



# Applied Computational Statistics

## Lab Manual

Department of Computer Science and  
Engineering  
The NorthCap University, Gurugram

# **Applied Computational Statistics**

## **Lab WorkBook**

### **CSL227**

**Ms.Sonal Saurabh**

**Dr. Srishti**



**Name: | Roll No:**

Department of Computer Science and Engineering

The NorthCap University, Gurugram- 122001, India

Session 2022-23

*Published by:*

**Department of Computer Science & Engineering  
School of Engineering and Technology  
The NorthCap University Gurugram**

- **Laboratory Manual is for Internal Circulation only**

#### © Authors

*No part of this Practical Record Book may be reproduced, used, stored without prior permission of The NorthCap University*

Copying or facilitating copying of lab work comes under cheating and is considered as use of unfair means. Students indulging in copying or facilitating copying shall be awarded zero marks for that particular experiment. Frequent cases of copying may lead to disciplinary action. Attendance in lab classes is mandatory.

Labs are open up to 7 PM upon request. Students are encouraged to make full use of labs beyond normal lab hours.

## PREFACE

Applied Computational Statistics Laboratory Manual is designed to meet the course and program requirements of NCU curriculum for B.Tech. fourth semester students of CSE branch. The concept of the lab work is to give brief practical experience for basic lab skills to students. It provides the space and scope for self-study so that students can come up with new and creative ideas.

The Lab manual is written on the basis of “teach yourself pattern” and expected that students who come with proper preparation should be able to perform the experiments without any difficulty. A brief introduction to each experiment with information about self-study material is provided. The laboratory exercises will help students to provide a hands-on each exercise that will help them to understand thoroughly. The students are expected to come thoroughly prepared for the lab. General disciplines, safety guidelines and report writing are also discussed.

The lab manual is a part of curriculum for the The NorthCap University, Gurugram. Teacher's copy of the experimental results and answer for the questions are available as sample guidelines.

We hope that lab manual would be useful to students of CSE branch and author requests the readers to kindly forward their suggestions / constructive criticism for further improvement of the work book.

Author expresses deep gratitude to Members, Governing Body-NCU for encouragement and motivation.

**Authors**  
**The NorthCap University**  
**Gurugram, India**

## CONTENTS

S.No.	Details	Page No.
1	<b>Introduction</b>	vi
2	<b>Lab Requirement</b>	vii
3	<b>General Instructions</b>	viii
4	<b>List of Experiments</b>	x
5	<b>List of Flip Experiments</b>	xi
6	<b>List of Projects</b>	xii
7	<b>Rubrics</b>	xiv
8	<b>Annexure 1 (Format of Lab Report)</b>	xv
9	<b>Annexure 2 (Format of Project Report)</b>	xlviii

## 1. INTRODUCTION

That 'learning is a continuous process' cannot be over emphasized. The theoretical knowledge gained during lecture sessions need to be strengthened through practical experimentation. Thus, practical makes an integral part of a learning process.

### COURSE OBJECTIVES:

- 1. Apply appropriate descriptive statistical and exploratory methods in the analysis of datasets.**
- 2. Visualize the data for the various analysis.**
- 3. Understand the probability mass function and various discrete distributions through application on real world examples.**
- 4. Understand the probability density function and various continuous distributions through application on real world examples.**
- 5. Understand and interpret statistical hypothesis test.**
- 6. Translate real-world problems into probability models using Bayesian Statistics.**

## 2. LAB REQUIREMENTS

S.No.	Requirements	Details
1	<b>Software Requirements</b>	Python 3.x, Numpy, Pandas, Matplotlib, Seaborn, statistics, sci-kit learn
2	<b>Operating System</b>	Windows 7 onwards or Linux (32 or 64 bit)
3	<b>Hardware Requirements</b>	4 GB RAM (Recommended) 2.60 GHz (Recommended)
4	<b>Required Bandwidth</b>	NA

### 3. GENERAL INSTRUCTIONS

#### a. General discipline in the lab

- Students must turn up in time and contact concerned faculty for the experiment they are supposed to perform.
- Students will not be allowed to enter late in the lab.
- Students will not leave the class till the period is over.
- Students should come prepared for their experiment.
- Experimental results should be entered in the lab report format and certified/signed by concerned faculty/ lab Instructor.
- Students must get the connection of the hardware setup verified before switching on the power supply.
- Students should maintain silence while performing the experiments. If any necessity arises for discussion amongst them, they should discuss with a very low pitch without disturbing the adjacent groups.
- Violating the above code of conduct may attract disciplinary action.
- Damaging lab equipment or removing any component from the lab may invite penalties and strict disciplinary action.

#### b. Attendance

- Attendance in the lab class is compulsory.
- Students should not attend a different lab group/section other than the one assigned at the beginning of the session.
- On account of illness or some family problems, if a student misses his/her lab classes, he/she may be assigned a different group to make up the losses in consultation with the concerned faculty / lab instructor. Or he/she may work in the lab during spare/extra hours to complete the experiment. No attendance will be granted for such case.

#### c. Preparation and Performance

- Students should come to the lab thoroughly prepared on the experiments they are assigned to perform on that day. Brief introduction to each experiment with information about self -study reference is provided on LMS.
- Students must bring the lab report during each practical class with written records of the last experiments performed complete in all respect.

- Each student is required to write a complete report of the experiment he has performed and bring to lab class for evaluation in the next working lab. Sufficient space in work book is provided for independent writing of theory, observation, calculation and conclusion.
- Students should follow the Zero tolerance policy for copying / plagiarism. Zero marks will be awarded if found copied. If caught further, it will lead to disciplinary action.
- Refer **Annexure 1** for Lab Report Format

#### 4. LIST OF EXPERIMENTS

Sr. No.	Title of the Experiment	Software used	Unit Covered	CO Covered	Time Required
1.	Introduction To Statistics using Python	Python	1	C01	2 Hours
2.	Demonstration of descriptive statistics	Python (Jupyter)	1	C01	2 hours
3.	Demonstrate Measure of Shape using Python	Python (Jupyter)	1	C01	2 hours
4.	Data visualization for qualitative and quantitative data	Python (Jupyter)	1	C01	2 hours
5.	Basics of Probability using Python	Python (Jupyter)	1	C01	2 hours
6.	Write python program to implement laws of Probability	Python (Jupyter)	2	C02	2 hours
7.	Implementation of Probability Mass Function for Discrete Random Variable using Python	Python (Jupyter)	3	C03	2 hours
8.	Write a python program to find expectation and variance in discrete distribution	Python (Jupyter)	3	C03	2 hours
9.	Write Python Program for Continuous Probability Distributions	Python (Jupyter)	3	C04	2 hours
10	Write python code to find critical value in normal distribution	Python (Jupyter)	4	C04	2 hours
11	Implement central limit theorem	Python (Jupyter)	5	C05	2 hours
12	Write a python program to conduct a hypothesis test for unknown variance.	Python (Jupyter)	5	C05	2 hours
13	Write a python program to conduct hypothesis testing for two populations	Python (Jupyter)	5	C05	2 hours
<b>Value Added Experiments</b>					
14	Write a python program to bootstrap your data.	Python (Jupyter)	5	C05	2 hours
15	Write a python program to find confidence interval for known variance.	Python (Jupyter)	5	C05	2 hours
16	Write a python program to conduct a hypothesis test for known variance	Python (Jupyter)	5	C05	2 hours

## 5. LIST OF FLIP EXPERIMENTS

Exp. No.	Title of the Experiment	Mapped CO
1.	Implement the ANOVA and ANCOVA analysis of real-world data	CO5
2.	Implement the Bayes Theorem in Python	CO2
3.	Competitions on Kaggle	CO1-6

## 6. LIST OF PROJECTS

Sr No.	Project Title	Mapped CO
1.	<p>2008_Swing_States dataset (Available on Kaggle) contains US swing states election results. There are 221 observations in the dataset with 6 features each (state, county, total_votes, dem_votes, rep_votes, dem_share). What is the highest dem_share percentage for different states? Find the country with highest dem_share in each state? Also check for null values. What is the average total_votes for each state? Also calculate their respective median. Calculate the variability and skewness (whether positive or negative or zero) for total_votes for each state. Also calculate the z - score of total_votes for each states. Calculate Interquartile range of total_votes of state PA.</p>	CO1
2.	<p>Suicide Rates Overview 1985 to 2016 dataset (Available on Kaggle). This compiled dataset pulled from four other datasets linked by time and place, and was built to find signals correlated to increased suicide rates among different cohorts globally, across the socio-economic spectrum. Our aim is to make general analysis of suicide rates and we will explain these rates in detail. Given the data, after removing null values, analyse and perform the EDA. How many males and females are there in dataset? Also find out how many males and females are there based on minimum and maximum gdp per capita? How many males and females are there in the data categorise by age parameter? Count the number of people based on Age values and plot. (Hint: Use groupby on age). How many generations are there in the data? Also based on sex parameter? plot a five number summary for three parameters - population, gdp per capita and year. Plot a line graph for the statistics of data. *(Hint: plot on data.describe() where kind = area). Draw two graphs on same canvas for generation parameter. First graph should be **Pie Chart** and Second graph should be Count plot. year on x-axis and suicides per100k of pop on y-axis (scatter plot). year on x-axis and gdp per capita on y-axis (reg plot). year on x-axis and gdp per capita on y-axis (kde plot).</p>	CO2

	year on x-axis and gdp per capita on y-axis (hex plot). Draw FacetGrid between gdp per capita and population. Visualise the graph between year and suicides using strip plot. Plot graph between each unique values of generation parameter. Also use violin plot between generation and population. Visualize the correlation between data.	
3.	Apply descriptive statistics on Insurance dataset.	CO1, CO2
4.	Do exploratory data analysis on University graduate admission dataset.	CO2
5.	Apply ANOVA (analysis of variance) on a given dataset.	CO3
6.	Apply ANCOVA (analysis of covariance) on a given dataset.	CO3
7.	Apply PCA (principal component analysis) on cancer dataset for dimensionality reduction.	CO4
8.	Do exploratory data analysis on IPL cricket dataset and calculate the chance of winning for each team	CO2
9.	Generate a pair of hypotheses for a given case study and statistically evaluate it	CO4, CO5
10.	Evaluate null hypothesis using p-value method	CO4, CO5

## 7. RUBRICS (Only for Lab components)

Marks Distribution	
<b>Continuous Evaluation (25 Marks)</b> Each experiment shall be evaluated for 5 marks and at the end of the semester proportional marks shall be awarded out of total 25.	<b>Project Evaluations (20 Marks)</b> Project shall be evaluated for 20 marks and at the end of the semester viva will be conducted related to the project.
<b>Viva and Reporting (25 Marks)</b> Following is the breakup of 25 marks for each <b>10 Marks:</b> Observation & conduct of experiment. Teacher may ask questions about experiment in mid-term viva. <b>10 Marks:</b> Observation & conduct of experiment. <b>5 Marks:</b> For report writing	

**Annexure 1**

# **CSL227**

## **Lab Work Book**



Faculty name: Dr. Sujata

Student Name: Piyush Gambhir

Roll No.: 21CSU349

Semester: Semester IV

Group: AL-3

Department of Computer Science and Engineering

The NorthCap University, Gurugram- 122001, India

Session 2022-23

## INDEX

Student name: Piyush Grambhua  
 Roll No.: 21CSU389

### INDEX

S.No.	Experiment	Date of Experiment	Date of Submission	Marks	Signature
1	Introduction to statistics using Python	20-01-23	20-01-23		20/1/23
2	Demonstration of descriptive statistical measures	03-02-23	03-02-23		2/2/23
3	Demonstrate measures of shape using Python	04-02-23	04-02-23		2/2/23
4	Data visualization for qualitative and quantitative data	10-02-23	17-02-23		12/2/23
5	Basics of probability using Python	24-02-23	24-02-23		24/2/23
6	To implement laws of probability	03-03-23	03-03-23		2/3/23
7	Implementation of probability mass function for discrete random variable using Python.	17-03-23	17-03-23		17/3/23
8	Write a python program to find expectation and variance in discrete distribution	17-03-23	17-03-23		17/3/23
9	Write Python Program for Continuous Probability Distributions	24-03-23	24-03-23		24/3/23
10	Write python code to find critical value in normal distribution	14-04-23	14-04-23		14/4/23
11	Write a python program to Implement central limit theorem	28-04-23	28-04-23		28/4/23
12	Write a python program to conduct a hypothesis test for unknown variance.	28-04-23	28-04-23		28/4/23
	13. Write a python program to conduct hypothesis testing for two populations	28-04-23	28-04-23		28/4/23

## EXPERIMENT NO. 1

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 20.01.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> Familiarization with statistics libraries in Python
<b>Outcome:</b> Revision of the libraries Numpy, Pandas, Matplotlib, Seaborn and exploration of Statistics and Scipy.Statistics
<b>Problem Statement:</b> Introduction to statistics using Python
<b>Background Study:</b> Python has libraries like Statistics and SciPy. Statistics which contain functions for several descriptive and inferential statistics tasks which can be of help to the students.
<b>Question Bank:</b> <ol style="list-style-type: none"><li>What is difference between differential and inferential statistics? Differential statistics deals with analyzing and describing data in a sample or population without making any generalizations to a larger population. It is concerned with describing the characteristics of a sample or population and includes measures such as mean, median, mode, variance, standard deviation, etc.  On the other hand, inferential statistics involves making inferences or generalizations about a population based on a sample of data. It is used to estimate parameters of the population using sample statistics, such as mean and standard deviation, and to test hypotheses about the population parameters. Inferential statistics is concerned with making predictions and drawing conclusions about the population based on the sample data.</li><li>What is the relevance of statistics in everyday life? Statistics is an important field because it helps us understand the general trends and patterns in a given data set. Statistics can be used for analyzing data and drawing conclusions from it. It can also be used for making predictions about future events and behaviors.</li><li>What are some of the common tasks you can accomplish using the statistics libraries you have studied in Python? It provides a wide range of tools and algorithms for machine learning and artificial intelligence tasks, including supervised, unsupervised, reinforcement, and deep learning.</li></ol> <b>Applications:</b> <ul style="list-style-type: none"><li>Pattern recognition.</li><li>Time-series prediction.</li><li>Reinforcement learning.</li><li>Natural language processing</li></ul>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Output

