

Statistics

What is statistics?

Statistics is a set of mathematical methods and tools that enable us to answer important questions about data.

It is divided into two categories: Descriptive and Inferential

Descriptive Statistics: (on population data)

This offers methods to summarize data by transforming raw observations into meaningful information that is easy to interpret and share.

Inferential Statistics: (on sample data → mostly used)

This offers methods to study experiments done on small samples of data chalk out the inferences to the entire population (entire domain).

Statistics

Descriptive

Inferential

- Measures of Central Tendency
- Measures of Variability
- Measures of Shape - Skewness
- Percentiles and Quartiles
- Frequency Distribution
- Covariance and Correlation

- * Central Limit Theorem
- Hypothesis Testing
 - Z-Test
 - T-Test
 - Chi Square Test

1. Measures of Central Tendency

mean
median
mode

1. Mean : (Average)

$$\text{Mean} = \frac{\text{Sum of all data}}{\text{Total no. of data}}$$

```
import numpy as np
```

```
import pandas as pd
```

```
ar = np.array([4, 5, 6, 2, 1, 8, 5, 6, 4, 7])
```

```
np.sum(ar)/len(ar) # manual process
```

```
np.mean(ar) # automatic process
```

" For solving real world problems. "

```
dataset = pd.read_csv("Students.csv") # read csv file
```

```
dataset.head(3) # shows first three rows
```

```
dataset["Age"].mean() # find mean using pandas
```

```
np.mean(dataset["Age"]) # find mean using numpy
```

" For graphical output "

```
import matplotlib.pyplot as plt # to create a graphs
import seaborn as sns # for advancements in graphs
```

```
mn = np.mean(dataset[["Age"]]) # store mean in mn
```

```
sns.histplot(x="Age", data=dataset, bins=[i for i in
range(10, 100, 10)])
```

bins is used to manage bins of graph columns.

This is column vs data graph.

```
plt.plot([mn for i in range(10, 250)], [i for i in
range(10, 250)], c="red")
```

Display mean in graph.

```
plt.show() # display graph with all figures
```

2. Median : (Middle Value)

5 Arrange data in ascending order.

6 Odd no. of data → middle value

Even no. of data → average of both middle values.

"For finding median"

```
dataset.isnull().sum() # displays total number of nulls in
# each columns in dataset.
```

```
dataset["Age"].fillna(dataset["Age"].mean(),
inplace=True)
```

filling up the null values by mean of dataset column.
 # (data cleaning uses measurement of central tendency)

Now, there is no null value in column. So, can find median
`np.median(dataset["Age"])` # find median using numpy.

`dataset["Age"].median()` # find median using pandas.

"For Graphical output."

`md = np.median(dataset["Age"])` # store median in md

`sns.histplot(x="Age", data=dataset, bins=[i for i in range(50, 100, 10)])` # Graph format
`plt.plot([md for i in range(10, 250)], [i for i in range(10, 250)], c="blue")`
 # Plot mean in graph.

`plt.show()` # Display graph with all figures

3 Mode : (Most Repeated Value)

- 6 Mode is a data with highest frequency (repetition).
- 6 Commonly used in categorical data. e.g. gender, class, marital status, sex, etc.

"For finding mode."

```
a
dataset = pd.read_csv("Students.csv")
dataset["Age"].mode()[0] # calculate mode using pandas.
# index 0 is used because it returns a series and the
# mode is the first data of series.
```

```
dataset["Age"].value_counts() # Shows frequencies of
# datas in column.
```

"To get Graphical output."

```
mo = dataset["Age"].mode()[0] # Store mode in mo
sns.histplot(x="Age", data=dataset, bins=[i for i in
range(30, 100, 10)]) # Graph format
plt.plot([mo for i in range(30, 250)], [i for i in range
(30, 250)], c="green") # Plot mode in graph
plt.show() # Display graph with all figures
```

"To display mean, median and mode in same graph"

```
sns.histplot(x="column", data=dataset, bins=[i for i in
range(30, 100, 10)]) # Graph format
plt.plot([mn for i in range(30, 250)], [i for i in range
(30, 250)], c="red", label="mean") # Plot for mean in range
plt.plot([md for i in range(30, 250)], [i for i in range
(30, 250)], c="blue", label="median") # Plot for median in range
plt.plot([mo for i in range(30, 250)], [i for i in range
(30, 250)], c="green", label="mode") # Plot for mode in range
```

```
# Plot median in graph  
plt.plot([m0 for i in range(0, 250)], [i for i in  
range(0, 250)], c="green", label="mode")  
  
# Plot mode in graph  
plt.legend() # Shows legend according to given labels  
plt.show() # Show graph with all figures
```

2. Measures of Variability

1. Range

- Range is the difference between the maximum and minimum values in the dataset.
i.e. Range = maximum - minimum.
- It provides a simple measure of the spread of the data, but it can be sensitive to outliers.

"For finding range of given dataset."

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt # not used  
import seaborn as sns # not used
```

```
dataset = pd.read_csv("Students.csv")
```

```
dataset.head(3) # display first three rows.
```

`min_r = dataset[["Age"]].min() # Returns min value
max_r = dataset[["Age"]].max() # Returns max value.`

`min_r, max_r # Prints min and max values
Example: (0.42, 80.0)`

`range = max_r - min_r # Gives range of column data.
range # Displays range`

2. Mean Absolute Deviation

↳ The mean absolute deviation of a dataset is the average distance between each datapoint and the mean.

↳ It gives us an idea about the variability in dataset.

↳ Mathematically, $MAD = \frac{\sum |x_i - \bar{x}|}{n}$

"Making choice based on data of sections"

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

`Sec-a = np.array([75, 65, 73, 68, 72, 67])`

`Sec-b = np.array([90, 47, 43, 96, 93, 51])`

```
# Created two arrays with some data using numpy.  
no = np.array([1, 2, 3, 4, 5, 6])  
# Array for y-axis values.
```

```
plt.figure(figsize = (10, 3)) # Fix figure size  
plt.scatter(sec_a, no, label = "sec A")
```

```
# Draw scatter graph of sec-a vs. no  
plt.scatter(sec_b, no, c = "blue", label = "sec B")
```

```
# Draw scatter graph of sec-b  
plt.plot([70, 70, 70, 70, 70, 70], no, c = "red",  
label = "mean")
```

```
# plot a line of mean from 0 to 6.  
plt.legend() # display legend according to labels  
plt.show() # show graph with all figures
```

" From the diagram, the less scattered data is less deviated and should be chosen.

But sometimes drawing graphs can be a tedious task or even not possible. So, we need to have some mathematical calculations to meet our needs.

The calculations can be done as follows:

Find mean of both sections

Mean absolute division is used when we have got the

same mean while figuring out for best.

i.e sec-a and sec-b has same mean

mean = np.mean(sec_b) # Same for sec-a

mean # display mean : 70.0

find mean absolute deviation for both.

$$\text{mad_a} = \text{np.sum}(\text{np.abs}(\text{sec_a} - \text{mean})) / \text{len}(\text{sec_a})$$

$$\text{mad_b} = \text{np.sum}(\text{np.abs}(\text{sec_b} - \text{mean})) / \text{len}(\text{sec_b})$$

`mad_a, mad_b` # display mean absolute deviation:
(3.3333333333333335, 23.0)

" We found that mean absolute deviation of sec-a is less than mean standard deviation of sec-b.
So, sec-a should be our choice."

" If we got same mean absolute deviation for both data then standard deviation is used."

3. Standard Deviation

- ↳ The standard deviation is the measure of the amount of variation or dispersion of a set of values.
- ↳ A low standard deviation indicates the values tends to be close to mean (also called the expected value) of the set, while a high standard deviation indicates that the values are spread out over a wider range.

↳ Mathematically,

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

where, σ = population standard deviation.

N = the size of the population

x_i = each value from the population

$\mu \in$ the population mean.

" Calculating standard deviation on previous lists."

np.std(sec-a), np.std(sec-b) # calculate and show;
 # (3.862..., 23.180...)

4. Variance

- 6 Variance is a measure of how data points differ from a mean.
- 6 According to layman, a variance is a measure of how far a set of data (numbers) are spread out from their mean (average) value.

6 Mathematically,

$$\text{Var} = \sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N}$$

" Calculating variance on previous lists."

np.var(sec-a), np.var(sec-b) # calculate variance and
 # display: (14.91..., 537.33...)

" Calculating variance and standard deviation on real

World data. If the variance and standard deviation is high it is a scattered data but as a data scientist we always prefer less scattered data for our work."

Example:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns # import necessary libraries.
```

```
dataset = pd.read_csv("Students.csv")
```

```
dataset.head(3)
```

```
dataset["Age"].var() # Calculate variance.
```

```
dataset["Age"].std() # Calculate standard deviation.
```

```
# Plot a column 'vs' data graph for visualization
sns.histplot(x="Age", data=dataset, bins=np.arange(0, 100, 10))
plt.show()
```

" If we don't want them to calculate all we have another method which will show {count, mean, std, min, 25%, 50%, 75%, max} of all the numerical columns "

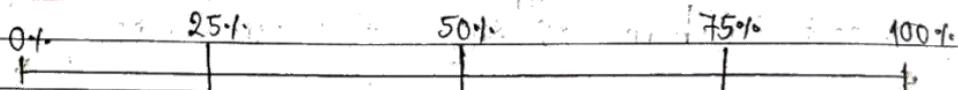
```
dataset.describe()
```

3. Percentiles and Quartiles

Percentiles

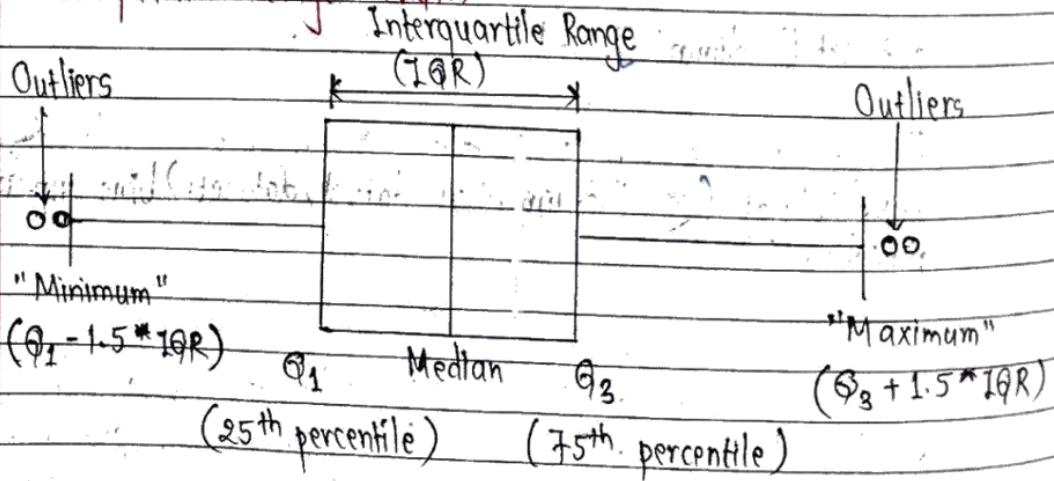
- 6 Percentiles are used in statistics to give you a number that describes the value that a given percent of the values are lower than. It ranges from 0% to 100%.

