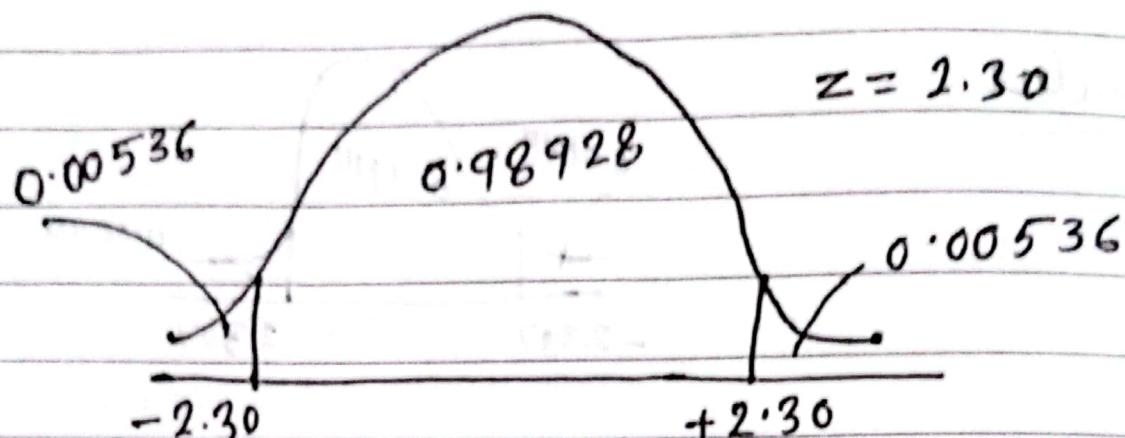
* $\text{Pvalue} \leq \text{Significance value}$

↓
Reject the Null Hypothesis

$\text{P-value} > \text{Significance value}$ → fail to Reject the Null Hypothesis



Mo	Tu	We	Th	Fr	Sa	Su
<input type="checkbox"/>						



$$1 - 0.98928$$

Here, $2.30 > 1.96$ {Reject the Null Hypothesis}

$$\text{Pvalue} = 0.00536 + 0.00536$$

If less than $\alpha \Rightarrow$ Reject the hypothesis

e.g.:

Average age of a college is 24 years with a standard deviation 1.5. Sample of 36 students, students mean is 25 years with $\alpha = 0.05$, CI = 95%. do the age vary?

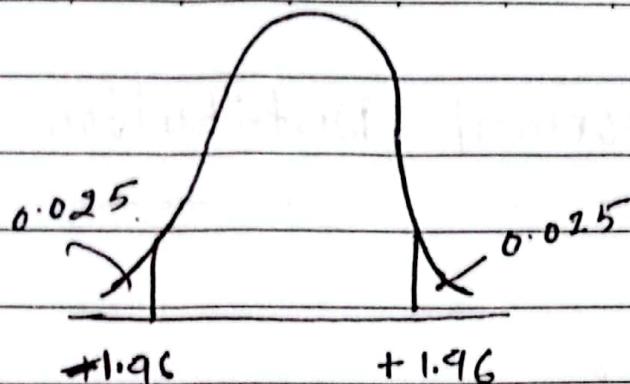
⇒ ①

$$H_0: \mu = 24 \quad \sigma = 1.5, n = 36, \bar{x} = 25, \alpha = 0.05$$

$$H_1: \mu \neq 24$$

② $\alpha = 0.05$

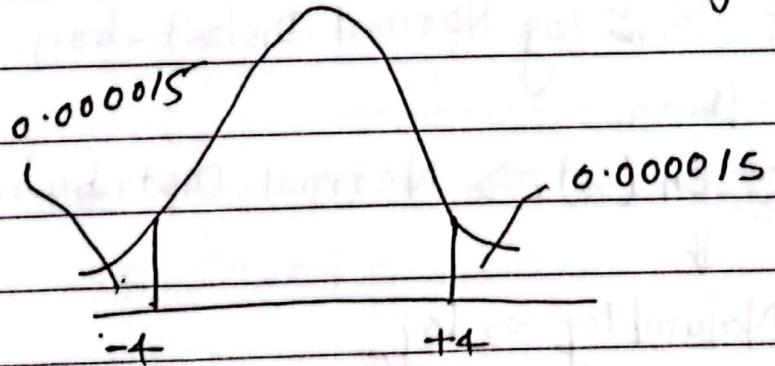
③



④

$$\begin{aligned} Z\text{-score} &= \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{25 - 24}{1.5} \times 1 \\ &= \frac{1 \times 1}{1.5} \\ &= 4 \end{aligned}$$