#### Experiment 1

Introduction to statistics using Python.

```

1 # Importing the libraries
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import numpy as np
[63]                                         Python

1 # Importing the dataset
2 dataset = pd.read_csv('diabetes_dataset.csv')
[64]                                         Python

1 # Checking the top 5 rows of the dataset
2 print("Top 5 rows of the dataset:")
3 print(dataset.head().to_markdown())
[71]                                         Python

... Top 5 rows of the dataset:
|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|-----:|-----:|-----:|-----:|-----:|-----:|-----:|-----:|-----:|
| 0 |       6 |    148 |         72 |        35 |       0 |  33.6 |      0.627 |    50 | Diabetic |
| 1 |       1 |     85 |         66 |        29 |       0 |  26.6 |      0.351 |    31 | Non-Diabetic |
| 2 |       8 |    183 |         64 |        0 |       0 |  23.3 |      0.672 |    32 | Diabetic |
| 3 |       1 |     89 |         66 |        23 |      94 |  28.1 |      0.167 |    21 | Non-Diabetic |
| 4 |       0 |    137 |         40 |        35 |     168 |  43.1 |      2.288 |    33 | Diabetic |

1 # Replacing 0 and 1 in the Outcome column with Non-Diabetic and Diabetic respectively
2 dataset['Outcome'] = dataset['Outcome'].replace(0, 'Non-Diabetic')
3 dataset['Outcome'] = dataset['Outcome'].replace(1, 'Diabetic')
4 # Checking the top 5 rows of the dataset
5 print("Top 5 rows of the dataset:")
6 print(dataset.head().to_markdown())
[66]                                         Python

... Top 5 rows of the dataset:
|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|-----:|-----:|-----:|-----:|-----:|-----:|-----:|-----:|-----:|
| 0 |       6 |    148 |         72 |        35 |       0 |  33.6 |      0.627 |    50 | Diabetic |
| 1 |       1 |     85 |         66 |        29 |       0 |  26.6 |      0.351 |    31 | Non-Diabetic |
| 2 |       8 |    183 |         64 |        0 |       0 |  23.3 |      0.672 |    32 | Diabetic |
| 3 |       1 |     89 |         66 |        23 |      94 |  28.1 |      0.167 |    21 | Non-Diabetic |
| 4 |       0 |    137 |         40 |        35 |     168 |  43.1 |      2.288 |    33 | Diabetic |

1 # Checking the information of the dataset
2 print("Information of the dataset:")
3 print(dataset.info())
[67]                                         Python

... Information of the dataset:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
 ---  -- 
 0   Pregnancies      768 non-null   int64  
 1   Glucose          768 non-null   int64  
 2   BloodPressure    768 non-null   int64  
 3   SkinThickness    768 non-null   int64  
 4   Insulin          768 non-null   int64  
 5   BMI              768 non-null   float64 
 6   DiabetesPedigreeFunction 768 non-null   float64 
 7   Age              768 non-null   int64  
 8   Outcome          768 non-null   object  
dtypes: float64(2), int64(6), object(1)
memory usage: 54.1+ KB
None

```

```
[68] 1 # Checking the statistical description of the dataset
2 print("Statistical description of the dataset:")
3 print(dataset.describe().to_markdown())
```

Python

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
count	768	768	768	768	768	768	768	768
mean	3.84505	120.895	69.1055	20.5365	79.7995	31.9926	0.471876	33.2409
std	3.36958	31.9726	19.3558	15.9522	115.244	7.88416	0.331329	11.7602
min	0	0	0	0	0	0	0.078	21
25%	1	99	62	0	0	27.3	0.24375	24
50%	3	117	72	23	30.5	32	0.3725	29
75%	6	140.25	80	32	127.25	36.6	0.62625	41
max	17	199	122	99	846	67.1	2.42	81

```
[69] 1 # Checking the number of null values in the dataset
2 null_values = dataset.isnull().sum()
3 print(null_values)
```

Python

```
[69] ... Pregnancies      0
Glucose          0
BloodPressure    0
SkinThickness    0
Insulin          0
BMI              0
DiabetesPedigreeFunction 0
Age              0
Outcome          0
dtype: int64
```

```
[70] 1 # Checking the number of unique values in the dataset
2 unique_values = dataset.nunique()
3 print(unique_values)
```

Python

```
[70] ... Pregnancies      17
Glucose          136
BloodPressure    47
SkinThickness    51
Insulin          186
BMI              248
DiabetesPedigreeFunction 517
Age              52
Outcome          2
dtype: int64
```

```
[72] 1 # Checking the number of Diabetic and Non-Diabetic patients
2 print("Number of Diabetic and Non-Diabetic patients:")
3 print(dataset['Outcome'].value_counts())
```

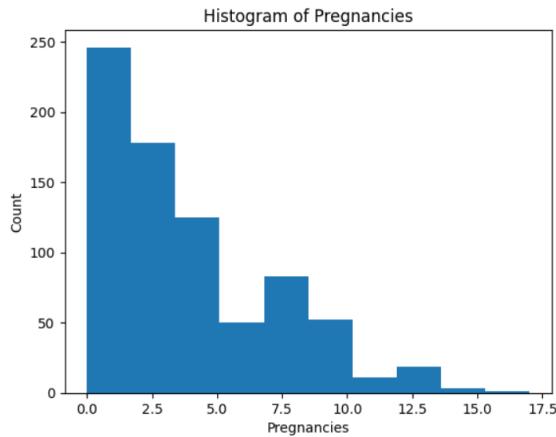
Python

```
[72] ... Number of Diabetic and Non-Diabetic patients:
Non-Diabetic     500
Diabetic        268
Name: Outcome, dtype: int64
```

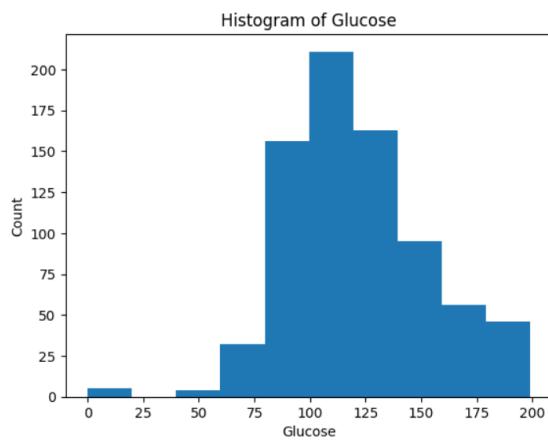
```
1 # Plotting Histograms for all the features in the dataset to check the distribution of the data
2 for i in dataset.columns:
3     plt.hist(dataset[i])
4     plt.xlabel(i)
5     plt.ylabel('Count')
6     plt.title("Histogram of " + i)
7     plt.show()
8
```

Python

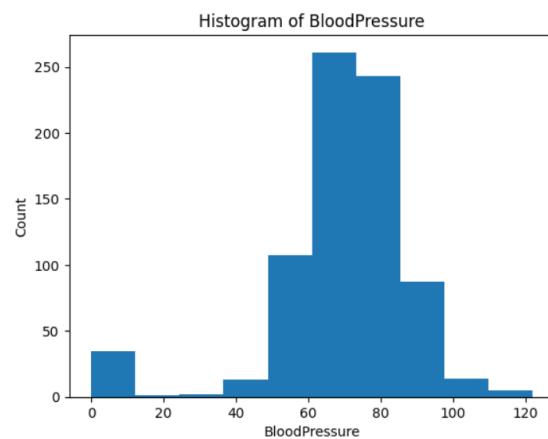
[78]



</>

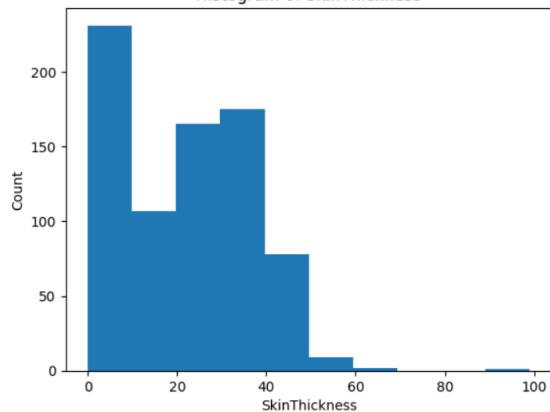


</>



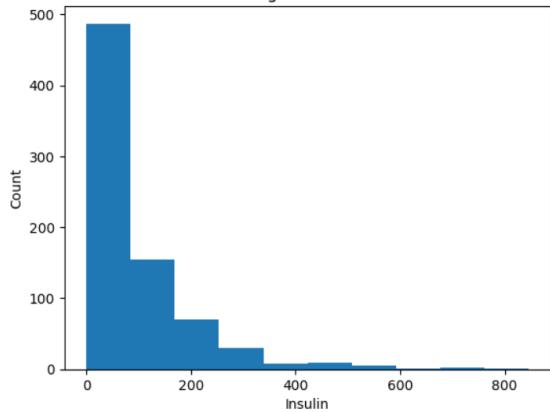
</>

Histogram of SkinThickness



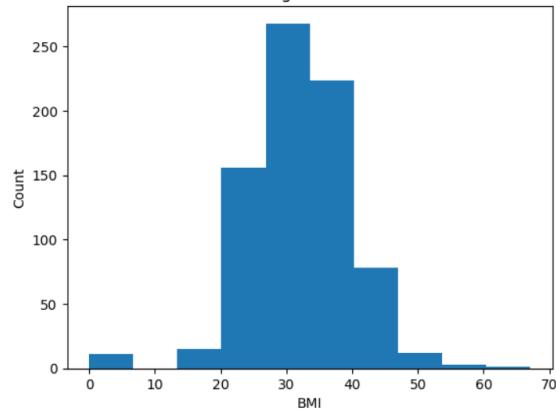
</>

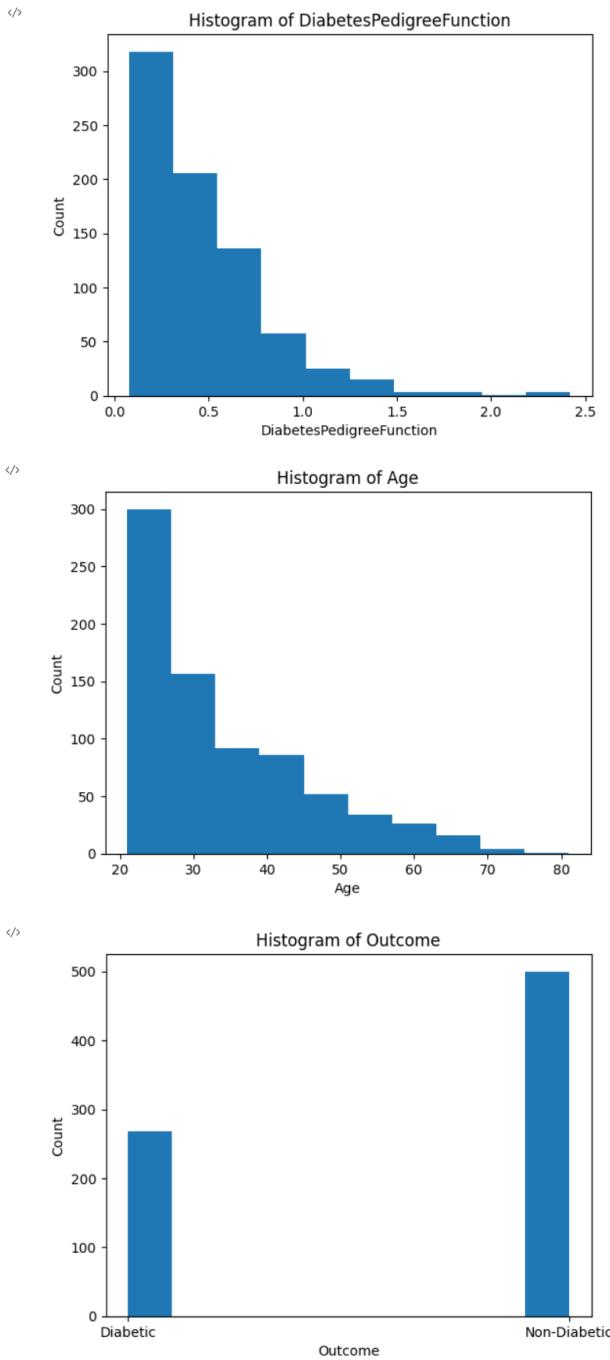
Histogram of Insulin



</>

Histogram of BMI





```

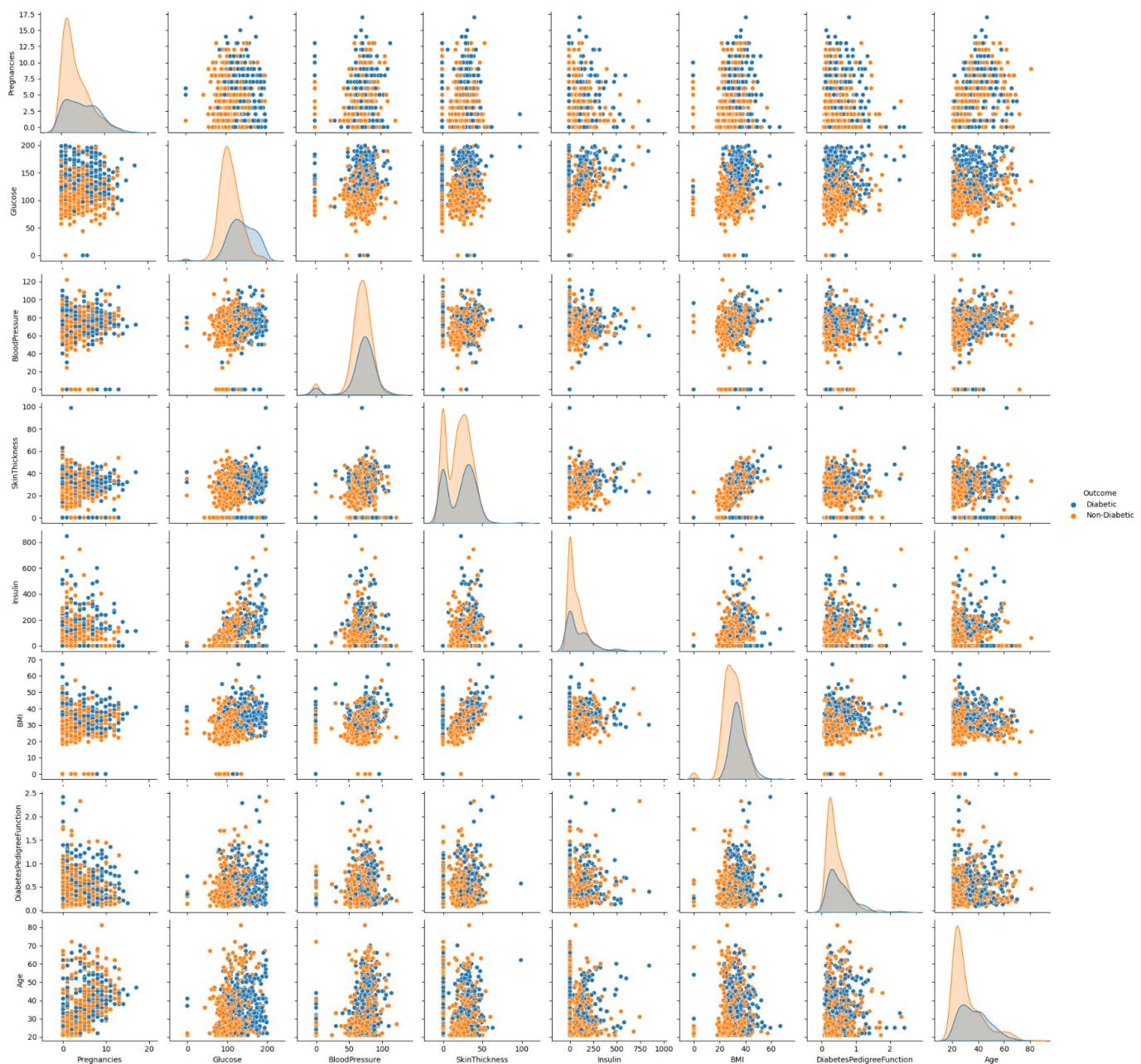
1 # Pairplot of the dataset
2 sns.pairplot(dataset, hue='Outcome')
3 plt.suptitle("Pairplot of all variables colored by Outcome", y=1.05, fontsize=16)
4 plt.show()