Quartiles



- 6 The equal four divisions of percentiles made are called quartiles.

Interquartile Range (IQR)



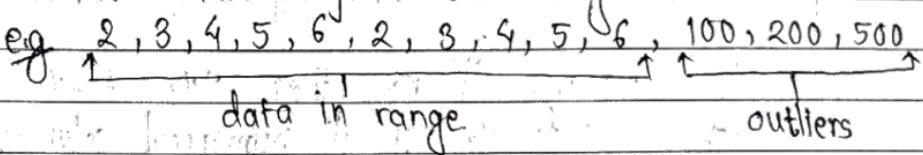
From diagram;

$$IQR = Q_3 - Q_1$$

NOTE :

- Percentage vs Percentile : Percentage works on fixed defined value whereas percentile works on undefined base value. Also, percentage shows obtained value whereas percentile shows obtained ranking.

- Outliers : The data beyond our range.

e.g. 

From diagram; the outliers are values beyond Minimum and Maximum which can be calculated using IQR.

" Using Percentiles and Quartiles to solve real world problems."

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
dataset = pd.read_csv("Students.csv")
```

```
dataset.head(3) # Display first three rows
```

```
dataset.isnull().sum() # Display total number of null
# values present in each column in dataset
```

```
dataset[["Agein"]].fillna(dataset[["Agein"]].mean(), inplace
= True)
```

filling up the null values by mean of dataset # column.

Calculating percentiles.

```
np.percentile(dataset["Age"], 50)
```

```
np.percentile(dataset["Age"], 25)
```

```
np.percentile(dataset["Age"], 75)
```

Displaying all percentiles 0%, 25%, 50%, 75%.

and 100% of all numerical columns at once.

```
dataset.describe()
```

"If you find a large gap between min and 25% or max and 75% there exist outliers. No. of outliers is directly proportional to the gap."

We can see outliers by using a box plot.

```
sns.boxplot(x = "Age", data = dataset)
```

"Percentiles are used to identify and remove outliers which working with machine learning algorithms."

4. Measures of Shape - Skewness

Skewness

G Skewness measures the asymmetry of the distribution.

- 6 A skewness value of 0 indicates a perfectly symmetrical distribution.
- 6 Positive skewness indicates that the distribution is skewed to the right (i.e. the tail is longer on the right), while negative skewness indicates a left skew (i.e. the tail is longer on the left).
- 6 Mathematically,

$$\text{Skewness} = \frac{\sum (x_i - \bar{x})^3}{(N-1) \cdot \sigma^3}$$

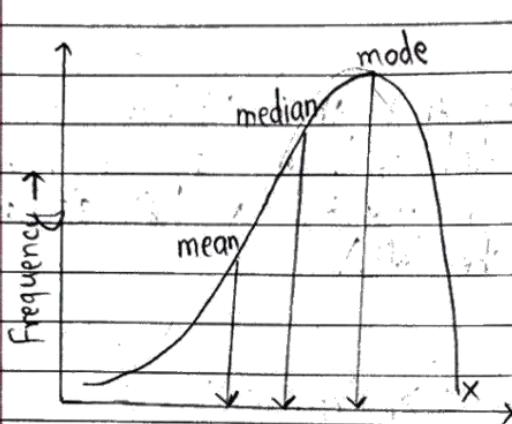
Frequency and Cumulative Distribution

- 6 A frequency distribution is a table that shows the number of occurrences of different values in a dataset.
- 6 A cumulative distribution shows the accumulation of frequencies up to a certain point. It is obtained by adding up the frequencies as you move through the values.

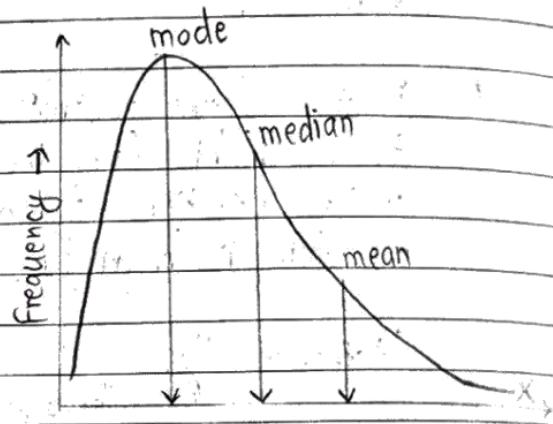
| Value Range | Frequency | Cumulative |
|-------------|-----------|--------------|
| 10 - 20 | 5 | 5 |
| 20 - 30 | 8 | 5 + 8 = 13 |
| 30 - 40 | 12 | 13 + 12 = 25 |
| 40 - 50 | 6 | 25 + 6 = 31 |

- 6 Histogram is used for visualization.

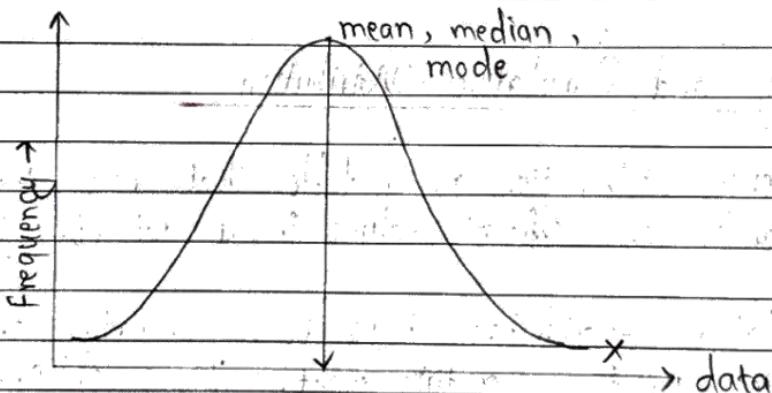
Skewness Charts.



(a) Negatively skewed



(b) Positively skewed



(c) Normal (no skew)

From diagrams :

- (a) When negatively skewed \rightarrow mean < median < mode
- (b) When positively skewed \rightarrow mode < median < mean
- (c) When no skew \rightarrow mean = mode = median

Working Practically in Jupyter :

import pandas as pd

```
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns
```

```
dataset = pd.read_csv("Students.csv")
```

```
dataset.head(3)
```

```
dataset["Age"].skew() # Returns skewness of column  
''' If skewness > 0 → positively skewed.  
If skewness < 0 → negatively skewed.  
If skewness = 0 → no skew. '''
```

```
# Plot a histogram for visualization.
```

```
sns.histplot(x="Age", data=dataset)  
plt.show()
```

```
''' Visualization by creating your own data. '''
```

```
data = np.random.normal(0, 100, 100) # Create a random  
# array of data.
```

```
data # Display created data.
```

```
df = pd.DataFrame({ "x": data }) # Create a column x where  
# all data is placed in different rows.
```

```
df["x"].skew() # Gives skewness
```

```
# Create a histogram for visualization.
```

```
sns.histplot(x="x", data=df)
```

plt.show()

" The normal (non-skew) data can be created as follows:

```
data = [2, 3, 3, 4, 4, 4, 5, 5, 5, 5, 6, 6, 6, 6, 6, 7, 7, 7, 7, 7, 8, 8, 8, 8, 8, 9, 9, 9, 9, 10, 10, 10, 10, 11, 11, 12] # non-skew data.
```

```
df = pd.DataFrame({ "x": data }) # Create a column  
# named x under which all data is placed.
```

```
df["x"].skew() # 0 skew
```

Histogram will be symmetrical.

```
sns.histplot(x="x", data=df, bins=[2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13])
```

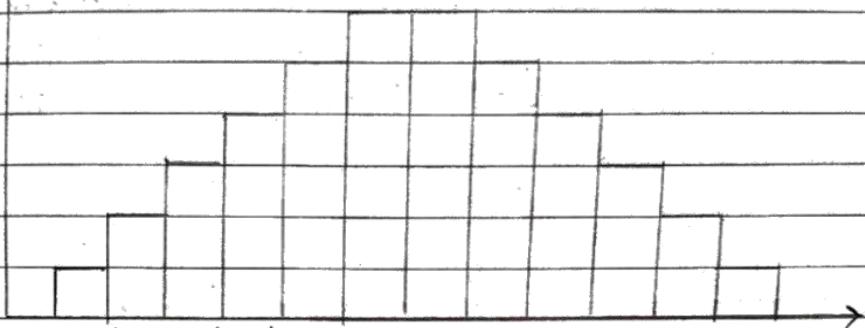
plt.show()

Mean, Median, Mode will be same.