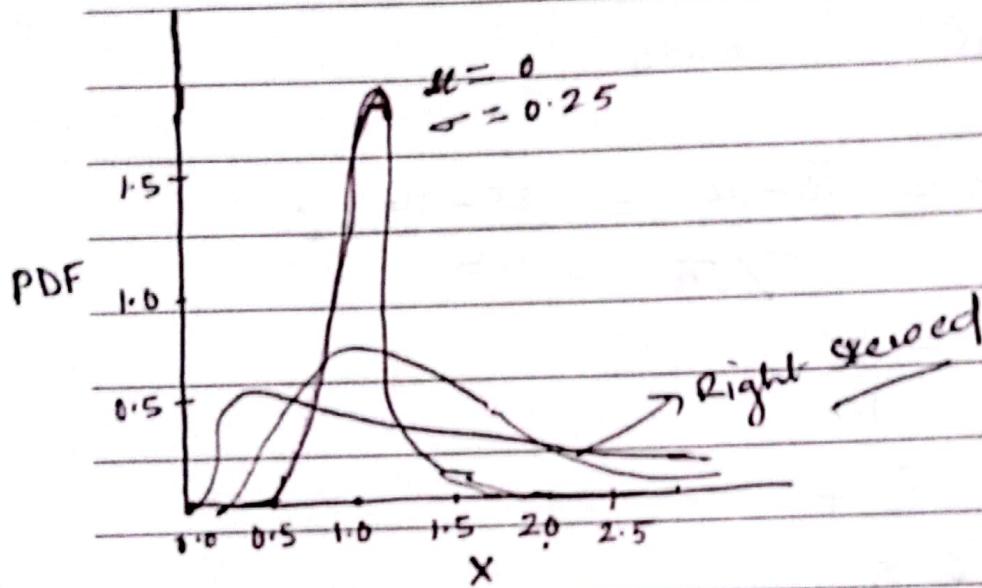
Now, $4 > 1.96 \rightarrow$ Reject Null Hypothesis.



$$\begin{aligned} \text{Pvalue} &= 0.000015 + 0.000015 \\ &= 0.00003 \end{aligned}$$

Now, $\text{Pvalue} \leq \text{Significance} \rightarrow$ Reject the Null Hypothesis

Log Normal Distribution



if $X \approx$ log Normal Distribution.

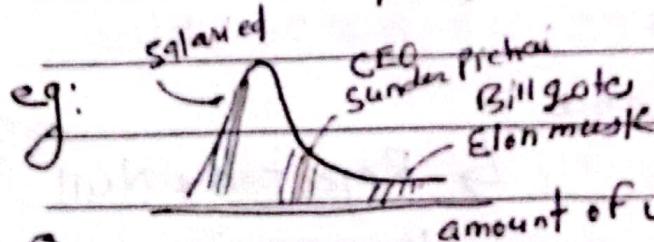
then,

$$y \approx \ln(x) \Rightarrow \text{Normal Distribution.}$$



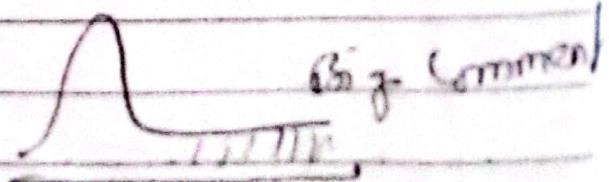
Natural log $\Rightarrow \log_e$

$X \approx \exp(y)$ (we can also write like this)



① Wealth distribution

② Comment in YT/FB





Mo Tu We Th Fr Sa Su

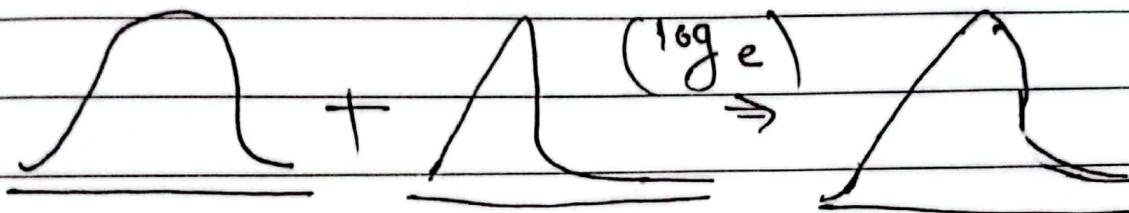
Date:

Where we used (Log Normal Distribution)

↳ Machine learning



Simple linear Regression



↓
Normal
Distribution

↓
log-Normal
Distribution.

↓
Normal Distribution

Now, we can use for the purpose of Training data model in machine learning and from this our algorithms will be efficiently Trained which we called

↓
Data Transformation Technique