```

[92]

Python

Pairplot of all variables colored by Outcome

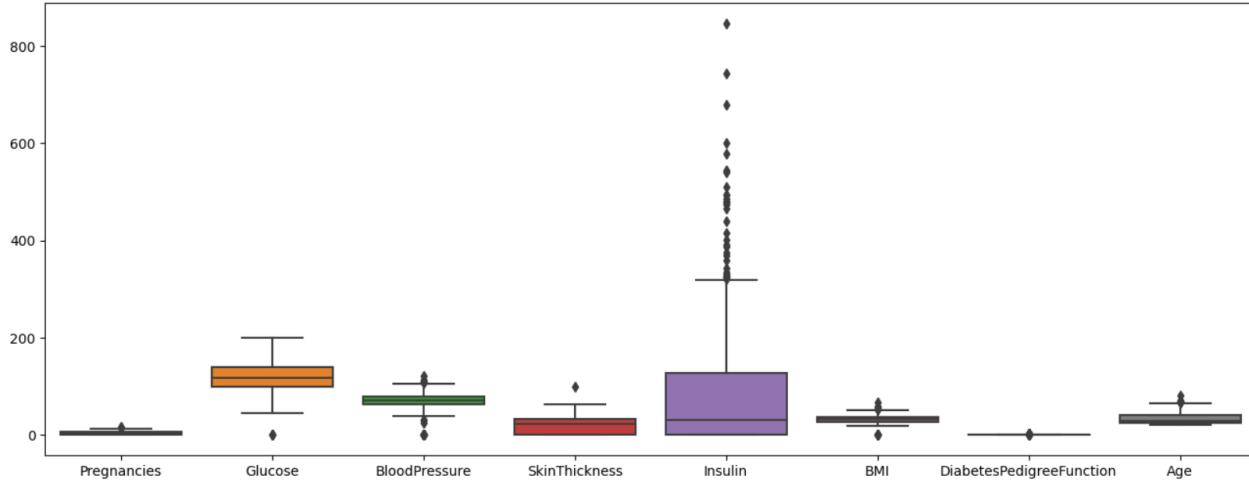


```
1 # Boxplot of the dataset
2 fig, ax = plt.subplots(figsize=(16,6))
3 sns.boxplot(data=dataset)
4 plt.title("Boxplot of all variables", fontsize=16)
5 plt.show()
```

[19] ✓ 0.3s

Python Python

Boxplot of all variables



## EXPERIMENT NO. 2

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 03.02.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> Find mean, median, mode, standard deviation for a given dataset
<b>Outcome:</b> Measure of central tendency, measure of variation and measure of position
<b>Problem Statement:</b> Demonstration of descriptive statistical measures
<b>Background Study:</b> Descriptive statistics are brief descriptive coefficients that summarize a given data set, which can be either a representation of the entire or a sample of a population. Descriptive statistics are broken down into measures of central tendency and measures of variability (spread). Measures of central tendency include the mean, median and mode, while measures of variability include standard deviation, variance, minimum and maximum variables.
<b>Question Bank:</b> 1. What is difference between mean mode medians? <i>The mean is the average of a set of data, the mode is the most common value in the set, and the median is the middle value of the set.</i>  2. What is advantage of median over mean? <i>The advantages of using the median over the mean are that the median is less sensitive to outliers and extreme values, and it can be used for non-normally distributed data.</i>  3. Is mean affected by presence of outlier in data. <i>Yes, the mean is affected by outliers in the data as it can significantly change its value.</i>  4. How to find the deviation of data in given dataset? <i>To find the deviation of data in a given dataset, you can follow these steps:</i> <i>1) Calculate the mean of the dataset.</i> <i>2) For each data point, subtract the mean from the data point.</i> <i>3) Square each of the differences obtained in step 2.</i> <i>4) Calculate the average of the squared differences obtained in step 3.</i> <i>5) Take the square root of the value obtained in step 4.</i> <i>This final value is the standard deviation of the dataset, which represents how much the data points deviate from the mean.</i>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 2

Demonstration of descriptive statistical measures.

```

1 # importing required libraries
2 import math
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 import scipy
9 from scipy import stats
10 from scipy.stats import zscore, kurtosis, variation, scoreatpercentile
11 from scipy.stats.mstats import gmean
[2] ✓ 2.6s                                         Python

1 # Read the data
2 data = pd.read_csv('iris_dataset.csv')
[3] ✓ 0.0s                                         Python

1 # Drop the Id column
2 data = data.drop('Id', axis=1)
[4] ✓ 0.1s                                         Python

1 # Print the first 5 rows of the data
2 print(data.head().to_markdown())
[5] ✓ 0.0s                                         Python
...
|   | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|-----:|-----:|-----:|-----:|:-----|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 4.9 | 3 | 1.4 | 0.2 | Iris-setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5 | 3.6 | 1.4 | 0.2 | Iris-setosa |

1 # Information about the data
2 print(data.info())
[6] ✓ 0.1s                                         Python
...
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   SepalLengthCm  150 non-null   float64
 1   SepalWidthCm   150 non-null   float64
 2   PetalLengthCm  150 non-null   float64
 3   PetalWidthCm   150 non-null   float64
 4   Species        150 non-null   object 
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
None

```

```

1 # Describe the data
2 print(data.groupby('Species').describe())
[7]   ✓ 0.1s
...
SepalLengthCm
Species      count    mean     std   min    25%  50%  75%  max
Iris-setosa    50.0  5.006  0.352490  4.3  4.800  5.0  5.2  5.8 \
Iris-versicolor 50.0  5.936  0.516171  4.9  5.600  5.9  6.3  7.0
Iris-virginica 50.0  6.588  0.635880  4.9  6.225  6.5  6.9  7.9

SepalWidthCm
Species      count    mean     ... PetalLengthCm      PetalWidthCm
Species      count    mean     ... 75% max      count
Species      ...
Iris-setosa    50.0  3.418  ...      1.575  1.9      50.0 \
Iris-versicolor 50.0  2.770  ...      4.600  5.1      50.0
Iris-virginica 50.0  2.974  ...      5.875  6.9      50.0

mean      std   min    25%  50%  75%  max
Species
Iris-setosa    0.244  0.107210  0.1  0.2  0.2  0.3  0.6
Iris-versicolor 1.326  0.197753  1.0  1.2  1.3  1.5  1.8
Iris-virginica  2.026  0.274650  1.4  1.8  2.0  2.3  2.5

[3 rows x 32 columns]

```

Mean, Median, Mode

```

1 # Mean of the data
2 print(data.groupby('Species').mean().to_markdown())
[8]   ✓ 0.0s
...
| Species | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|:-----:|:-----:|:-----:|:-----:|:-----:|
| Iris-setosa | 5.006 | 3.418 | 1.464 | 0.244 |
| Iris-versicolor | 5.936 | 2.770 | 4.600 | 1.326 |
| Iris-virginica | 6.588 | 2.974 | 5.552 | 2.026 |

1 # Median of the data
2 print(data.groupby('Species').median().to_markdown())
[9]   ✓ 0.0s
...
| Species | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|:-----:|:-----:|:-----:|:-----:|:-----:|
| Iris-setosa | 5 | 3.4 | 1.5 | 0.2 |
| Iris-versicolor | 5.9 | 2.8 | 4.35 | 1.3 |
| Iris-virginica | 6.5 | 3 | 5.55 | 2 |

1 # Mode of the data
2 data.groupby('Species').value_counts()
[10]  ✓ 0.0s
...
  Species SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm
Iris-setosa  4.9       3.1       1.5       0.1       3
             4.3       3.0       1.1       0.1       1
             4.4       2.9       1.4       0.2       1
             5.0       3.4       1.5       0.2       1
                           1.6       0.4       1
                           ..
Iris-virginica 6.5       3.0       5.8       2.2       1
                 3.2       5.1       2.0       1
                 6.7       2.5       5.8       1.8       1
                 7.9       3.0       5.2       2.3       1
                 7.9       3.8       6.4       2.0       1
Name: count, Length: 147, dtype: int64

```

#### Geometric Mean and Harmonic Mean

```

1 # geometric mean of the data
2 print(data.groupby('Species').apply(gmean).to_markdown())
3
4 # Harmonic mean of the data
5 print(data.groupby('Species').apply(scipy.stats.hmean).to_markdown())
[11] ✓ 0.0s
... c:\Users\mainp\AppData\Local\Programs\Python\Python311\Lib\site-packages\tabulate\_init_.py:107: FutureWarning: elementwise comparison failed; returning scalar instead, but in the fut
 (len(row) >= 1 and row[0] == SEPARATING_LINE)
c:\Users\mainp\AppData\Local\Programs\Python\Python311\Lib\site-packages\tabulate\_init_.py:108: FutureWarning: elementwise comparison failed; returning scalar instead, but in the fut
 or (len(row) >= 2 and row[1] == SEPARATING_LINE)
c:\Users\mainp\AppData\Local\Programs\Python\Python311\Lib\site-packages\tabulate\_init_.py:107: FutureWarning: elementwise comparison failed; returning scalar instead, but in the fut
 (len(row) >= 1 and row[0] == SEPARATING_LINE)
c:\Users\mainp\AppData\Local\Programs\Python\Python311\Lib\site-packages\tabulate\_init_.py:108: FutureWarning: elementwise comparison failed; returning scalar instead, but in the fut
 or (len(row) >= 2 and row[1] == SEPARATING_LINE)
| Species | 0
|-----|---|
| Iris-setosa | [4.99384106 3.3969062 1.45373856 0.22346247]
| Iris-versicolor | [5.91397936 2.75187353 4.23308089 1.31118738]
| Iris-virginica | [6.55579452 2.95701356 5.52578887 2.00721413]
| Species | 0
|-----|---|
| Iris-setosa | [4.9816814 3.37529787 1.44318022 0.20449898]
| Iris-versicolor | [5.89193597 2.73301053 4.20449935 1.29602855]
| Iris-virginica | [6.52735823 2.94008724 5.50023704 1.98786729]

```

#### Standard Deviation

```

1 # Standard Deviation of the data
2 print(data.groupby('Species').std().to_markdown())
[12] ✓ 0.0s
... | Species | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-----|---|---|---|---|
| Iris-setosa | 0.35249 | 0.381024 | 0.173511 | 0.10721 |
| Iris-versicolor | 0.516171 | 0.313798 | 0.469911 | 0.197753 |
| Iris-virginica | 0.63588 | 0.322497 | 0.551895 | 0.27465 |

```

```

1 # Standard Deviation of the data using numpy
2 print(data.groupby('Species').agg(np.std).to_markdown())
[13] ✓ 0.0s
... | Species | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-----|---|---|---|---|
| Iris-setosa | 0.35249 | 0.381024 | 0.173511 | 0.10721 |
| Iris-versicolor | 0.516171 | 0.313798 | 0.469911 | 0.197753 |
| Iris-virginica | 0.63588 | 0.322497 | 0.551895 | 0.27465 |

```

```

1 # Standard Deviation Using Formula
2 variance = data.groupby('Species').var()
3 standard_deviation = np.sqrt(variance)
4 print(standard_deviation.to_markdown())
[14] ✓ 0.0s
... | Species | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-----|---|---|---|---|
| Iris-setosa | 0.35249 | 0.381024 | 0.173511 | 0.10721 |
| Iris-versicolor | 0.516171 | 0.313798 | 0.469911 | 0.197753 |
| Iris-virginica | 0.63588 | 0.322497 | 0.551895 | 0.27465 |

```

Variance

```

1 # Variance of the data
2 print(data.groupby('Species').var().to_markdown())
[15] ✓ 0.0s
... | Species      | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-----|-----|-----|-----|-----|
| Iris-setosa | 0.124249 | 0.14518 | 0.0301061 | 0.0114939 |
| Iris-versicolor | 0.266433 | 0.0984694 | 0.220816 | 0.0391061 |
| Iris-virginica | 0.404343 | 0.104004 | 0.304588 | 0.0754327 |

1 # Variance of the data using numpy
2 print(data.groupby('Species').agg(np.var).to_markdown())
[16] ✓ 0.0s
... | Species      | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-----|-----|-----|-----|-----|
| Iris-setosa | 0.124249 | 0.14518 | 0.0301061 | 0.0114939 |
| Iris-versicolor | 0.266433 | 0.0984694 | 0.220816 | 0.0391061 |
| Iris-virginica | 0.404343 | 0.104004 | 0.304588 | 0.0754327 |

1 # Variance of the data using formula
2 standard_deviation = data.groupby('Species').std()
3 variance = np.square(standard_deviation)
4 print(variance.to_markdown())
[17] ✓ 0.0s
... | Species      | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-----|-----|-----|-----|-----|
| Iris-setosa | 0.124249 | 0.14518 | 0.0301061 | 0.0114939 |
| Iris-versicolor | 0.266433 | 0.0984694 | 0.220816 | 0.0391061 |
| Iris-virginica | 0.404343 | 0.104004 | 0.304588 | 0.0754327 |