```
df["x"].mean(); df["x"].median(); df["x"].mode()  
() [0] # (7.0, 7.0, 7)
```

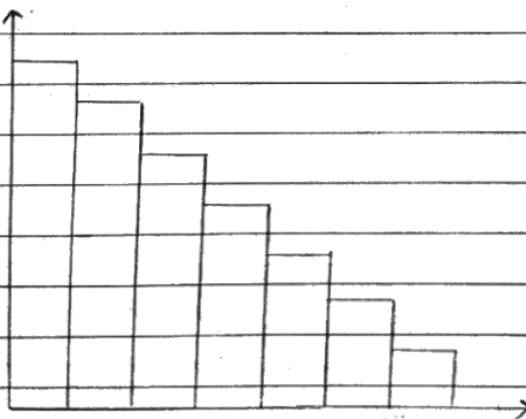
" Mean, Median and Modes will not be same in negatively and positively skewed charts (data).
For negatively skewed: mean < median < mode
For positively skewed: mode < median < mean "

Symmetric (normal) vs skewed and normal distributions:

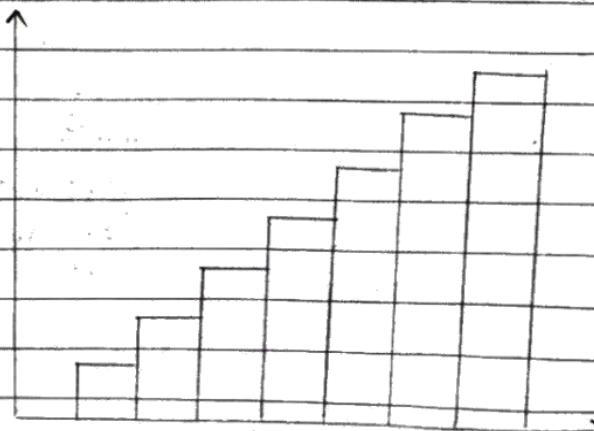


Normal distribution

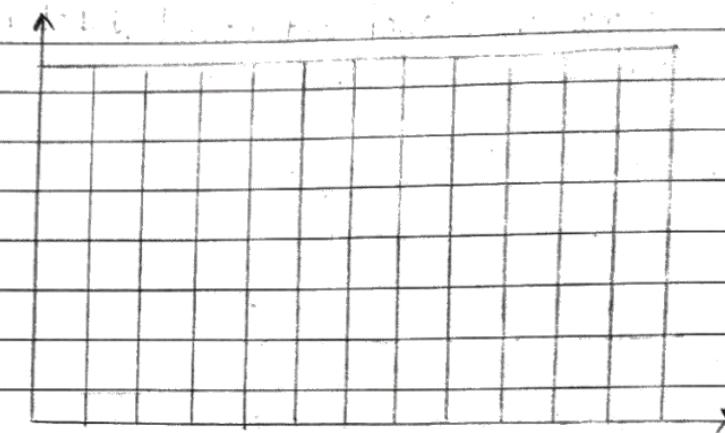
(unimodal, symmetric, the "bell curve")



Right-skewed distribution
(Positively skewed)

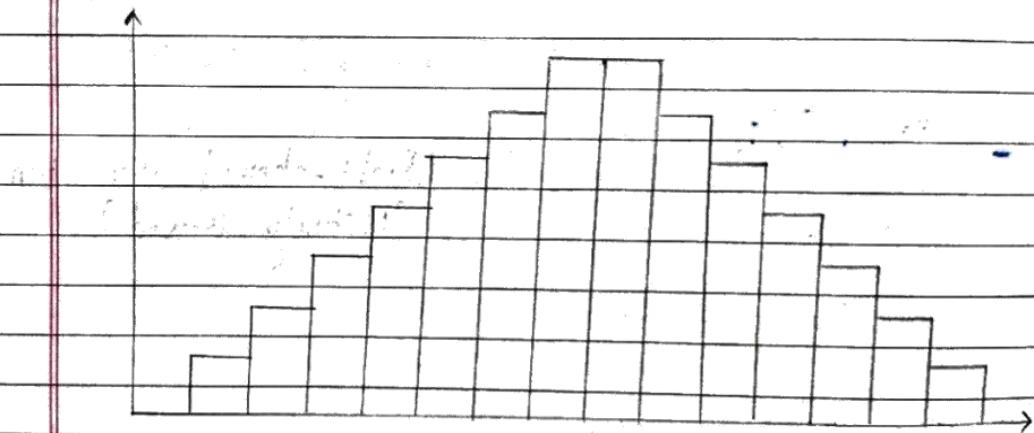


Left-skewed distribution
(Negatively skewed)

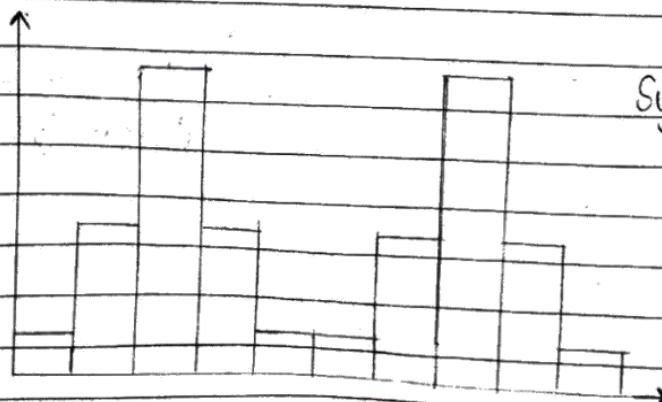


Uniform distribution
(equal spread,
no peaks)

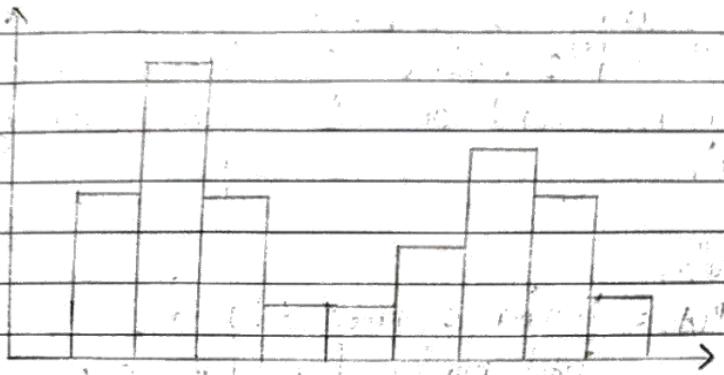
Unimodal vs Bimodal distributions.



Normal distribution
(unimodal, symmetric, the "bell curve")



Symmetrical bimodal
distribution
(two modes)



Non-symmetric bimodal distribution
(two modes)

5. Probability

Random Variables : A Random variable x is a function that assigns a real number to each outcome in the sample space of a random experiment.

1. Discrete Random Variable : A random variable that takes on a countable number of distinct values. e.g. dice roll

2. Continuous Random Variable : A random variable that can take on any value within a given range or interval. e.g. height of people in a sample.
This donot have any fix no. of possible outcomes.

Probability

6 Probability measures the likelihood of a particular outcome

or event occurring. It is typically expressed as a number between 0 and 1, where 0 indicates impossibility (event will not occur) and 1 indicates certainty (event will occur).

G Mathematically,

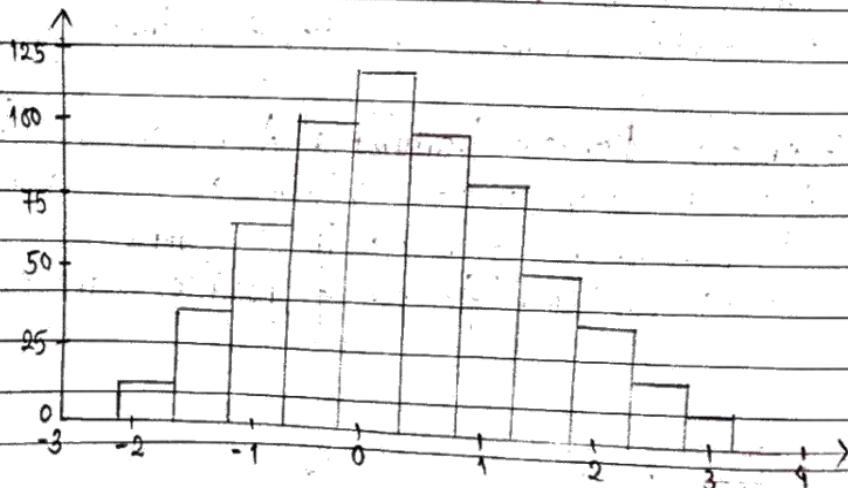
$$P(A) = \frac{\text{Number of times } A \text{ occurs}}{\text{Total number of possible outcomes}}$$

Probability Distribution

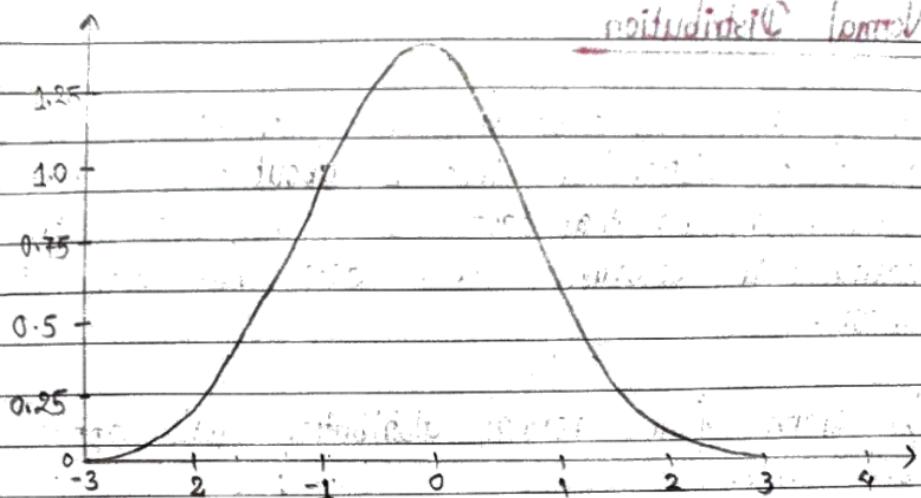
G Probability distributions describe how the probabilities of different outcomes are distributed over the sample space of a random variable.

G Two types of probability distribution are:

- (a) continuous probability distributions
- (b) discrete probability distributions



: (a) Discrete probability distribution



(b) Continuous probability distribution.

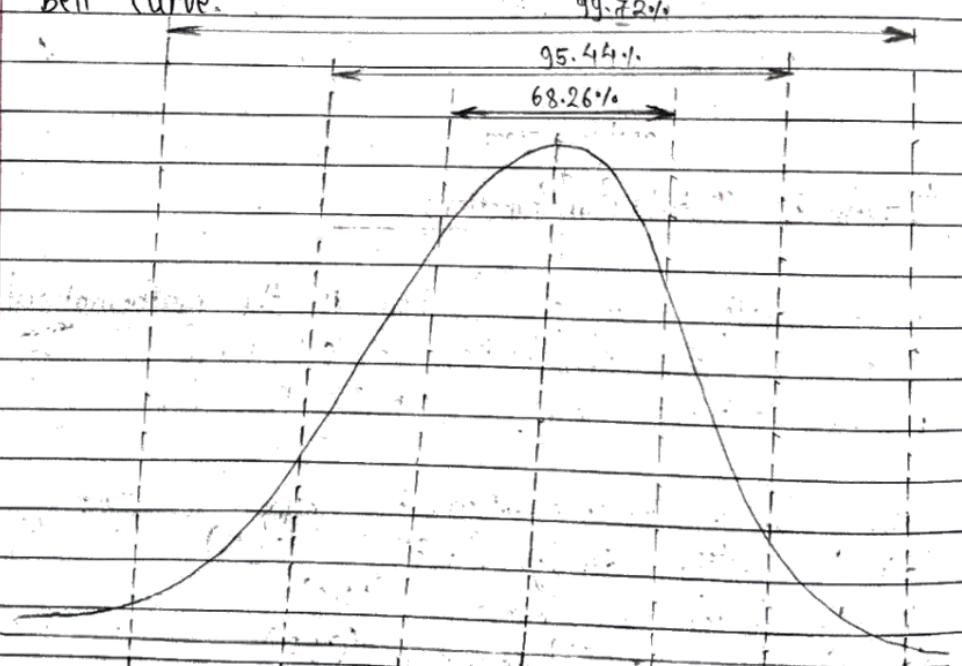
Probability Distribution Functions

- 6 A probability distribution function is the mathematical function that gives the probabilities of occurrence of different possible outcomes for an experiment.
- 6 Types of probability distribution functions are:
- (a) Probability density function (PDF)
 - (b) Probability mass function (PMF)
 - (c) Cumulative density function (CDF)