```

## EXPERIMENT NO. 3

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 04.02.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> Find skewness and kurtosis for a given dataset
<b>Outcome:</b> Understand skewness and kurtosis and their implementation
<b>Problem Statement:</b> Demonstrate measures of shape using Python
<b>Background Study:</b> Descriptive statistics are brief descriptive coefficients that summarize a given data set, which can be either a representation of the entire or a sample of a population. Descriptive statistics are broken down into measures of central tendency and measures of variability (spread). Measures of central tendency include the mean, median and mode, while measures of variability include standard deviation, variance, minimum and maximum variables, and kurtosis and skewness.
<b>Question Bank:</b> <ol style="list-style-type: none"><li>1. Can we identify outlier from stem-and-leaf plot? Yes, we can identify outliers from a stem-and-leaf plot, although it may not be as straightforward as identifying them from other types of plots like box plots or histograms. In a stem-and-leaf plot, outliers are typically defined as values that are much smaller or much larger than the other values in the data set. This can be assessed by looking for any "leaves" that are significantly separated from the rest of the data points.</li><li>2. What is kernel density estimation? Kernel Density Estimation (KDE) is a non-parametric technique used for estimating the probability density function (PDF) of a random variable. It is particularly useful when the underlying distribution of the data is unknown, or when a smooth estimate of the density is needed.</li><li>3. How do we decide upon the number of bins in histogram. The number of bins in a histogram can have a significant impact on the way the data is presented and interpreted, so it is important to choose an appropriate number of bins.</li><li>4. Which plot is used for analysis of bivariate data? Scatter plot is the most common plot used for analyzing the relationship between two continuous variables. Other plots include line plots, bubble plots, contour plots, and heatmaps.</li></ol>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 3

Demonstrate measures of shape using Python.

```
1 # importing required libraries
2 import math
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 import scipy
9 from scipy import stats
10 from scipy.stats import zscore, kurtosis, variation, scoreatpercentile
11 from scipy.stats.mstats import gmean
12
```

[1] ✓ 1.6s

Python

```
1 # Read the data
2 data = pd.read_csv('iris_dataset.csv')
3
```

[2] ✓ 0.0s

Python

```
1 # Drop the Id column
2 data = data.drop('Id', axis=1)
3
```

[3] ✓ 0.0s

Python

```
1 # Print the first 5 rows of the data
2 print(data.head().to_markdown())
3
```

[4] ✓ 0.0s

Python

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5	3.6	1.4	0.2	Iris-setosa

```
1 # Information about the data
2 print(data.info())
3
```

[5] ✓ 0.0s

Python

```
... <class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   SepalLengthCm  150 non-null   float64 
 1   SepalWidthCm   150 non-null   float64 
 2   PetalLengthCm  150 non-null   float64 
 3   PetalWidthCm   150 non-null   float64 
 4   Species       150 non-null   object  
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
None
```

```

1 # Describe the data
2 print(data.groupby('Species').describe())
3
[6] ✓ 0.1s
...
SepalLengthCm
   count    mean      std    min    25%   50%   75%   max
Species
Iris-setosa      50.0  5.006  0.352490  4.3  4.800  5.0  5.2  5.8 \
Iris-versicolor  50.0  5.936  0.516171  4.9  5.600  5.9  6.3  7.0
Iris-virginica   50.0  6.588  0.635880  4.9  6.225  6.5  6.9  7.9

SepalWidthCm
   count    mean      ... PetalLengthCm      PetalWidthCm
   ...          ...    75%   max           count
Species
Iris-setosa      50.0  3.418  ...     1.575  1.9      50.0 \
Iris-versicolor  50.0  2.770  ...     4.600  5.1      50.0
Iris-virginica   50.0  2.974  ...     5.875  6.9      50.0

mean      std    min    25%   50%   75%   max
Species
Iris-setosa     0.244  0.107210  0.1  0.2  0.2  0.3  0.6
Iris-versicolor 1.326  0.197753  1.0  1.2  1.3  1.5  1.8
Iris-virginica  2.026  0.274650  1.4  1.8  2.0  2.3  2.5

[3 rows x 32 columns]

```

Z-Score

```

1 # Z-Score of the data
2 print(data.groupby('Species').apply(zscore).to_markdown())
3
[7] ✓ 0.0s
...
Output exceeds the size\_limit. Open the full output data in a text editor
|-----|-----|-----|-----|-----|
| ('Iris-setosa', 0) | 0.269382 | 0.217394 | -0.372597 | -0.414578 |
| ('Iris-setosa', 1) | -0.303771 | -1.10818 | -0.372597 | -0.414578 |
| ('Iris-setosa', 2) | -0.876924 | -0.577951 | -0.95478 | -0.414578 |
| ('Iris-setosa', 3) | -1.1635 | -0.843065 | 0.209586 | -0.414578 |
| ('Iris-setosa', 4) | -0.0171946 | 0.482509 | -0.372597 | -0.414578 |
| ('Iris-setosa', 5) | 1.12911 | 1.27785 | 1.37395 | 1.46987 |
| ('Iris-setosa', 6) | -1.1635 | -0.0477207 | -0.372597 | 0.527645 |
| ('Iris-setosa', 7) | -0.0171946 | -0.0477207 | 0.209586 | -0.414578 |
| ('Iris-setosa', 8) | -1.73665 | -1.3733 | -0.372597 | -0.414578 |
| ('Iris-setosa', 9) | -0.303771 | -0.843065 | 0.209586 | -1.3568 |
| ('Iris-setosa', 10) | 1.12911 | 0.747624 | 0.209586 | -0.414578 |
| ('Iris-setosa', 11) | -0.590348 | -0.0477207 | 0.791769 | -0.414578 |
| ('Iris-setosa', 12) | -0.590348 | -1.10818 | -0.372597 | -1.3568 |
| ('Iris-setosa', 13) | -2.02323 | -1.10818 | -2.11915 | -1.3568 |
| ('Iris-setosa', 14) | 2.27542 | 1.54297 | -1.53696 | -0.414578 |
| ('Iris-setosa', 15) | 1.98884 | 2.60343 | 0.209586 | 1.46987 |
| ('Iris-setosa', 16) | 1.12911 | 1.27785 | -0.95478 | 1.46987 |
| ('Iris-setosa', 17) | 0.269382 | 0.217394 | -0.372597 | 0.527645 |
| ('Iris-setosa', 18) | 1.98884 | 1.01274 | 1.37395 | 0.527645 |
| ('Iris-setosa', 19) | 0.269382 | 1.01274 | 0.209586 | 0.527645 |
| ('Iris-setosa', 20) | 1.12911 | -0.0477207 | 1.37395 | -0.414578 |
| ('Iris-setosa', 21) | 0.269382 | 0.747624 | 0.209586 | 1.46987 |
| ('Iris-setosa', 22) | -1.1635 | 0.482509 | -2.70133 | -0.414578 |
...
| ('Iris-virginica', 146) | -0.457514 | -1.4847 | -1.01035 | -0.463423 |
| ('Iris-virginica', 147) | -0.139796 | 0.0814395 | -0.644278 | -0.095627 |
| ('Iris-virginica', 148) | -0.616373 | 1.33435 | -0.278211 | 1.00776 |
| ('Iris-virginica', 149) | -1.09295 | 0.0814395 | -0.827312 | -0.831219 |

```

Skewness

```
1 # Skewness
2 print(data.groupby('Species').skew().to_markdown())
3
[8] ✓ 0.0s
```

Species	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
Iris-setosa	0.120087	0.107053	0.0718461	1.19724
Iris-versicolor	0.105378	-0.362845	-0.606508	-0.0311796
Iris-virginica	0.118015	0.365949	0.549445	-0.129477

Python

Kurtosis

```
1 # Kurtosis
2 data_agg = data.groupby('Species').agg({'SepalLengthCm': kurtosis, 'SepalWidthCm': kurtosis,
3                                         'PetalLengthCm': kurtosis, 'PetalWidthCm': kurtosis}).reset_index()
4 print(data_agg.to_markdown())
5
[9] ✓ 0.0s
```

	Species	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
0	Iris-setosa	-0.345765	0.685134	0.813665	1.29648
1	Iris-versicolor	-0.598827	-0.448272	-0.0744018	-0.487833
2	Iris-virginica	-0.0879422	0.519766	-0.256472	-0.661348

Python

## EXPERIMENT NO. 4

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 10.02.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> Data visualization for qualitative and quantitative data
<b>Outcome:</b> Students will be familiarized with the use of python libraries such as matplotlib and seaborn
<b>Problem Statement:</b> Write a python program to use visualize qualitative and quantitative data
<b>Background Study:</b> Pareto Chart, Pie chart, bar charts and boxplot are some examples of visualization for qualitative data stem-and-leaf chart, count plot, histogram, frequency polygon and scatter plot are some examples of visualization for qualitative data
<b>Question Bank:</b> 1. What kind of data is generated using boxplot. A box plot is a graphical representation of the distribution of numerical data through their quartiles, allowing the observation of the spread and skewness of the data, as well as the presence of outliers. Box plots are particularly useful for comparing data between different groups or datasets.  2. Can we plot boxplot in both matplotlib and seaborn library? Yes, both Matplotlib and Seaborn libraries support the creation of box plots. In Matplotlib, you can create a box plot using the boxplot() function from the pyplot module. In Seaborn, you can create a box plot using the boxplot() function from the seaborn module.

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 4 (Qualitative)

Data visualization for qualitative data.

```

1 # importing required libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 %matplotlib inline
[1] ✓ 2.7s

```

```

1 # Read the data
2 housing = pd.read_csv('housing_dataset.csv')
[2] ✓ 0.1s

```

```

1 # Printing the First 5 Rows of the Data
2 print(housing.head().to_markdown())
3
[3] ✓ 0.0s

```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280	565	259	3.8462	342200	NEAR BAY

```

1 # Describing the Data
2 print(housing.describe().to_markdown())
[4] ✓ 0.1s

```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value
count	20640	20640	20640	20640	20433	20640	20640	20640	20640
mean	-119.57	35.6319	28.6395	2635.76	537.871	1425.48	499.54	3.87067	206856
std	2.00353	2.13595	12.5856	2181.62	421.385	1132.46	382.33	1.89982	115396
min	-124.35	32.54	1	2	1	3	1	0.4999	14999
25%	-121.8	33.93	18	1447.75	296	787	280	2.5634	119600
50%	-118.49	34.26	29	2127	435	1166	409	3.5348	179700
75%	-118.01	37.71	37	3148	647	1725	605	4.74325	264725
max	-114.31	41.95	52	39320	6445	35682	6082	15.0001	500001

```

1 # Information about the Data
2 print(housing.info())
[5] ✓ 0.0s

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
 #   Column      Non-Null Count  Dtype  
 --- 
 0   longitude    20640 non-null   float64
 1   latitude     20640 non-null   float64
 2   housing_median_age  20640 non-null   float64
 3   total_rooms   20640 non-null   float64
 4   total_bedrooms 20433 non-null   float64
 5   population    20640 non-null   float64
 6   households    20640 non-null   float64
 7   median_income 20640 non-null   float64
 8   median_house_value 20640 non-null   float64
 9   ocean_proximity 20640 non-null   object 
dtypes: float64(9), object(1)
memory usage: 1.6+ MB
None

```

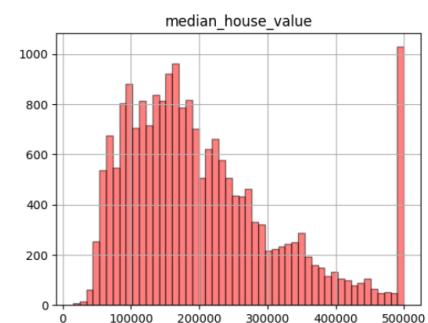
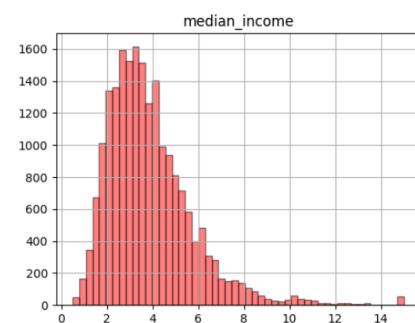
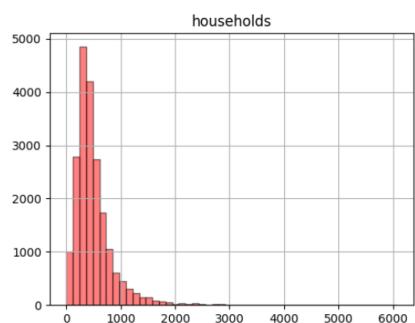
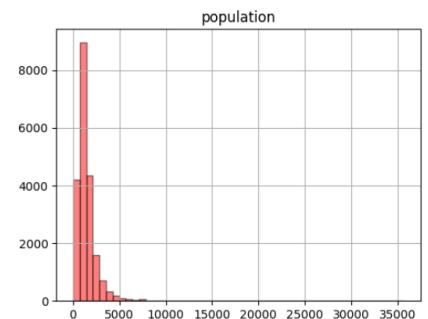
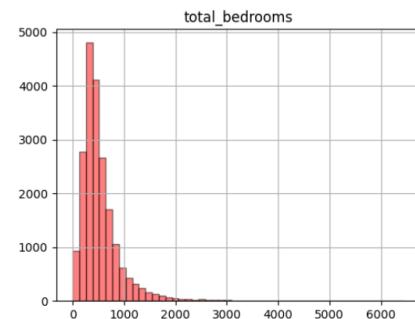
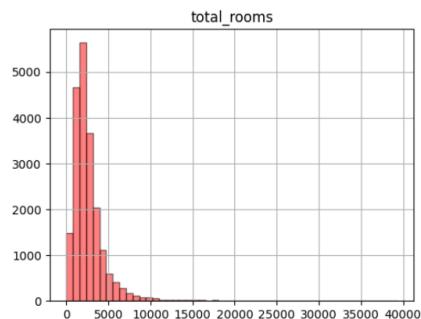
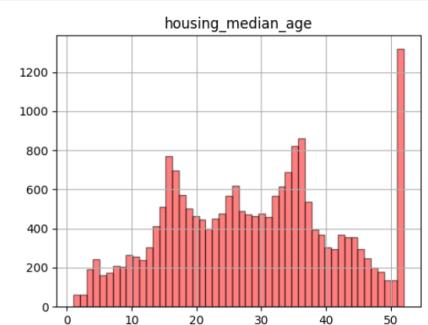
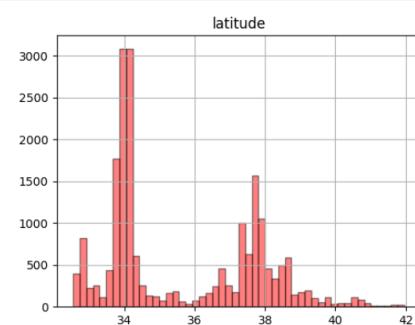
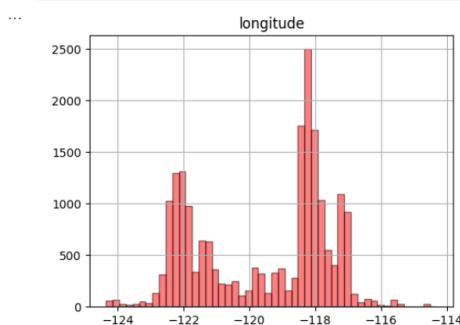
```
[6] 1 # Checking for Missing Values
2 print(housing.isnull().sum())
✓ 0.1s
```

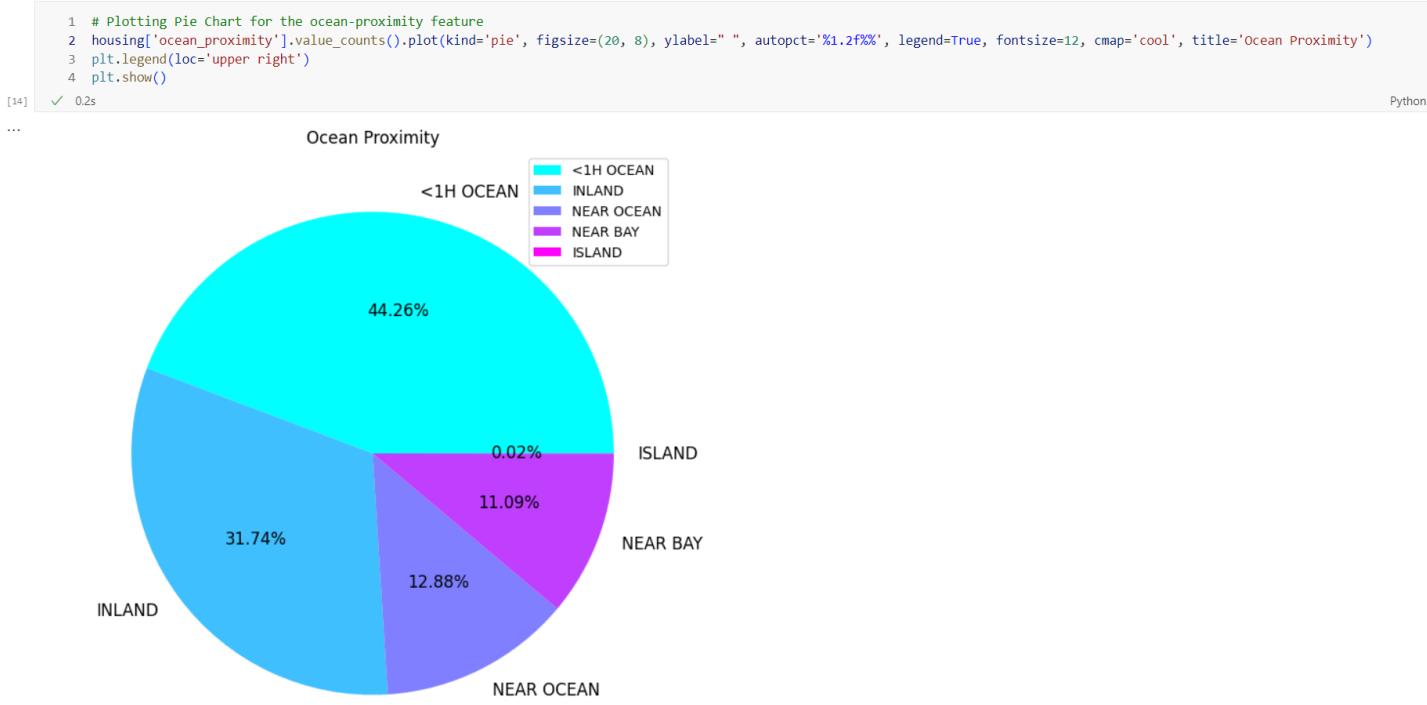
Python

```
[...]
longitude      0
latitude       0
housing_median_age   0
total_rooms     0
total_bedrooms  207
population      0
households      0
median_income    0
median_house_value 0
ocean_proximity 0
dtype: int64
```

```
[7] 1 # Plotting Histograms for each Numerical Feature
2 housing.hist(bins=50, figsize=(20, 15), color='red',
3 |           alpha=0.5, edgecolor='black')
4 plt.show()
✓ 2.4s
```

Python







## Experiment 4 (Quantitative)

Data visualization for quantitative data.

```
1 # importing required libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import pylab
7
8 %matplotlib inline
```

[1] ✓ 1.9s

Python

```
1 # Loading the dataset
2 data = pd.read_csv('heart_dataset.csv')
```

[2] ✓ 0.0s

Python

```
1 # Prining the first 5 rows of the dataset
2 print(data.head().to_markdown())
```

[3] ✓ 0.0s

Python

	age	sex	cp	trestbps	chol	fbfs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	52	1	0	125	212	0	1	168	0	1	2	2	3	0
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3	0
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3	0
3	61	1	0	148	203	0	1	161	0	0	2	1	3	0
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2	0

```
1 # Information about the dataset
2 print(data.info())
3
```

[4] ✓ 0.0s

Python

```
... <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1025 entries, 0 to 1024
Data columns (total 14 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   age      1025 non-null   int64  
 1   sex      1025 non-null   int64  
 2   cp       1025 non-null   int64  
 3   trestbps 1025 non-null   int64  
 4   chol     1025 non-null   int64  
 5   fbs      1025 non-null   int64  
 6   restecg  1025 non-null   int64  
 7   thalach  1025 non-null   int64  
 8   exang    1025 non-null   int64  
 9   oldpeak  1025 non-null   float64 
 10  slope    1025 non-null   int64  
 11  ca       1025 non-null   int64  
 12  thal    1025 non-null   int64  
 13  target   1025 non-null   int64  
dtypes: float64(1), int64(13)
memory usage: 112.2 KB
None
```

```
1 # Check for missing values
2 print(data.isnull().sum())
[5] ✓ 0.0s
```

Python

```
... age      0
sex      0
cp      0
trestbps 0
chol      0
fbs      0
restecg  0
thalach   0
exang      0
oldpeak   0
slope      0
ca      0
thal      0
target     0
dtype: int64
```

```
1 # Dropping the missing values
2 data.dropna(inplace=True)
[6] ✓ 0.0s
```

Python

```
1 # Information about the dataset
2 print(data.info())
[7] ✓ 0.0s
```

Python

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1025 entries, 0 to 1024
Data columns (total 14 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   age      1025 non-null   int64  
 1   sex      1025 non-null   int64  
 2   cp       1025 non-null   int64  
 3   trestbps 1025 non-null   int64  
 4   chol      1025 non-null   int64  
 5   fbs       1025 non-null   int64  
 6   restecg   1025 non-null   int64  
 7   thalach   1025 non-null   int64  
 8   exang     1025 non-null   int64  
 9   oldpeak   1025 non-null   float64 
 10  slope     1025 non-null   int64  
 11  ca        1025 non-null   int64  
 12  thal      1025 non-null   int64  
 13  target    1025 non-null   int64  
dtypes: float64(1), int64(13)
memory usage: 112.2 KB
None
```

```
1 # Checking the first 5 rows of the dataset
2 print(data.head().to_markdown())
3
```

Python

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	52	1	0	125	212	0	1	168	0	1	2	2	3	0
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3	0
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3	0
3	61	1	0	148	203	0	1	161	0	0	2	1	3	0
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2	0

```
1 # Extract the age (relevant column) for the stem and leaf plot
2 age_data = data['age'].head(10).tolist()
3
4 # Convert the values in the data list to integers
5 age_list = [int(value) for value in age_data]
[9] ✓ 0.0s
```

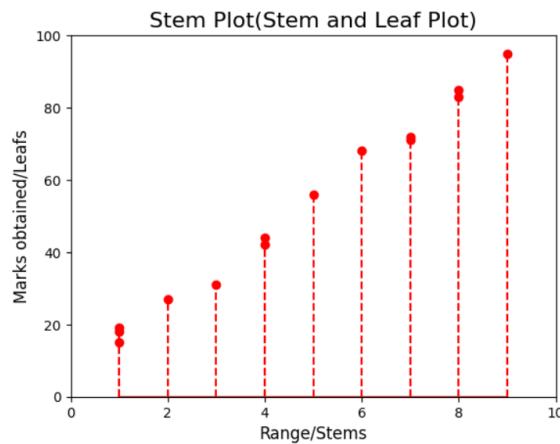
Python

## Leaf and Stem Plot

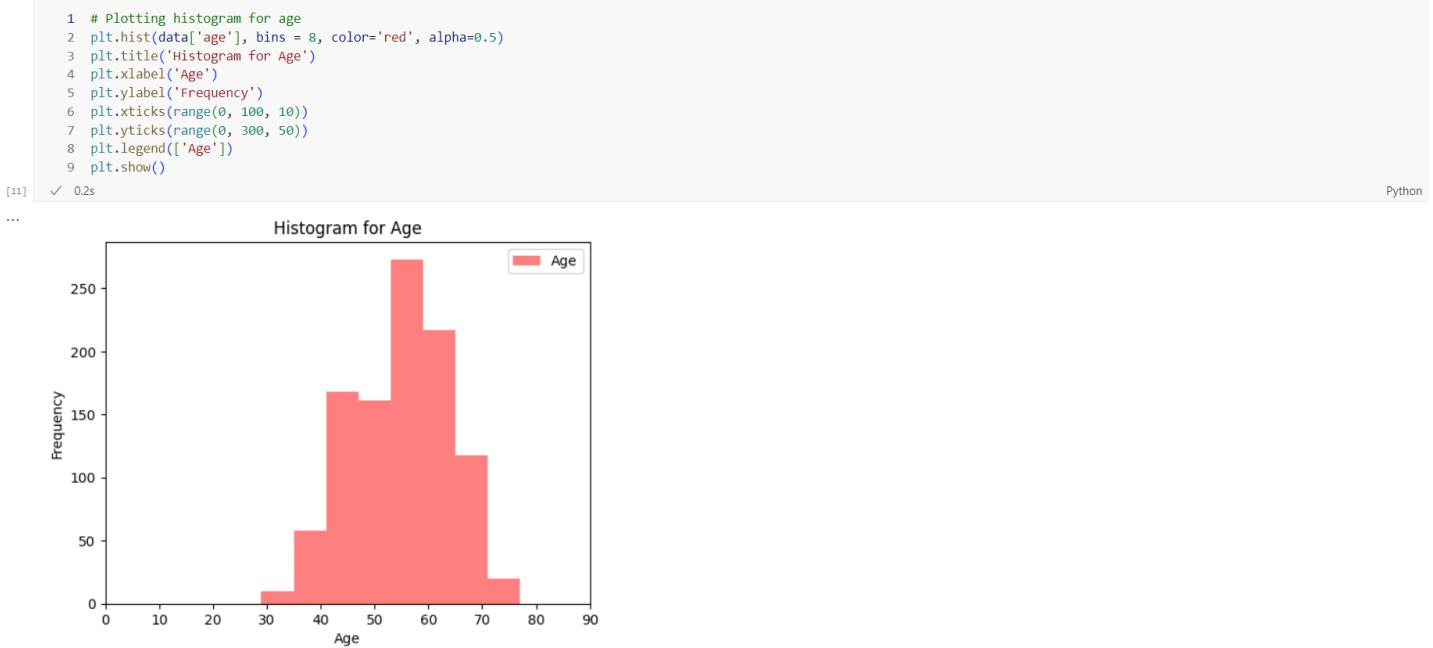
```
1 # Create the stem and leaf plot
2
3 # marks obtained by students in an examination
4 marks = [15, 18, 19, 27, 31, 42, 44, 56, 68, 71, 72, 83, 85, 95]
5 # corresponding stems
6 stems = [1, 1, 1, 2, 3, 4, 4, 5, 6, 7, 7, 8, 8, 9]
7
8 # set the x axis and y axis limits
9 pylab.xlim([0, 10])
10 pylab.ylim([0, 100])
11
12 # Provide a title for the stem plot
13 plt.title('Stem Plot(Stem and Leaf Plot)', fontsize=16)
14
15 # Give x and y label for the stem plot
16 plt.xlabel('Range/Stems', fontsize=12)
17 plt.ylabel('Marks obtained/Leafs', fontsize=12)
18
19 # plot the stem plot using matplotlib
20 markerline, stemlines, baseline = plt.stem(
21 |     stems, marks, '--')
22
23 plt.setp(markerline, color='red')
24 plt.setp(stemlines, color='red')
25
26 # display the stem plot
27 plt.show()
28
```

[10] ✓ 0.2s

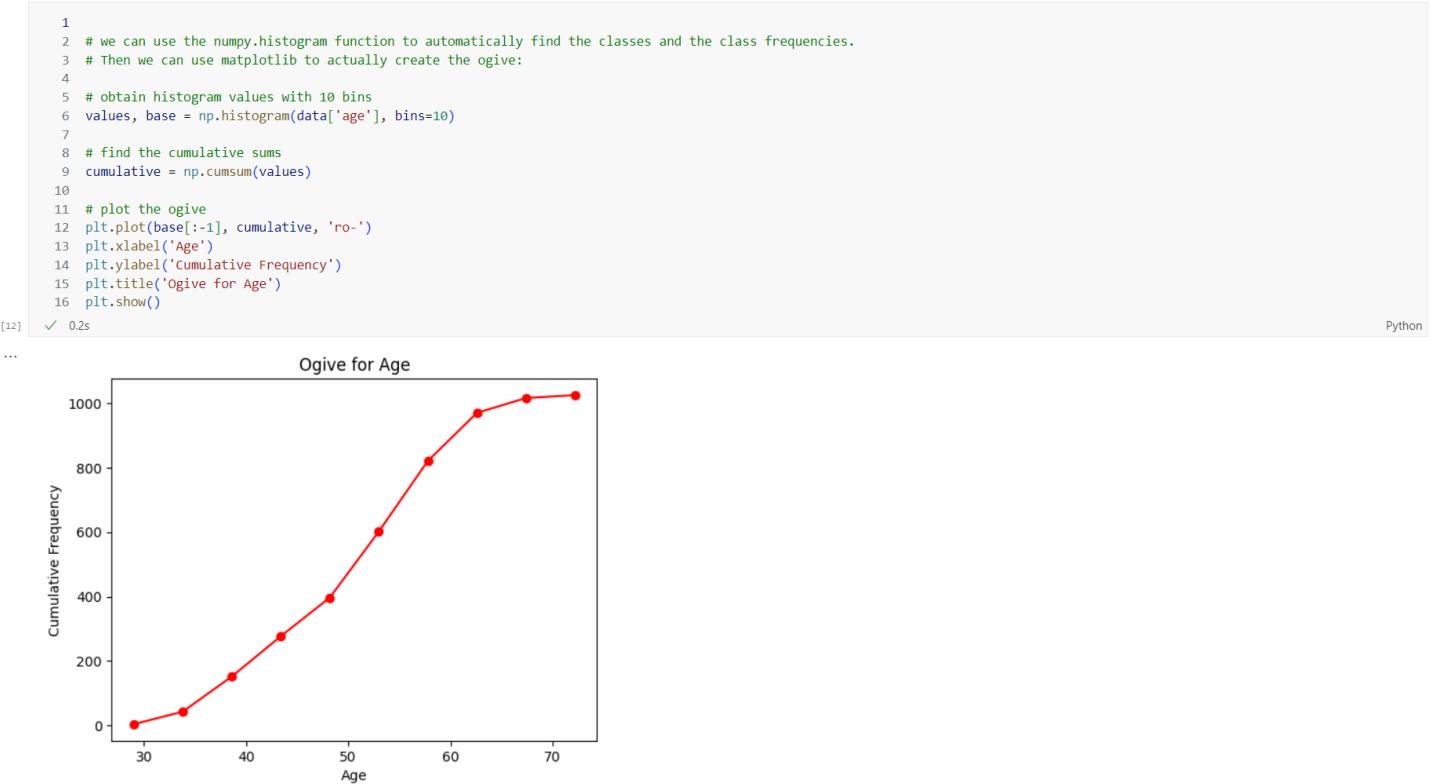
Python



### Histogram: For Univariate Variable



### Ogive /Cumulative Frequency Polygon Chart for Univariate Quantitative data

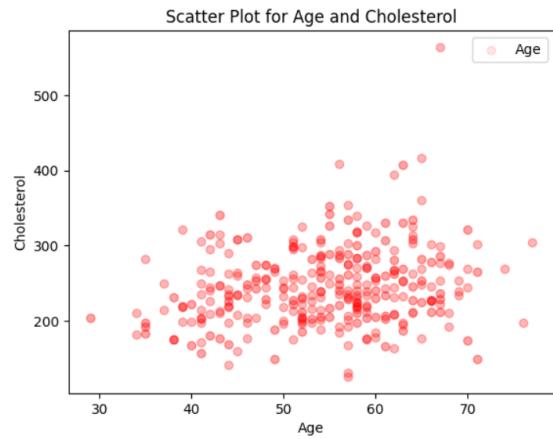


## Scatter Plot