- (a) **Probability density function (PDF)** : This function is defined on continuous data. Density means total number of data in a particular range.
- (b) **Probability mass function (PMF)** : This function is defined on discrete data. Mass means frequency of data.
- (c) **Cumulative density function (CDF)** : This shows the percentage of data before and after a certain point.

6. Normal Distribution

It is known as the Gaussian distribution, is a probability distribution that is symmetric about the mean, showing that data near the mean are more frequent in occurrence than data far from the mean.

In graph form, normal distribution will appear as a bell curve.

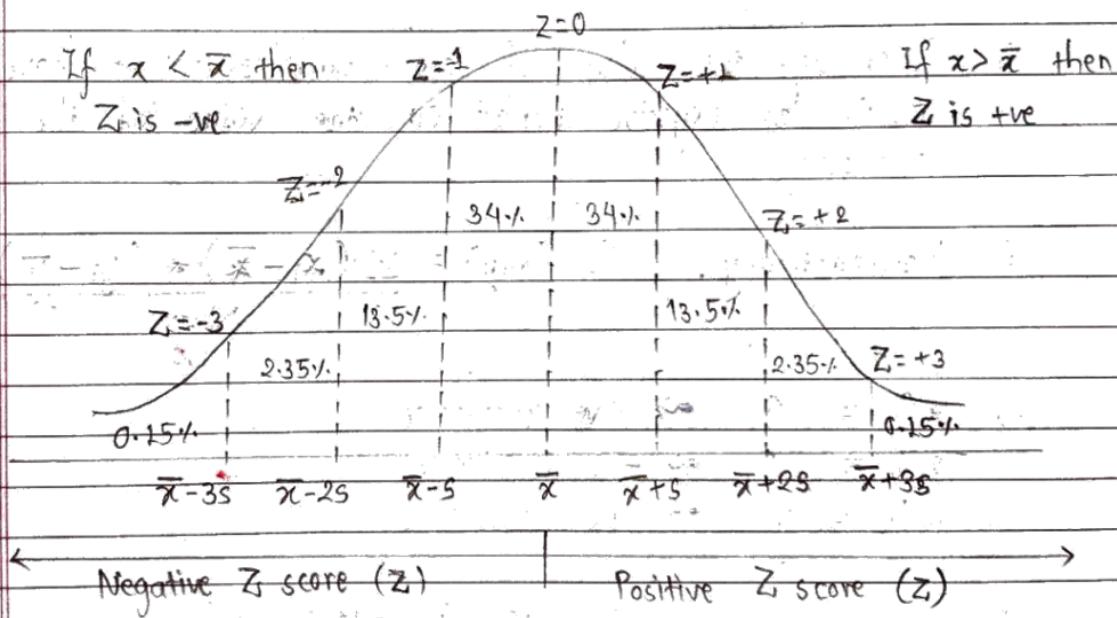


$$\begin{matrix} 2.14\% & 13.59\% & 34.13\% & 34.13\% & 13.59\% & 2.14\% \\ \mu-3\sigma & \mu-2\sigma & \mu-\sigma & \mu & \mu+\sigma & \mu+2\sigma & \mu+3\sigma \end{matrix}$$

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

Standard Normal Distribution

- ↳ The standard normal distribution, also known as the Z distribution or Z-score, is a special case of the normal distribution.
- ↳ mean (μ) of 0 and a standard deviation (σ) of 1.



In diagram : \bar{x} = mean score
 x = individual score.

7. Covariance and Correlation.

Covariance

- ↳ Covariance signifies the direction of the linear relationship between the two variables. By direction we mean if the variables are directly proportional or inversely proportional to each other.
- ↳ Increasing the value of one variable might have a positive or negative impact on the value of the other variable.
- ↳ Mathematically, $\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x}) * (y_i - \bar{y})}{N}$
- ↳ $x \uparrow \Rightarrow y \uparrow \rightsquigarrow$ +ve covariance.
 $x \downarrow \Rightarrow y \uparrow \rightsquigarrow$ -ve covariance.
 $x \uparrow \downarrow \Rightarrow y \rightsquigarrow$ no covariance.
- ↳ The range of covariance is $-\infty$ to ∞ .

Correlation

- ↳ Correlation is a method of statistical evaluation used to study the strength of a relationship between two, numerically measured, continuous variables.
- ↳ Mathematically, Correlation = $\frac{\text{Cov}(x, y)}{\sigma_x * \sigma_y}$

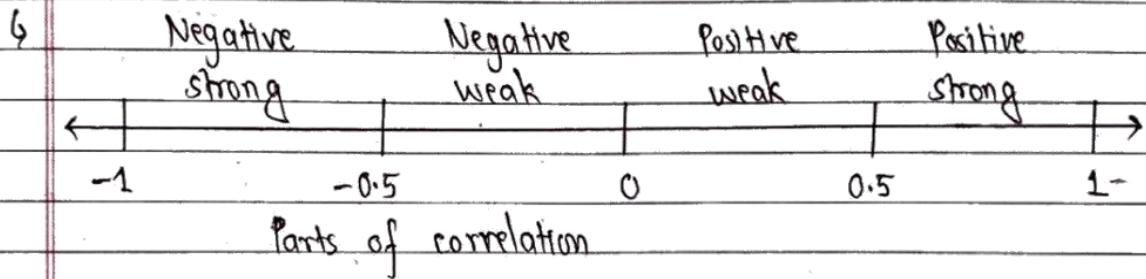
Where,

cov is covariance

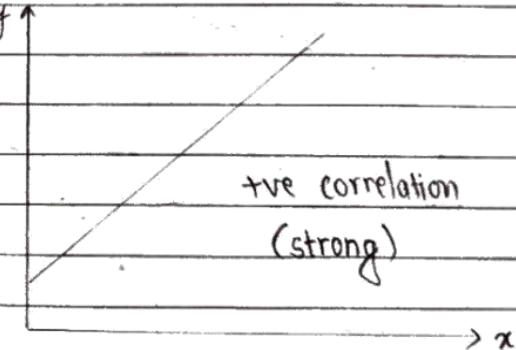
σ_x is standard deviation of x

σ_y is standard deviation of y .

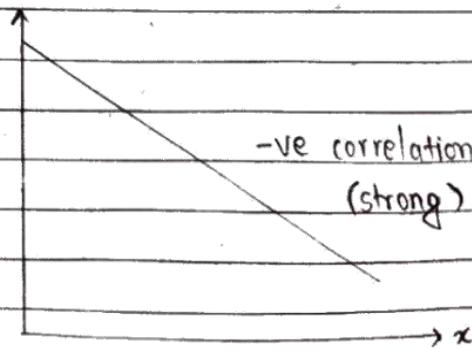
- 6) The range of correlation is -1 to 1 which is converted from range of covariance for better analysis.



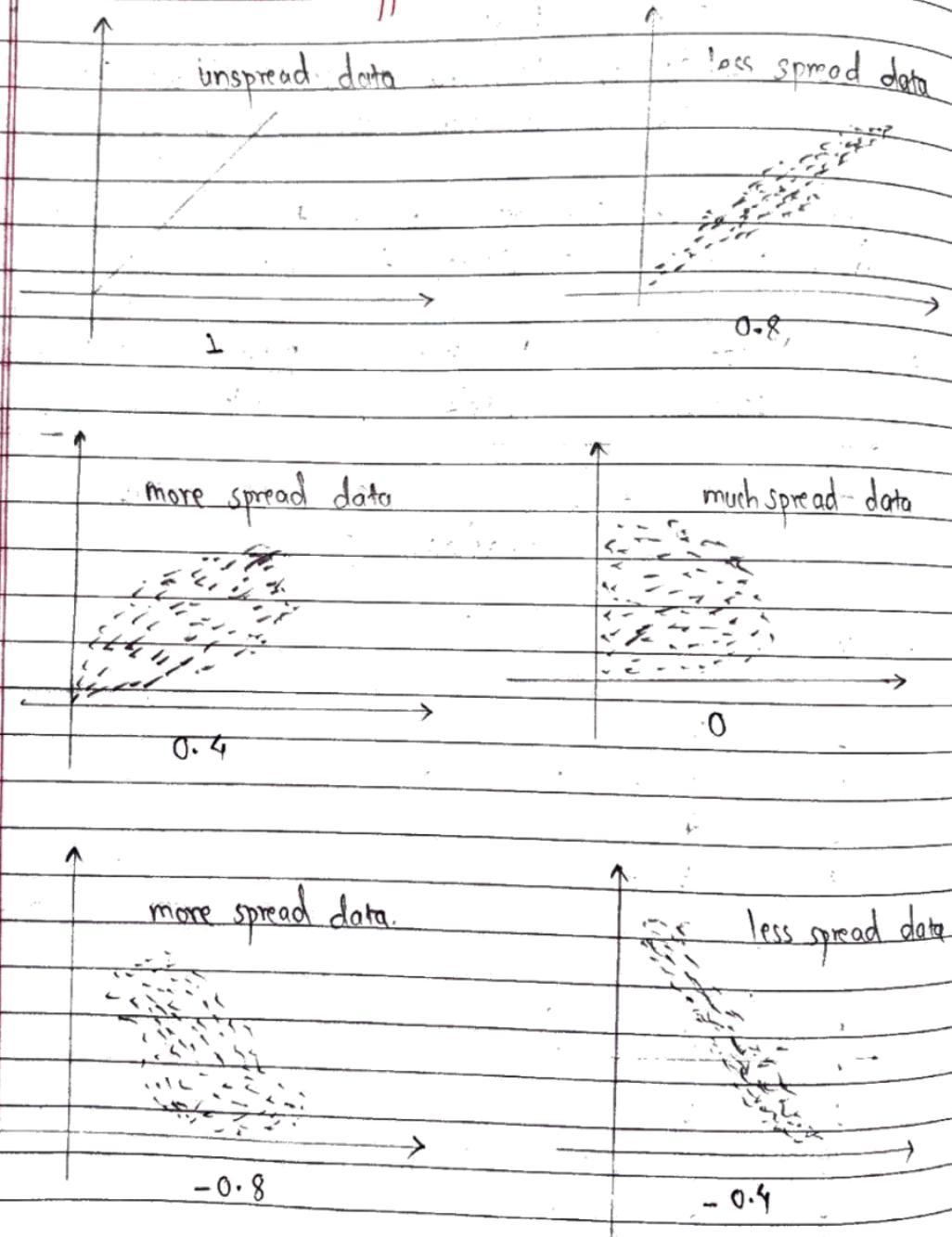
| x | y |
|-----|-----|
| 1 | 2 |
| 2 | 3 |
| 3 | 4 |
| 4 | 5 |
| 5 | 6 |

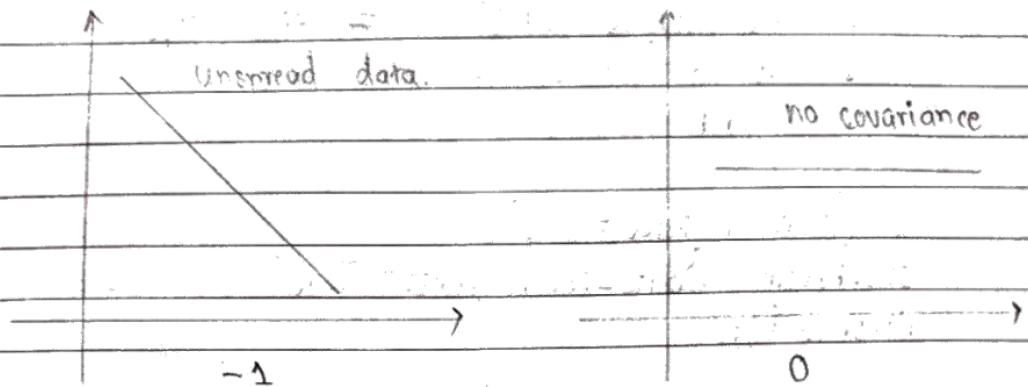


| x | y |
|-----|-----|
| 1 | 5 |
| 2 | 4 |
| 3 | 3 |
| 4 | 2 |
| 5 | 1 |



Pearson Correlation Coefficient:





Working Practically in Jupyter:

```
import pandas as pd
```

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

```
import seaborn as sns
```

```
dataset = pd.read_csv("Students.csv")
```

```
dataset.head(3) # display first three head rows.
```

```
dataset.isnull().sum() # displays sum of null values  
# in each columns.
```

```
dataset.info() # You can check data types using this.
```

```
data_corr = dataset.select_dtypes(["float64", "int64"]).corr()  
# find correlation on selected datatypes.
```

```
data_corr # displays what data_corr variable stores.
```

```
data_cov = dataset.select_dtypes(["float64", "int64"]).cov()
```

```
data_cov # for covariance
```

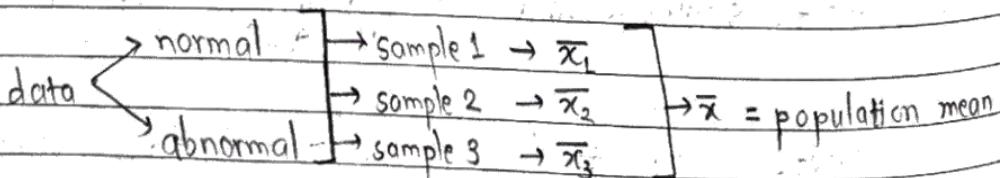
```
plt.figure(figsize = (4,3)) # adjust size  
sns.heatmap(data_corr, annot = True)  
plt.show();
```

```
plt.figure(figsize = (4,3))  
sns.heatmap(data_cov, annot = True)  
plt.show();
```

Correlation is more used than covariance in machine learning.

8. Central Limit Theorem

The central limit theorem (CLT) states that when plotting a sample distribution of means the means of the sample will be equal to the population mean and the sample distribution will approach normal distribution with variance equal to standard error.



There are few assumptions behind the CLT:

- The sample data must be sampled and selected randomly from the population.
- There shouldn't be any multicollinearity in the sampled data which is one sample should not influence the other samples.

- The sample size should be no more than 10% of the population. Generally, sample size greater than 30 ($n > 30$) is considered good.

Working practically in Jupyter:

```
import numpy as np  
import seaborn as sns  
import matplotlib.pyplot as plt  
import pandas as pd
```

```
pop_data = [np.random.randint(10, 100) for i in range  
(10000)]
```

```
# Creates a list of 10000 integers generated randomly.  
pop_table = pd.DataFrame ({ "Pop-data": pop_data })  
# Creates a data frame with column Pop-data and  
# rows pop-data.
```

```
np.mean(pop_data) # find mean of pop-data.
```

```
plt.figure(figsize=(5,5))
```

```
sns.kdeplot(x="Pop-data", data=pop_table)
```

```
plt.show();
```

```
sample_mean = []
```

```
for no_sample in range(60): # select 60 samples
```

```
sample_data = []
```

```
for data in range(500): # Select 500 data in
```

```
# each sample
```

```
sample_data.append(np.random.choice(pop_data))
```

```
sample_mean.append(np.mean(sample_data))
```

np.mean (sample_mean)

sample_M = pd.DataFrame ({ "Sample-Mean": sample_mean})

plt.figure (figsize = (5,5))

sns.kdeplot (x = "Sample-Mean", data = sample_M)

plt.show ()

We find mean of pop-data and sample_mean is
nearly equal which proves CLT.

9. Hypothesis Testing

- ↳ It is a part of statistical analysis, where we test the assumptions made regarding a population parameter.
- ↳ It is generally used when we were to compare a single group with an external standard and two or more groups with each other.

| Product A | Product B |
|-----------|-----------|
| 75% | 85% |

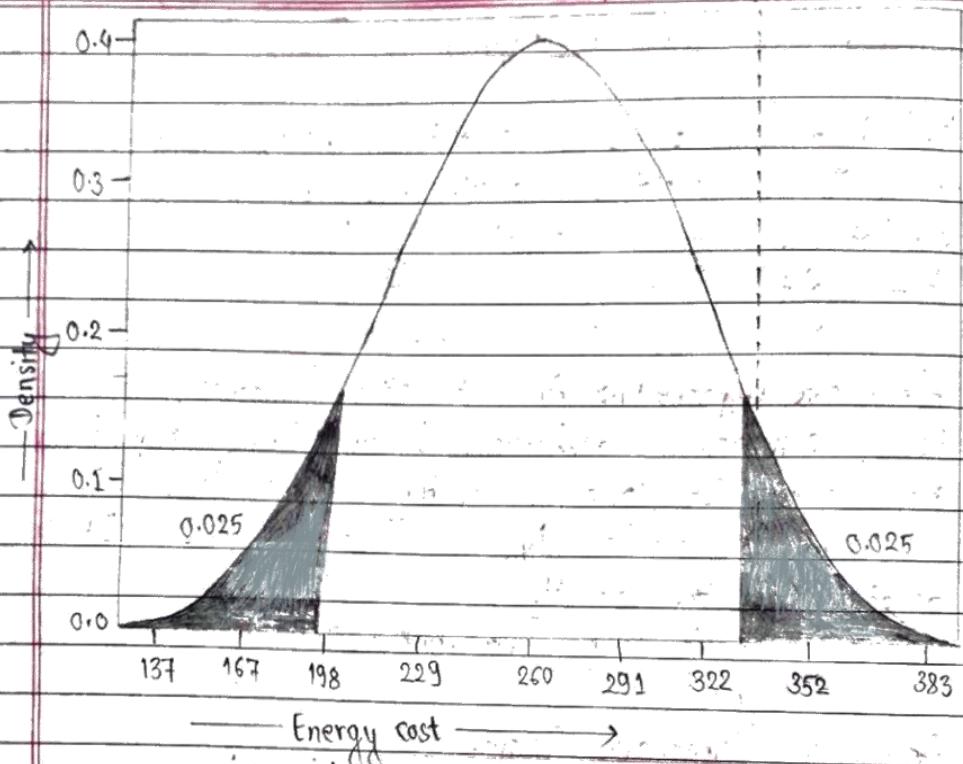
Termonology Used in Hypothesis Testing

- **Null Hypothesis :** is a statistical theory that suggests there is no statistical significance exists between the populations. It is denoted by H_0 and read as H-naught.
- **Alternative Hypothesis :** An alternative hypothesis suggests there is a significant difference between the population parameters. It could be greater or smaller. Basically, it is the contrast of "the" Null Hypothesis. It is denoted by H_a or H_1 .

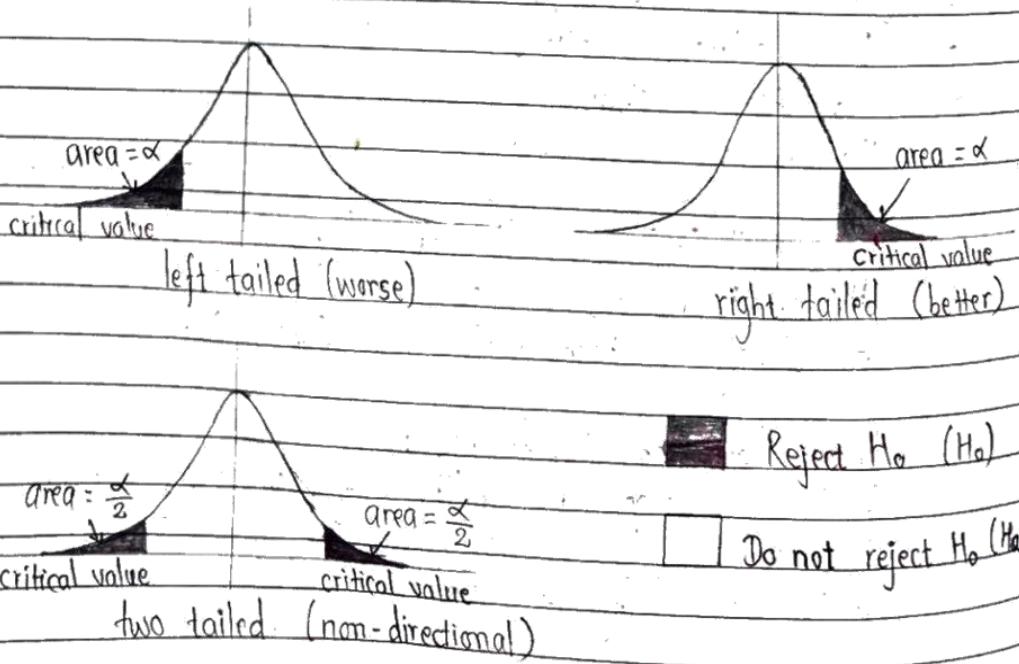
Steps of Hypothesis Testing

- State null (H_0) and alternative (H_1) hypothesis.
- Choose level of significance (α).
- Find critical values.
- Find test statistic.
- Draw your conclusion.
- **State null (H_0) and alternative (H_1) hypothesis :**
 H_0 must always contains equality (=).
 H_1 always contain differences ($\neq, >, <$).
- **Choose level of significance :**
It is a fixed probability of wrongly rejecting a True Null Hypothesis.
It is denoted by α (alpha).

330.5



Visual representation of level of significance

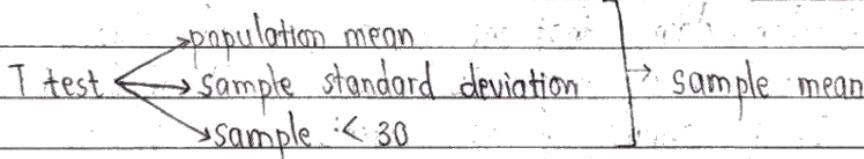
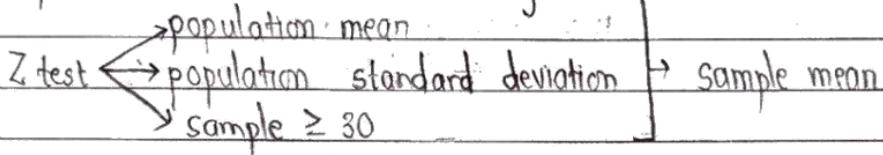


o Find critical values:

Using the test statistic find critical values like z value, t value, etc.

o Find test statistic:

Find the required test like: Z-test, T-test, Chi-square test, etc. These test gives critical values.



Chi-square test \rightarrow checks goodness and relations between two variables (independance)

o Draw your conclusion:

Conclude which one to be chosen either H_0 or H_a for better performance.

1. Z Test

↳ Data should be normally distributed.

- ↳ Given :
- population mean
 - population standard deviation
 - sample mean
 - no. of sample ≥ 30

6 Mathematically,

$$Z_{\text{test}} = \frac{\bar{x} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$

Where, \bar{x} = sample mean

μ = population mean

σ = population standard deviation

n = no. of sample.

Example 1:

A teacher claims that the mean score of students in his class is greater than 82 with a standard deviation of 20. If a sample of 81 students was selected with a mean score of 90.

Solution:

$H_0 : \mu \neq 82$ (Opposing / null hypothesis)

$H_a : \mu > 82$ (Supporting / alternative hypothesis)

Assume that $\alpha = 0.05$

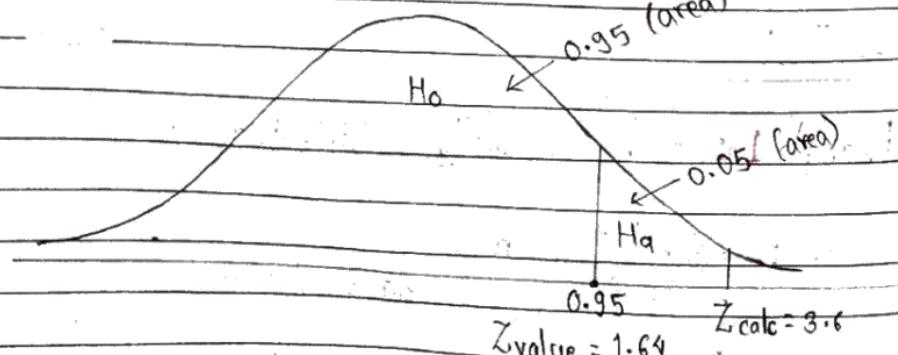
Given, $\bar{x} = 90$

$\sigma = 20$

$n = 81$

Total area under curve is always

1 being a probability distribution.



It is right tailed so, Z-value should be of $(1-\alpha)$

Z value = sum of row and column z values of $(1-\alpha)$ area
 from same tailed z-table (choose close value)

$$= 1.6 + 0.04 = 1.64$$

$$Z_{\text{calc}} = \frac{\bar{x} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)} = \frac{90 - 82}{\left(\frac{20}{\sqrt{81}}\right)} = \frac{8 \times 9}{20} = 3.6$$

$\therefore Z_{\text{calc}} > Z_{\text{value}}$, the teacher was right.

Working practically in Jupyter:

```
import scipy.stats as st # library with constants
# stats has statistical analysis
```