```
1 # Plotting Scatter Plot
2 plt.scatter(data['age'], data['chol'], color='red', alpha = 0.1)
3 plt.title('Scatter Plot for Age and Cholesterol')
4 plt.xlabel('Age')
5 plt.ylabel('Cholesterol')
6 plt.legend(['Age', 'Weight'])
7 plt.show()
```

[13] ✓ 0.3s

Python



## EXPERIMENT NO. 5

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 24.02.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> Write python program to implement Probability
<b>Problem Statement:</b> Basics of probability using Python
<b>Background Study:</b> Various types of events, conditional probability, permutation and combination
<b>Question Bank:</b>
<ol style="list-style-type: none"><li>1. What are mutually exclusive events? Mutually exclusive events are events that cannot occur at the same time. In other words, if one event happens, the other event cannot happen. For example, if you flip a coin, the result can be either heads or tails, but not both. Therefore, "getting heads" and "getting tails" are mutually exclusive events. Similarly, if you roll a dice, getting a 1 and getting a 6 are also mutually exclusive events. In probability theory, the probability of the union of mutually exclusive events is equal to the sum of the probabilities of each individual event.</li><li>2. What is the importance of studying probability? Relate with real-life situations and explain. It helps us make informed decisions based on the likelihood of different outcomes. Probability is used in a wide range of real-life situations, from weather forecasting to sports betting to business decision-making. By understanding the probability of different events occurring, we can make more accurate predictions, assess risks more effectively, and make better-informed decisions. Probability is a fundamental concept in mathematics and statistics, and its applications extend to many different fields. Whether we are trying to predict the weather, assess risks in business or medicine, or make smart bets in sports, probability helps us make sense of the world around us and make better decisions based on the available data.</li></ol>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

DRILL-9

Five Minute  
Mystery



#### The Case of the Two Classes

The Head First Health Club prides itself on its ability to find a class for everyone. As a result, it is extremely popular with both young and old.

The Health Club is wondering how best to market its new yoga class, and the Head of Marketing wonders if someone who goes swimming is more likely to go to a yoga class. "Maybe we could offer some sort of discount to the swimmers to get them to try out yoga."

The CEO disagrees. "I think you're wrong," he says. "I think that people who go swimming and people who go to yoga are independent. I don't think people who go swimming are any more likely to do yoga than anyone else."

They ask a group of 96 people whether they go to the swimming or yoga classes. Out of these 96 people, 32 go to yoga and 72 go swimming. 24 people are exceptionally eager and go to both.

**So who's right? Are the yoga and swimming classes dependent or independent?**

$$P(\text{Yoga}) = \frac{32}{96} = \frac{1}{3}$$

$$P(\text{Swim}) = \frac{72}{96} = \frac{3}{4}$$

$$P(Y \cap S) = \frac{24}{96} = \frac{1}{4}$$

$$P(Y) \times P(S) = \frac{1}{3} \times \frac{3}{4} = \frac{1}{4}$$

$$= P(Y \cap S)$$

Manic Mango selects one of the volunteers at random to ask if she enjoyed playing the game, and she says she did. Given that the volunteer enjoyed playing the game, what's the probability that she played game 2? Use Bayes' Theorem.

Hint: What's the probability of someone choosing game 2 and being satisfied?  
What's the probability of someone being satisfied overall? Once you've found these, you can use Bayes Theorem to obtain the right answer.

$$P(\text{Game2} / \text{Enjoyed})$$

$$\rightarrow P(\text{G}_2 / E) = \frac{P(\text{G}_2) \cdot P(E | \text{G}_2)}{P(\text{G}_2) \cdot P(E | \text{G}_2) + P(\text{G}_1) \cdot P(E | \text{G}_1)}$$

$$= \frac{0.2 \times 0.7}{(0.2 \times 0.7) + (0.8 \times 0.6)}$$

$$= \frac{0.14}{0.14 + 0.48}$$

$$= \frac{0.14}{0.62} = 0.226$$

*long exercise*



### LONG Exercise

### DRILL-7, 8

The Manic Mango games company is testing two brand-new games. They've asked a group of volunteers to choose the game they most want to play, and then tell them how satisfied they were with game play afterwards.

80 percent of the volunteers chose Game 1, and 20 percent chose Game 2. Out of the Game 1 players, 60 percent enjoyed the game and 40 percent didn't. For Game 2, 70 percent of the players enjoyed the game and 30 percent didn't.

Your first task is to fill in the probability for this scenario.

$$P(\text{Game 1}) = .80$$

$$P(\text{Game 2}) = .20$$

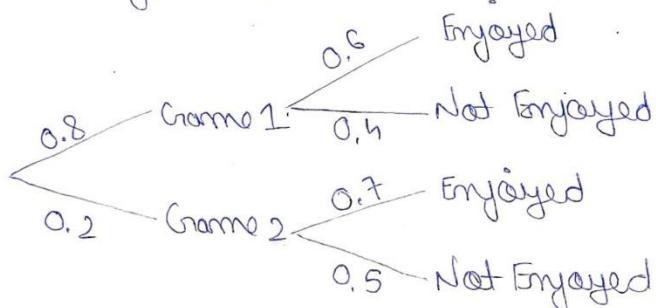
Game 1

$$P(\text{Enjoy}) = .60$$

$$P(\text{Enjoy}) = .70$$

$$P(\text{Didn't Enjoy}) = .40 \quad P(\text{Didn't Enjoy}) = .30$$

Tree Diagram:



*probability magnets*



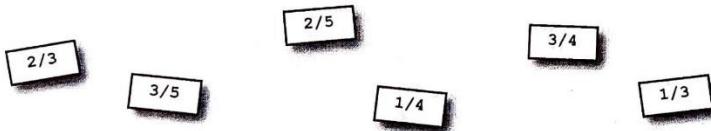
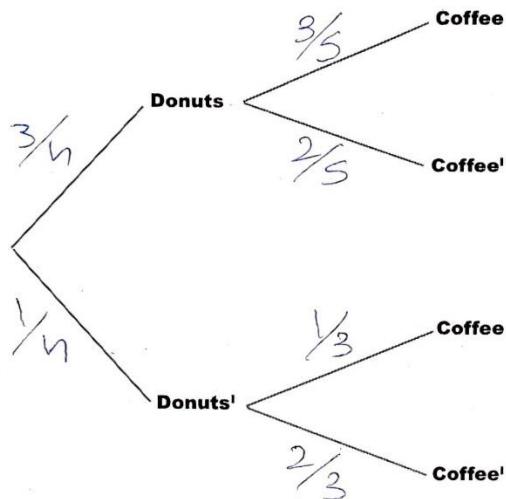
## Probability Magnets

Duncan's Donuts are looking into the probabilities of their customers buying donuts *and* coffee. They drew up a probability tree to show the probabilities, but in a sudden gust of wind, they all fell off. Your task is to pin the probabilities back on the tree. Here are some clues to help you.

$$P(\text{Donuts}) = 3/4$$

$$P(\text{Coffee} | \text{Donuts}) = 1/3$$

$$P(\text{Donuts} \cap \text{Coffee}) = 9/20$$



### DRILL-6



Exercise

We haven't quite finished with Duncan's Donuts! Now that you've completed the probability tree, you need to use it to work out some probabilities.

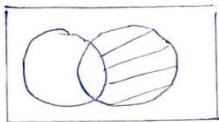
$$P(D) = \frac{3}{n} \quad P(C \cap D) = \frac{1}{3} \quad P(D \cap C) = \frac{9}{20}$$

1.  $P(\text{Donuts})$

$$\begin{aligned} P(D') &= 1 - P(D) \\ &= 1 - \frac{3}{n} \\ &= \frac{n-3}{n} \end{aligned}$$

2.  $P(\text{Donuts} \cap \text{Coffee})$

$$\frac{1}{3} \times \frac{1}{n} = \frac{1}{12}$$



3.  $P(\text{Coffee} \mid \text{Donuts})$

$$P(C' \mid D) = \frac{2}{5}$$

4.  $P(\text{Coffee})$  ← Hint: How many ways are there of getting coffee? (You can get coffee with or without donuts.)

$$\begin{aligned} P(C) &= P(D' \cap C) + \\ &\quad P(D \cap C) \\ &= \frac{1}{12} + \frac{9}{20} = \frac{8}{15} \end{aligned}$$

5.  $P(\text{Donuts} \mid \text{Coffee})$

$$P(D \mid C) = \frac{9/20}{8/15} = \frac{27}{32}$$

Hint: maybe some of your other answers can help you.

$$n(\text{Basketball}) = 18 + 6 = 28$$

$$n(\text{Football}) = 12 + n = 16$$

$$n(\text{Baseball}) = 10 + 6 = 16$$

$$P(\text{Baseball} \cap \text{Football}) = 0$$

∴ Basketball and Football are mutually exclusive.

$$P(\text{Baseball} \cup \text{Football} \cup \text{Basketball}) = 1$$

So, the events for baseball, football and basketball are exhaustive.



### DRILL-1

50 sports enthusiasts at the Head First Health Club are asked whether they play baseball, football, or basketball. 10 only play baseball. 12 only play football. 18 only play basketball. 6 play baseball and basketball but not football. 4 play football and basketball but not baseball.

Draw a Venn diagram for this probability space. How many enthusiasts play baseball in total? How many play basketball? How many play football?

Are any sports' rosters mutually exclusive? Which sports are exhaustive (fill up the possibility space)?

$$\text{Total} = 50$$

$$n P(B) = 10$$

$$n P(F) = 12$$

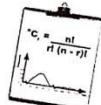
$$n P(K) = 18$$

$$n (B \cap K) = 6$$

$$n (F \cap K) = 4$$

B - Baseball

K - Basketball



### Vital Statistics

A or B

To find the probability of getting event A or B, use

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

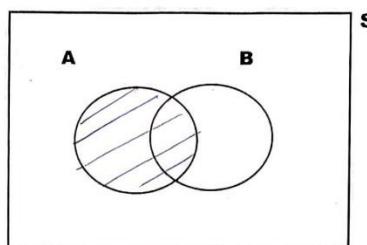
$\cup$  means OR

$\cap$  means AND

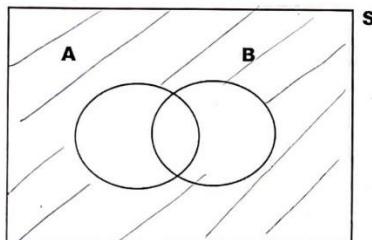
### DRILL-3



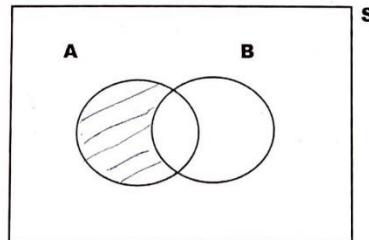
**BE the probability**  
Your job is to play like you're the probability and shade in the area that represents each of the following probabilities on the Venn diagrams.



$$P(A \cap B) + P(A \cap B')$$



$$P(A' \cap B')$$



$$P(A \cup B) - P(B)$$

 Sharpen your pencil DRTLL-1,2

### **calculating probabilities**

Let's find the probability of getting a black or even (assume 0 and 00 are not even).

1. What's the probability of getting a black?

$$P(\text{Getting Black}) = \frac{18}{38} = 0.474$$

2. What's the probability of getting an even number?

$$P(\text{Even Number}) = \frac{18}{38} = 0.474$$

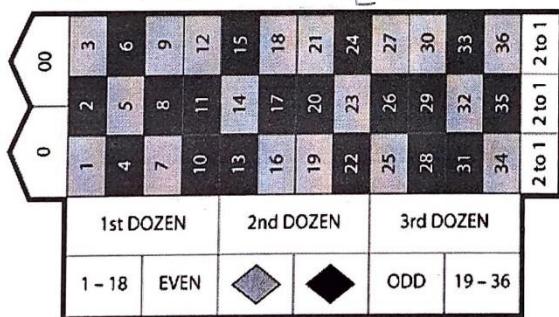
3. What do you get if you add these two probabilities together?

$$\text{Adding Probabilities} = 0.947 \left[ \frac{18}{38} + \frac{18}{38} \right]$$

4. Finally, use your roulette board to count all the holes that are either black or even, then divide by the total number of holes. What do you get?

$$P\left(\frac{\text{Black or Even}}{\text{Total Number of Holes}}\right) = \frac{26}{38} = 0.684$$

$$\left[ \frac{18}{38} + \frac{18}{36} - \frac{10}{38} = \frac{26}{38} \right]$$



we can say that the classes are independent.

## EXPERIMENT NO. 6

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 03.03.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> Understand and apply the laws of probability
<b>Problem Statement:</b> To implement laws of probability
<b>Background Study:</b> Laws of probability-General and special laws of addition and multiplication, concepts of conditional probability and Bayes theorem.
<b>Question Bank:</b>
1. What is the difference between independent and dependent events? <i>In probability theory, independent events are events that do not affect each other, meaning the occurrence of one event does not change the probability of the other event happening. In contrast, dependent events are events where the occurrence of one event affects the probability of the other event happening. In other words, the outcome of one event is dependent on the outcome of the other event.</i>
2. What is the relationship between conditional probability ad bayes theorem? <i>Bayes' theorem is a way to calculate conditional probabilities. It states that the probability of an event A given that another event B has occurred is equal to the probability of event B given event A multiplied by the prior probability of event A, divided by the prior probability of event B. In mathematical notation, this can be expressed as:</i> <i><math>P(A B) = P(B A) * P(A) / P(B)</math></i>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