```
import numpy as np
```

sample_mean = 90

population_mean = 82

population_std = 20

n = 81

alpha = 0.05 # Assume

$Z_{\text{calc}} = (\text{sample-mean} - \text{population-mean}) / (\text{population-std} / \text{np.sqrt}(n))$

Z_{calc}

$Z_{\text{table}} = \text{st.norm.ppf}(1-\alpha)$ # norm contains functions
 # related to normal distribution

Z_{table}

if $Z_{\text{table}} < Z_{\text{calc}}$:

print ("Teacher is right")

else:

print ("Teacher is wrong")

Example 2:

Scenario :

Imagine you work for an e-commerce company and your team is responsible for analyzing customer purchase data. You want to determine whether a new website design has lead to a significant increase in the average purchase amount compared to the old design.

Data :

You have collected data from a random sample of 30 customers who made purchase on the old website design and 30 customers who made purchase on the new website design. You have the sample means, sample standard deviations, and sample sizes for both groups.

- Old design data = [45.2, 42.8, 38.9, 48.5, 41.0, 44.6, 40.5, 42.7, 39.8, 41.4, 44.3, 39.7, 42.1, 40.6, 43.0, 42.2, 41.5, 39.6, 44.0, 43.1, 38.7, 43.9, 42.0, 41.9, 42.8, 43.7, 41.3, 40.9, 42.5, 41.6]
- New design data = [48.5, 49.1, 50.2, 47.8, 48.7, 49.9, 48.0, 50.5, 49.8, 49.6, 48.2, 48.9, 49.7, 50.3, 49.4, 50.1, 48.6, 48.3, 49.0, 50.0, 48.4, 49.3, 49.5, 48.8, 50.6, 50.4, 48.1, 49.2, 50.7, 50.8]
- Population standard deviation = 2.5

Solution :

$$H_0 : \text{new website} = \text{old website}$$

$$H_a : \text{new website} > \text{old website}$$

Assuming no change in population mean i.e. $\mu_{\text{new}} = \mu_{\text{old}}$
 Given, $\sigma = 2.5$

$$n = 30$$

We need, \bar{x}_{new} and \bar{x}_{old} → calculate from given samples.

Now, assuming $\alpha = 0.05$:

Find Z value from Z table, Z value = 1.645

And,

$$Z_{\text{calc}} = \frac{(\bar{x}_{\text{new}} - \bar{x}_{\text{old}}) - (\mu_{\text{new}} - \mu_{\text{old}})}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$

Since, $\mu_{\text{new}} = \mu_{\text{old}}$

$$\therefore Z_{\text{calc}} = \frac{\bar{x}_{\text{new}} - \bar{x}_{\text{old}}}{\left(\frac{\sigma}{\sqrt{n}}\right)} = 1.645$$

Hence, you can draw your conclusion that H_0 is right
 i.e. new website is better than old website.

Working practically in Jupyter:

import scipy.stats as st

import numpy as np

Old-design-data = np.array([<Old design data>])

New-design-data = np.array([<New design data>])

np.array() converts lists to array for fast
 # calculation

pop_std = 2.5

n_sp = len(New-design-data) # 30

alpha = 0.05 # Assume

`mean_new = np.mean (New-design-data)`
`mean_old = np.mean (Old-design-data)`

`z_cal = (mean_new - mean_old) / (pop_std / np.sqrt(n_sp))`
`z_cal`

`z_table = st.norm.ppf (1-alpha)`
`z_table`

`if z_calc > z_table :`

`print (" New website is better than old. ")`
`else :`

`print (" New website is same as old. ")`

2. T Test

↳ Data should be normally distributed.

↳ Given : • population mean
 • sample mean
 • sample standard deviation
 • no. of sample ≤ 30

↳ Mathematically,

$$t_{\text{test}} = \frac{\bar{x} - \mu}{\left(\frac{s}{\sqrt{n}}\right)}$$

Where, \bar{x} = sample mean

μ = population mean

s = standard deviation of sample
 n = no. of sample

Example 1:

A manufacturer claims that the average weight of a bag of potato chips is 150 grams. A sample of 25 bags is taken and the average weight is found to be 148 grams, with a standard deviation of 5 grams. Test the manufacturer's claim using a one-tailed t-test with a significance level of 0.05.

Solution:

$H_0: \mu = 150 \text{ gm}$ (null hypothesis)

$H_a: \mu < 150 \text{ gm}$ (alternative hypothesis)

Given,

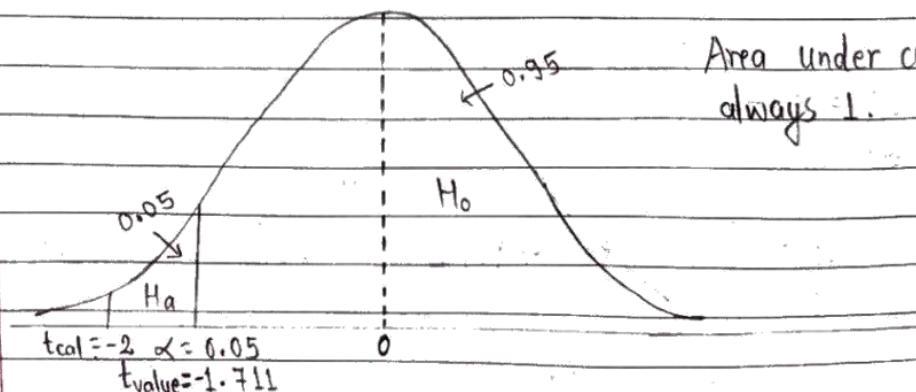
$$\alpha = 0.05$$

$$\mu = 150$$

$$\bar{x} = 148$$

$$s = 5$$

$$n = 25$$



Area under curve is always 1.

It is left tailed so, tvalue should be of α .

degree of freedom, $df = n - 1 = 25 - 1 = 24$.

t_{value} = the value that is present in df row and α column in t-table.

= -1.711 (-ve sign indicates it is left tailed)

$$t_{cal} = \frac{\bar{x} - \mu}{\left(\frac{s}{\sqrt{n}} \right)} = \frac{148 - 150}{\left(\frac{5}{\sqrt{25}} \right)} = \frac{148 - 150}{1} = -2$$

Hence, the claims of the company is wrong as we find $t_{cal} < t_{value}$.

Working practically in Jupyter:

```
import scipy.stats as st
```

```
import numpy as np
```

```
alpha = 0.05
```

```
pop_mean = 150
```

```
Sam_mean = 148
```

```
Sam_std = 5
```

```
n = 25
```

$$df = n - 1$$

$$df$$

$$t_{-calc} = (\text{Sam_mean} - \text{pop_mean}) / (\text{Sam_std} / \text{np.sqrt}(n))$$

$$t_{-table} = st.t.ppf(\alpha, df)$$

$$t_{-table}$$

if $t_{\text{calc}} > t_{\text{table}}$:

print ("The claim of manufacturer is right")

else:

print ("The claim of manufacturer is wrong")

Example 2:

A company wants to test whether there is a difference in productivity between two teams. They randomly select 20 employees from each team and record their productivity scores. The mean productivity for Team A is 80 with a standard deviation of 5, while the mean productivity score for Team B is 75 with a standard deviation of 6. Test at a 5% level of significance whether there is a difference in productivity between two teams.

Solution:

H_0 : Productivity of A - Productivity of B = 0

H_a : Productivity of A - Productivity of B $\neq 0$ (two tailed)

Given,

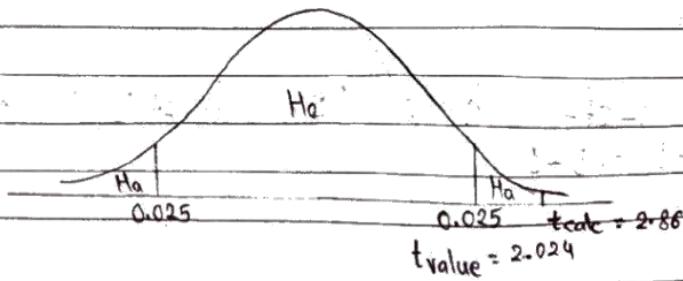
$$2\alpha = 5\% = 0.05 \Rightarrow \alpha = 0.025 \quad [:\text{two tailed}]$$

$$n_A = n_B = 20$$

$$\bar{x}_A = 80 \text{ and } s_A = 5$$

$$\bar{x}_B = 75 \text{ and } s_B = 6$$

Assume that $\mu_A = \mu_B$.



$$\text{degree of freedom, } df = n_A - 1 + n_B - 1 \\ = 20 - 1 + 20 - 1 \\ = 38$$

$t\text{value} = -2.024$ [From t-table]

$$t_{\text{calc}} = \frac{(\bar{x}_A - \bar{x}_B) - (\mu_A - \mu_B)}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}$$

Since, $\mu_A = \mu_B$

$$\therefore t_{\text{calc}} = \frac{\bar{x}_A - \bar{x}_B}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}} = 2.86$$

Since, $t_{\text{calc}} > t\text{value}$ there is some difference between productivity of A and productivity of B.

Working practically in Jupyter:

```
import numpy as np
from scipy.stats import t
```

$$\text{level-of-significance} = 5/100 \# 5\%$$

$$n_A = n_B = 20$$

$$\text{mean_A} = 80$$

$$\text{std_A} = 5$$

$$\text{mean_B} = 75$$

$$\text{std_B} = 6$$

$$\alpha = \text{level-of-significance} / 2 \# \text{It is two tailed}$$

$$df = n_A - 1 + n_B - 1$$

```

t-table = t.ppf(alpha, df) # Testing for left tail
t-calc = (mean_A - mean_B) / (np.sqrt((np.square(std_A)
                                         / n_A) + (np.square(std_B) / n_B)))
print("t-table =", t-table, "and t-calc =", t-calc)

if t-calc > t-table:
    print("Productivity of A != Productivity of B")
else:
    print("Productivity of A = Productivity of B")

```

Example 3:

A company wants to test whether a new training program improves the typing speed of its employees. The typing speed of 20 employees was recorded before and after the training program. The data is given below. Test at a 5% level of significance whether the training program has an effect on the typing speed of the employees.

- Before : 50, 60, 45, 65, 55, 70, 40, 75, 80, 65, 70, 60, 50, 55, 45, 75, 60, 50, 65, 70
- After : 60, 70, 55, 75, 65, 80, 50, 85, 90, 70, 75, 65, 55, 60, 50, 80, 65, 55, 70, 75

Solution:

$$H_0 : \text{Speed before training} - \text{Speed after training} = 0$$

$$H_a : \text{Speed before training} - \text{Speed after training} \neq 0$$

Given,

$$\alpha = 5\% = 0.05 \Rightarrow \frac{\alpha}{2} = 0.025 \quad [\text{two tailed test}]$$

$$n = 20$$

degree of freedom, $df = n - 1 = 20 - 1 = 19$

Assume that $\mu_{\text{after}} = \mu_{\text{before}}$ [\because not given]

Calculate \bar{x}_{before} , \bar{x}_{after} , s_{before} and s_{after} from the given sample data.

t-value = 2.093 [from t-table]

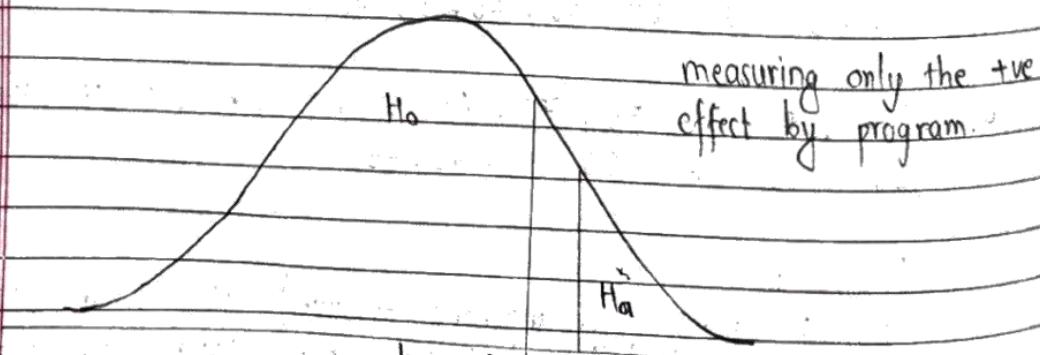
$$t\text{-calc} = (\bar{x}_{\text{after}} - \bar{x}_{\text{before}}) - (\mu_{\text{after}} - \mu_{\text{before}})$$

$$\sqrt{\frac{s_{\text{after}}^2}{n_{\text{after}}} + \frac{s_{\text{before}}^2}{n_{\text{before}}}}$$

$\therefore \mu_{\text{after}} = \mu_{\text{before}}$ and $n_{\text{after}} = n_{\text{before}} = n$

$$\therefore t\text{-calc} = \frac{\bar{x}_{\text{after}} - \bar{x}_{\text{before}}}{\sqrt{\frac{s_{\text{after}}^2}{n} + \frac{s_{\text{before}}^2}{n}}} = 2.061$$

Clearly, $2.061 < 2.093$ i.e. $t\text{-calc} < t\text{-value}$ so, the training program doesn't have any effect on typing speed of employees.



two tailed but analyzing only right tail.

Working practically in Jupyter:

```
import numpy as np
from scipy.stats import t
```

before = np.array ([< before data >])

after = np.array ([< after data >])

level_of_significance = 5 / 100 # 5%

n = 20

alpha = level_of_significance / 2 # It is two tailed

df = n - 1

mean_before = np.mean (before)

mean_after = np.mean (after)

std_before = stnd.std (before)

std_after = np.std (after)

t_table = t.ppf (1 - alpha, df) # Testing for right tail

t_calc = (mean_after - mean_before) / np.sqrt ((np.square (std_after) + np.square (std_before)) / n)

print (f"t-table = {t_table} and t-calc = {t_calc}")

if t_calc > t_table :

print ("The training program is effective")

else :

print ("The training program is not effective")

3. Chi Square Test

↳ Applicable in two types of situations :

i. Goodness

ii. Independance

↳ Mathematically,

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Where, χ^2 = chi square

O_i = observed value

E_i = expected value

Example 1:

A fair dice is rolled 120 times and the following results are obtained :

Face 1 : 22 times

Face 2 : 17 times

Face 3 : 20 times

Face 4 : 26 times

Face 5 : 22 times

Face 6 : 13 times

Test at 5% level of significance whether the dice is fair.

Solution:

H_0 : dice is fair

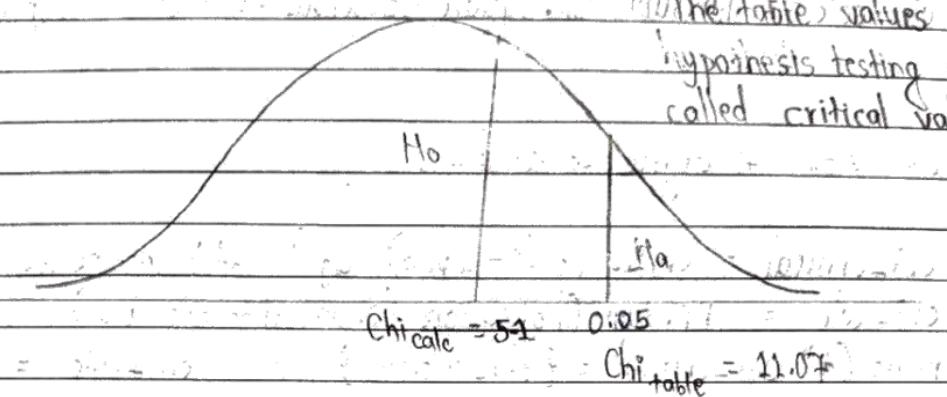
H_a : dice is not fair

In dice, number of faces, $n = 6$

Degree of freedom, $df = n-1 = 6-1 = 5$

level of significance, $\alpha = 5\% = 0.05$

The table values for the hypothesis testing are also called critical values.



The expected value for each face when a dice is rolled 120 times is 20 for each.
The observed values are given above.

Chi_{table} = the value that is present in df row and
& column in chi square table.
= 11.07

$$\text{Chi}_{\text{calc}} = \sum \frac{(O_i - E_i)^2}{E_i} = 5.1$$

Clearly Chi_{calc} < Chi_{table} so, the dice is fair i.e.
H₀ is right.

Working practically in Jupyter

```
import numpy as np
from scipy.stats import chi2
```

observed = np.array([22, 17, 20, 26, 22, 18])

expected = np.array([20, 20, 20, 20, 20, 20])

$\alpha = 0.05$ # Given level of significance is 5%
 $n = 6$ # Total number of possibilities when a dice is
 # rolled

$df = n - 1$ # degree of freedom

$\text{chi-critical} = \text{chi2.ppf}(1 - \alpha, df)$ # Right tailed.

$\text{chi-calc} = \text{np.sum}(\text{np.square}(\text{observed} - \text{expected}) / \text{expected})$
 $\text{print}(f"\text{chi-critical} = \{\text{chi-critical}\} \text{ and chi-calc} = \{\text{chi-calc}\})$

if $\text{chi-calc} > \text{chi-critical}$:

print ("Dice is not fair")

else:

print ("Dice is fair")

Example 2.

A study was conducted to investigate whether there is a relationship between gender and the preferred genre of music. A sample of 235 people was selected, and the data collected is shown below. Test at a 5% level of significance whether there is significant association between gender and music preference.

| | Pop | Hip Hop | Classical | Rock |
|--------|-----|---------|-----------|------|
| Male | 90 | 45 | 25 | 10 |
| Female | 35 | 30 | 20 | 30 |

Solution:

H_0 : No association

H_a : Has association

Given:

no. of Sample, $n = 235$

level of significance, $\alpha = 5\% = 0.05$

no. of rows, rows = 2

no. of columns, columns = 4

Now,

$$\begin{aligned} df &= (\text{rows} - 1) * (\text{columns} - 1) \\ &= (2 - 1) * (4 - 1) \\ &= 3 \end{aligned}$$

So,

$\chi^2_{\text{critical}} = 7.815$ [From chi-square table.]

We are given the observed values in the given table.
To find the expected value for each observed value proceed as follows:

- i.) find the sum of each column and row.
- ii.) expected value $(i, j) = \frac{\sum R_i \times \sum C_j}{\text{No. of sample}}$

So,

$$\sum R_1 = 120 \text{ and } \sum R_2 = 115$$

$$\sum C_1 = 75, \sum C_2 = 75, \sum C_3 = 45 \text{ and } \sum C_4 = 40$$

We're,

$$\text{observed values} = [40, 45, 25, 10, 35, 30, 20, 30]$$

for expected values:

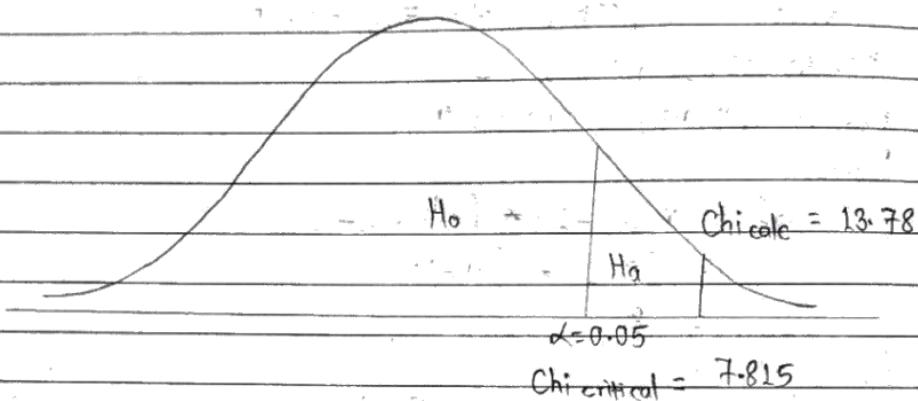
$$\text{expected value } (1,1) = \frac{\sum R_1 \times \sum C_1}{n} = 38.29$$

$$\text{expected value } (1,2) = \frac{\sum R_1 \times \sum C_2}{n} = 38.29$$

Similarly, all the expected values can be calculated.
Finally,

$$\text{expected values} = [38.29, 38.29, 22.97, 20.42, 36.70, 36.70, 22.02, 19.57]$$

$$\chi^2_{\text{calc}} = \sum \frac{(O_i - E_i)^2}{E_i} = 13.78$$



Clearly $\chi^2_{\text{calc}} > \chi^2_{\text{critical}}$ so, there is association between gender and music preference.

Working practically in Jupyter:

```
import numpy as np
from scipy.stats import chi2
```

$\alpha = 0.05$ # Given level of significance is 5%
 $n = 235$ # Sample of 235 people was selected.

rows = 2 # From table
columns = 4

row1 = np.array([40, 45, 25, 10])
row2 = np.array([35, 30, 20, 30])

df = (rows - 1) * (columns - 1) # Calculate degree of freedom

```
chi_critical = chi2.ppf (alpha, df)
print ("chi_critical = ", chi_critical)
```

```
sum_row1 = np.sum (row1)
```

```
sum_row2 = np.sum (row2)
```

```
sum_rows = np.array ([sum_row1, sum_row2]) # Create an
# array of sum of rows.
```

```
sum_rows
```

```
sum_columns = row1 + row2. # find sum of all columns at
once in an array.
sum_columns
```

```
# find expected values
```

```
expected = []
```

```
for i in sum_rows :
```

```
    for j in sum_columns :
```

```
        expected.append (i*j/n)
```

```
expected
```

```
# Get observed values
```

```
observed = np.append (row1, row2)
```

```
observed
```

```
chi_calc = np.sum (np.square (observed - expected) / expected )
```

```
if chi_calc > chi_critical :
```

```
    print ("There is association")
```

```
else:
```

```
    print ("There is no association")
```