```
"""
Experiment 6
Write Python Program to Implement Laws of Probability
"""

# Addition Law of Probability
def addition_law(a, b):
    return a + b - (a * b)

# Multiplication Law of Probability
def multiplication_law(a, b):
    return a * b

# Conditional Probability
def conditional_probability(a, b):
    return (a * b) / b

# Total Probability Theorem
def total_probability_theorem(probability):

    event = int(input("Enter the Event Number: "))
    probability_event = 0
    for i in range(len(probability)):
        probability_event_given_i = conditional_probability(probability[event-1],
probability[i])
        probability_event += probability[i] * probability_event_given_i

    print("\nThe Probability if the {} is : {}".format(event, probability_event))

# Complement of an Event
def complement(a):
    return 1 - a

def main():
    # Menu for user to choose which law of probability to use
    print("1. Addition Law of Probability")
    print("2. Multiplication Law of Probability")
    print("3. Conditional Probability")
```

```
print("4. Complement of an Event")
print("5. Total Probability Theorem")
print("6. Exit")
choice = int(input("\nEnter Your Choice: "))

# Addition Law of Probability
if choice == 1:
    a = float(input("Enter the probability of event A: "))
    b = float(input("Enter the probability of event B: "))

    if(a > 1 or b > 1):
        print("Probability cannot be greater than 1")
        exit()
    print("\n")
    print("P(A or B) = ", addition_law(a, b))

# Multiplication Law of Probability
elif choice == 2:
    a = float(input("Enter the probability of event A: "))
    b = float(input("Enter the probability of event B: "))
    if(a > 1 or b > 1):
        print("Probability cannot be greater than 1")
        exit()
    print("\n")
    print("P(A and B) = ", multiplication_law(a, b))

# Conditional Probability
elif choice == 3:
    a = float(input("Enter the probability of event A: "))
    b = float(input("Enter the probability of event B: "))
    if(a > 1 or b > 1):
        print("Probability cannot be greater than 1")
        exit()
    print("\n")
    print("P(A | B) = ", conditional_probability(a, b))

# Complement of an Event
elif choice == 4:
    a = float(input("Enter the probability of event A: "))
    print("P(not A) = ", complement(a))

# Total Probability Theorem
elif choice == 5:
    probability = []
    number_of_events = int(input("Enter the number of events: "))
```

```

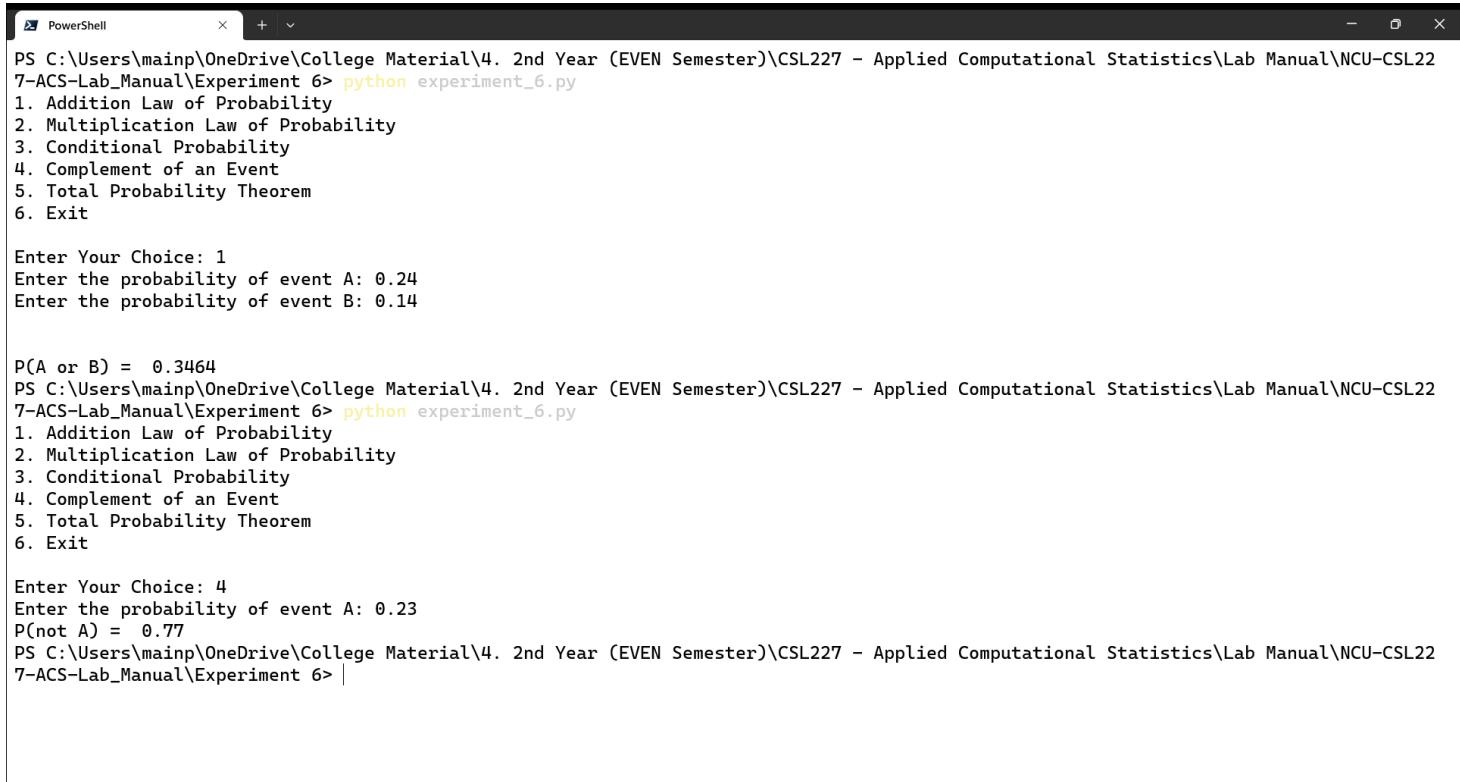
for i in range(number_of_events):
    inp = float(input("Enter the probability of event {}: ".format(i + 1)))
    if(inp > 1):
        print("Probability cannot be greater than 1")
        exit()
    probability.append(inp)

# Exit
elif choice == 6:
    exit()

# Invalid Choice
else:
    print("Invalid Choice")

main()

```



```

PS C:\Users\mainp\OneDrive\College Material\4. 2nd Year (EVEN Semester)\CSL227 - Applied Computational Statistics\Lab Manual\NCU-CSL227-ACS-Lab_Manual\Experiment 6> python experiment_6.py
1. Addition Law of Probability
2. Multiplication Law of Probability
3. Conditional Probability
4. Complement of an Event
5. Total Probability Theorem
6. Exit

Enter Your Choice: 1
Enter the probability of event A: 0.24
Enter the probability of event B: 0.14

P(A or B) = 0.3464
PS C:\Users\mainp\OneDrive\College Material\4. 2nd Year (EVEN Semester)\CSL227 - Applied Computational Statistics\Lab Manual\NCU-CSL227-ACS-Lab_Manual\Experiment 6> python experiment_6.py
1. Addition Law of Probability
2. Multiplication Law of Probability
3. Conditional Probability
4. Complement of an Event
5. Total Probability Theorem
6. Exit

Enter Your Choice: 4
Enter the probability of event A: 0.23
P(not A) = 0.77
PS C:\Users\mainp\OneDrive\College Material\4. 2nd Year (EVEN Semester)\CSL227 - Applied Computational Statistics\Lab Manual\NCU-CSL227-ACS-Lab_Manual\Experiment 6>

```

## EXPERIMENT NO. 7

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab Manual (github.com)</a>
<b>Date:</b> 17.03.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> Write a python program to compute PMF in discrete distribution
<b>Outcome:</b> Explore the functions related to implementation of discrete distributions
<b>Problem Statement:</b> Implementation of probability mass function for discrete random variable using Python.
<b>Background Study:</b> The Probability Mass Function (PMF) is also called a <b>probability function</b> or <b>frequency function</b> which characterizes the distribution of a discrete random variable. Let X be a discrete random variable of a function, then the probability mass function of a random variable X is given by
$P_x(x) = P(X=x)$ , For all x belongs to the range of X
It is noted that the probability function should fall on the condition:
<ul style="list-style-type: none"> <li>• <math>P_x(x) \geq 0</math> and</li> <li>• <math>\sum_{x \in \text{Range}(X)} P_x(x) = 1</math></li> </ul>
The discrete distributions under consideration are geometric distribution, binomial distribution, poisson distribution, hypergeometric distribution
<b>Question Bank:</b>
1. What is meant by random variable? Explain with example.  In probability theory, a random variable is a variable that takes on different values based on the outcomes of a random event. It can be thought of as a numerical representation of the outcomes of a random process.  For example, suppose we flip a fair coin two times and count the number of heads that appear. The random variable in this case is the number of heads, which can take on values of 0, 1, or 2.
2. How to calculate the probability for the random variable?  To calculate the probability for a random variable, we first need to define the probability distribution of the random variable, which specifies the probability of the random variable taking on each possible value.
3. Explain the probability mass function of various discrete distributions?  The probability mass function (PMF) is a function that specifies the probability distribution of a discrete random variable. It gives the probability that the random variable takes on each possible value.

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 7

Implementation of probability mass function for discrete random variable using Python.

```

1 # importing required libraries
2 from scipy.stats import bernoulli
3 from scipy.stats import binom
4 from scipy.stats import geom
5 from scipy.stats import hypergeom
6 from scipy.stats import poisson

```

[1] ✓ 1.4s

Python

#### 1. Bernoulli Distribution

```

1 # defining parameters
2 n = 10 # number of trials
3 p = 0.5 # probability of success
4
5 print('Mean' , bernoulli.mean(p))
6 print('Variance' , bernoulli.var(p))

```

[2] ✓ 0.0s

Python

... Mean 0.5  
Variance 0.25

#### 2. Binomial Distribution

```

1 # defining parameters
2 n = 10 # number of trials
3 x = 7 # number of successes
4 p = 0.2 # probability of success
5
6 print("Mean: ", binom.mean(n, p))
7 print("Variance: ", binom.var(n, p))
8 print("Probability Mass Function: ", binom.pmf(x, n, p))
9 print("Cumulative Distribution Function: ", binom.cdf(x,n,p))

```

[3] ✓ 0.1s

Python

... Mean: 2.0  
Variance: 1.6  
Probability Mass Function: 0.0007864320000000000  
Cumulative Distribution Function: 0.9999220736

#### 3. Geometric Distribution

```

1 # defining parameters
2 n = 10 # number of trials
3 p = 0.2 # probability of success
4
5 print("Mean: ", geom.mean(p))
6 print("Variance: ", geom.var(p))
7 print("Probability Mass Function: ", geom.pmf(3, p))
8 print("Cumulative Distribution Function: ", geom.cdf(3, p))

```

[4] ✓ 0.0s

Python

... Mean: 5.0  
Variance: 20.0  
Probability Mass Function: 0.12800000000000003  
Cumulative Distribution Function: 0.488

#### 4. Hypergeometric Distribution

```
1 # defining parameters
2 n = 10 # number of trials
3 x = 7 # number of successes
4 M = 20 # number of balls in urn
5 N = 10 # number of balls drawn
6
7 print("Mean: ", hypergeom.mean(M, n, N))
8 print("Variance: ", hypergeom.var(M, n, N))
9 print("Probability Mass Function: ", hypergeom.pmf(x, M, n, N))
10 print("Cumulative Distribution Function: ", hypergeom.cdf(x, M, n, N))

[5] ✓ 0.1s
```

Python

```
... Mean: 5.0
Variance: 1.3157894736842106
Probability Mass Function: 0.07794063521617701
Cumulative Distribution Function: 0.9884929312173895
```

#### 5. Poisson Distribution

```
1 # defining parameters
2 x = 5 # number of events
3 Lambda = 2/5 # lambda parameter
4
5 print("Mean: ", poisson.mean(Lambda))
6 print("Variance: ", poisson.var(Lambda))
7 print("Probability mass function: ", poisson.pmf(x, Lambda))
8 print("Cumulative distribution function: ", poisson.cdf(x, Lambda))

[6] ✓ 0.0s
```

Python

```
... Mean: 0.4
Variance: 0.4
Probability mass function: 5.720064392837453e-05
Cumulative distribution function: 0.999995957316826
```

## EXPERIMENT NO. 8

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 17.03.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> To find expectation and variance in a discrete distribution
<b>Outcome:</b> Knows how to compute mean and variance in a given discrete distribution
<b>Problem Statement:</b> Write a python program to find expectation and variance in discrete distribution
<b>Background Study:</b> Mean and variance are measure of centre and variation in data, respectively. We need to compute them for a given distribution, whether discrete or continuous
<b>Question Bank:</b>
<b>1.</b> How do we compute the mean of a discrete distribution? To compute the mean (or expected value) of a discrete distribution, we can use the formula: $E(X) = \sum x * P(X=x)$ , where $x$ is a possible outcome and $P(X=x)$ is its corresponding probability. This means that we multiply each possible outcome by its probability, and then sum up these products to get the expected value.
<b>2.</b> How do we compute the variance in a discrete distribution? To compute the variance of a discrete distribution, we can use the formula: $Var(X) = \sum (x - E(X))^2 * P(X=x)$ , where $x$ is a possible outcome, $E(X)$ is its expected value, and $P(X=x)$ is its corresponding probability. This means that we subtract each possible outcome from its expected value, square this difference, multiply it by its probability, and then sum up these products to get the variance.

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 8

Write a python program to find expectation and variance in discrete distribution.

```
1 # importing required libraries
2 from scipy.stats import rv_discrete
[1] ✓ 1.25                                         Python

1 # define the distribution
2 x = [1, 2, 3, 4]
3 p = [0.2, 0.3, 0.1, 0.4]
4
5 # create a discrete random variable
6 random_variable = rv_discrete(values=(x, p))
7
8 # calculate the expectation
9 expectation = random_variable.expect()
10
11 # calculate the variance
12 variance = random_variable.var()
13
14 print("Expectation: ", expectation)
15 print("Variance: ", variance)
[2] ✓ 0.0s                                         Python

... Expectation: 2.7
Variance: 1.4099999999999984
```

## EXPERIMENT NO. 9

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 24.03.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> To understand various types of continuous probability distributions
<b>Outcome:</b> Students will be familiarized with normal and uniform distribution.
<b>Problem Statement:</b> Write Python Program for Continuous Probability Distributions
<b>Background Study:</b> Probability density function and cumulative distribution function for the normal and uniform distribution is computed
<b>Question Bank:</b>
<p>1. What is the mean and variance of the uniform distribution?  <b>The mean and variance of the uniform distribution are:</b>        • Mean: <math>(a + b) / 2</math>, where <math>a</math> and <math>b</math> are the lower and upper bounds of the distribution.        • Variance: <math>(b - a)^2 / 12</math>        Note that the uniform distribution is a continuous probability distribution that assigns equal probability to all values within a given interval <math>(a, b)</math>.</p>
<p>2. What is the mean and the variance of the normal distribution?  <b>The mean and variance of the normal distribution are:</b>        • Mean: <math>\mu</math>, where <math>\mu</math> is the mean of the population        • Variance: <math>\sigma^2</math>, where <math>\sigma</math> is the standard deviation of the population        Note that the normal distribution is a continuous probability distribution that is symmetric and bell-shaped, and it is characterized by two parameters, the mean and the standard deviation. The mean determines the center of the distribution, while the standard deviation determines the spread of the distribution.</p>
<p>3. How do we compute probability density function and cumulative distribution function for the normal distribution?        The probability density function (PDF) of the normal distribution is <math>f(x) = (1 / \sigma\sqrt{2\pi}) * e^{-(x-\mu)^2/(2\sigma^2)}</math>, and the cumulative distribution function (CDF) can be computed by standardizing the variable and using the CDF of the standard normal distribution.</p>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 9

Write Python Program for Continuous Probability Distributions.

##### 1. Uniform Distribution

```
1 # importing required libraries
2 from scipy.stats import uniform
3
4 # setting the values of a and b
5 a = 32
6 b = 37
7
8 # printing the mean, variance, PDF, CDF
9 print("Mean: ", uniform.mean(a, b))
10 print("Variance: ", uniform.var(a, b))
11 print("PDF: ", uniform.pdf(a, b))
12 print("CDF: ", uniform.cdf(a, b))
13
14
[1] ✓ 0.7s
```

Mean: 50.5  
Variance: 114.08333333333333  
PDF: 0.0  
CDF: 0.0

Python

##### 2. Normal Distribution

```
1 # importing required libraries
2 from scipy.stats import norm as normal
3
4 # setting the values of mean and standard deviation
5 x = 1.6
6 mean = 8
7 sd = 5
8
9 # printing the mean, variance, PDF, CDF
10 print("Mean: ", normal.mean(x))
11 print("Variance: ", normal.var(x))
12 print("PDF: ", normal.pdf(x))
13 print("CDF: ", normal.cdf(x))
14
15
[2] ✓ 0.0s
```

Mean: 1.6  
Variance: 1.0  
PDF: 0.11092083467945554  
CDF: 0.945200708300442

Python

### 3. Chi-Square Distribution

```

1 # importing required libraries
2 from scipy.stats import chi2 as chi
3
4 # setting the values of degrees of freedom, variable, critical value
5 nu = 6 # degrees of freedom
6 x = 3 # variable
7 alpha = 0.24 # critical value
8
9 # printing the mean, variance, PDF, CDF, survival function, percent point function, critical value
10 print("Mean: ", chi.mean(nu))
11 print("Variance: ", chi.var(nu))
12 print("PDF: ", chi.pdf(x, nu))
13 print("CDF: ", chi.cdf(x, nu))
14 print("Survival Function: ", chi.sf(x, nu))
15 print("Percent Point Function: ", chi.ppf(alpha, nu))
16 print("Critical Value: ", chi.isf(alpha, nu))

[3] ✓ 0.0s
...
Mean: 6.0
Variance: 12.0
PDF: 0.12551071508349182
CDF: 0.1911531694619418
Survival Function: 0.8088468305380582
Percent Point Function: 3.3789420922671023
Critical Value: 7.974194328808714

```

Python

### 4. The t-Distribution

```

1 # importing required libraries
2 from scipy.stats import t
3
4 # setting the values of degrees of freedom, variable, critical value
5 nu = 68 # degrees of freedom
6 x = 1 # variable
7 alpha = 0.925 # critical value
8
9 # printing the mean, variance, PDF, CDF, survival function, percent point function, critical value
10 print("Expectation: ", t.mean(nu))
11 print("Variance: ", t.var(nu))
12 print("Probability Distribution Function: ", t.pdf(x, nu))
13 print("Cumulative Distribution Function: ", t.cdf(x, nu))
14 print("Survival Function (1-cdf): ", t.sf(x, nu))
15 print("Percent Point Function (inverse of cdf): ", t.ppf(alpha, nu))
16 print("Inverse Survival Function (inverse of (1-cdf)): ", t.isf(alpha, nu)) # useful for finding the critical value

[4] ✓ 0.0s
...
Expectation: 0.0
Variance: 1.0303030303030303
Probability Distribution Function: 0.2402023823138007
Cumulative Distribution Function: 0.8395721066953445
Survival Function (1-cdf): 0.1604278933046555
Percent Point Function (inverse of cdf): 1.455979454119329
Inverse Survival Function (inverse of (1-cdf)): -1.455979454119329

```

Python

## 5. The F-Distribution

```
1 # importing required libraries
2 from scipy.stats import f
3
4
5 # setting the values of degrees of freedom, variable, critical value
6 nu1 = 24 # degree of freedom 1
7 nu2 = 16 # degree of freedom 2
8 x = 2 # variable
9 alpha = 0.925 # critical value
10
11 # printing the mean, variance, PDF, CDF, survival function, percent point function, critical value
12 print("Expectation: ", f.mean(nu1, nu2))
13 print("Variance: ", f.var(nu1, nu2))
14 print("Probability Distribution Function: ", f.pdf(x ,nu1, nu2))
15 print("Cumulative Distribution Function: ", f.cdf(x, nu1, nu2))
16 print("Survival Function (1-cdf): ", f.sf(x, nu1, nu2))
17 print("Percent Point Function (inverse of cdf): ", f.ppf(alpha, nu1, nu2))
18 print("Inverse Survival Function (inverse of (1-cdf)): ", f.isf(alpha, nu1, nu2)) # Useful for finding the critical value
[5] ✓ 0.0s
...
Expectation: 1.1428571428571428
Variance: 0.34467120181405897
Probability Distribution Function: 0.14612805411889215
Cumulative Distribution Function: 0.92524818763759
Survival Function (1-cdf): 0.07745718123624103
Percent Point Function (inverse of cdf): 2.0171044619502414
Inverse Survival Function (inverse of (1-cdf)): 0.5257867781767004
```

Python

## EXPERIMENT NO. 10

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 14.04.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> To understand how to compute critical value of standardized test statistic.
<b>Outcome:</b> Students will compute critical value from the standard normal distribution table
<b>Problem Statement:</b> Write python code to find critical value in normal distribution
<b>Background Study:</b> The normal distribution is the most common type of continuous distribution for which we need to compute the standardized test statistic and the its critical value.
<b>Question Bank:</b> <ol style="list-style-type: none"><li>What is standardized test statistic? A standardized test statistic is a value that measures how many standard deviations a sample statistic (such as the sample mean or sample proportion) is away from the expected value under the null hypothesis. It is often used in hypothesis testing to determine the statistical significance of the sample statistic.</li><li>How do we compute the critical value of the standardized test statistic? The critical value of the standardized test statistic is the value that separates the rejection region from the non-rejection region in a hypothesis test. It is determined based on the level of significance, the degrees of freedom, and the type of test (one-tailed or two-tailed).</li></ol>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 10

Write python code to find critical value in normal distribution.

```
1 # import the required libraries
2 from scipy.stats import norm
3
4 # Set the probability value (alpha)
5 alpha = 0.05
6
7 # Set the mean and standard deviation of the normal distribution
8 mu = 0
9 sigma = 1
10
11 # calculate the critical value
12 critical_value = norm.ppf(1 - alpha/2, loc=mu, scale=sigma)
13
14 # Print the critical value
15 print("The critical value is:", critical_value)
16
```

[2] ... The critical value is: 1.959963984540054

Python

## EXPERIMENT NO. 11

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b> 28.04.2023
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> To understand central limit theorem and its applications
<b>Outcome:</b> Students will be able to apply central limit theorem on the sampling distribution
<b>Problem Statement:</b> Write a python program to Implement central limit theorem
<b>Background Study:</b> When sample size approaches 30, central limit theorem kicks in
<b>Question Bank:</b>
1. What is central limit theorem? The central limit theorem states that the distribution of the sample means approaches a normal distribution as the sample size increases, regardless of the underlying distribution of the population from which the samples are taken.
2. How is it related to sample size? The central limit theorem is directly related to sample size. It states that as the sample size increases, the distribution of the sample means approaches a normal distribution, regardless of the underlying distribution of the population from which the samples are taken. This means that the larger the sample size, the more accurate the estimate of the population mean will be, and the more confident we can be in our results. In other words, as the sample size increases, the variability of the sample means decreases, and the distribution becomes more concentrated around the true population mean.
3. Can we apply central limit theorem for both small and large sample size? The central limit theorem is typically only applicable for large sample sizes. While there is no specific cut-off for what constitutes a "large" sample size, a general rule of thumb is that the sample size should be at least 30 in order for the central limit theorem to apply. For sample sizes smaller than 30, the central limit theorem may not hold, and other techniques may be necessary to estimate population parameters.

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 11



Write a python program to Implement central limit theorem.

```
1 # importing required libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
[1] ✓ 1.6s
```

Python

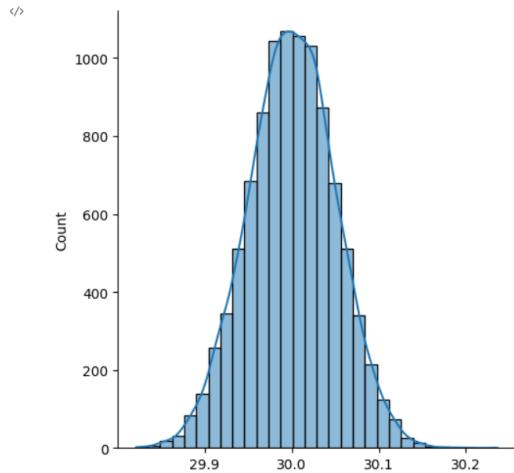
Central Limit Theorem

The central limit theorem states that if you have a population with mean  $\mu$  and standard deviation  $\sigma$  and take sufficiently large random samples from the population with replacement , then the distribution of the sample means will be approximately normally distributed.

```
1 # implementing the central limit theorem
2 mu = 30
3 sigma = 5
4
5 samples = np.zeros(10000)
6
7 for i in range(10000):
8     samples[i] = np.random.normal(mu, sigma, 10000).mean()
9
10 print("Samples mean: ", samples.mean())
11
12 sns.displot(samples, bins=30, kde=True)
[14] ✓ 3.4s
```

Python

<seaborn.axisgrid.FacetGrid at 0x2a248f9a410>



## EXPERIMENT NO. 12

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual (github.com)</a>
<b>Date:</b>
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> To understand the steps of hypothesis testing for single population
<b>Outcome:</b> Ability to test hypothesis on the population
<b>Problem Statement:</b> Write a python program to conduct a hypothesis test for unknown variance.
<b>Background Study:</b> Hypothesis testing allows the statistical evaluation of the claim made
<b>Question Bank:</b>
<ol style="list-style-type: none"><li>1. What is hypothesis testing? Hypothesis testing is a statistical method for making decisions or drawing conclusions about a population based on a sample of data. It involves formulating two hypotheses - a null hypothesis (<math>H_0</math>) and an alternative hypothesis (<math>H_a</math>) - and testing whether the sample data provides sufficient evidence to reject the null hypothesis in favor of the alternative hypothesis.</li><li>2. What are the steps of hypothesis testing? Step 1: State the Null Hypothesis. Note. Step 2: State the Alternative Hypothesis. Step 3: Set <math>\alpha</math> Step 4: Collect Data. Step 5: Calculate a test statistic. Step 6: Construct Acceptance / Rejection regions. Step 7: Based on steps 5 and 6, draw a conclusion about <math>H_0</math>.</li><li>3. What is level of significance? The level of significance, denoted by the symbol alpha (<math>\alpha</math>), is the maximum allowable probability of making a Type I error in a statistical hypothesis test. A Type I error occurs when the null hypothesis is rejected when it is actually true, which means that the test has produced a false positive result.</li></ol>

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 12

Write a python program to conduct a hypothesis test for unknown variance.

```
1 # importing required libraries
2 import math
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 import scipy.stats as stats
[21] ✓ 0.2s
```

Python

#### One Sample T-Test

```
1 # Set the random seed for reproducibility
2 np.random.seed(8)
3
4 # Generate a population age distribution with a mean of 40 and a size of 200,000
5 population_ages1 = stats.poisson.rvs(loc=18, mu=40, size=150000)
6 population_ages2 = stats.poisson.rvs(loc=18, mu=10, size=50000)
7 population_ages = np.concatenate((population_ages1, population_ages2))
8
9 # Generate a sample from a specific state (e.g., Minnesota) with a mean of 35 and a sample size of 50
10 minnesota_ages1 = stats.poisson.rvs(loc=18, mu=35, size=30)
11 minnesota_ages2 = stats.poisson.rvs(loc=18, mu=10, size=20)
12 minnesota_ages = np.concatenate((minnesota_ages1, minnesota_ages2))
13
14 # Print the mean age of the population and the sample
15 print("Mean age of the population:", population_ages.mean())
16 print("Mean age of the sample:", minnesota_ages.mean())
17
[22] ✓ 0.1s
```

Python

... Mean age of the population: 50.48598  
Mean age of the sample: 42.92

```
1 # Calculate the t-statistic and p-value for the one-sample t-test
2 t_stat, p_val = stats.ttest_1samp(a=minnesota_ages, popmean=population_ages.mean())
3
4 print("The t-statistic is {} and the p-value is {}".format(t_stat, p_val))
[23] ✓ 0.0s
```

Python

... The t-statistic is -4.155007645348416 and the p-value is 0.00012990129848405887.

```

1 # Check the quantiles for a 95% confidence level with 49 degrees of freedom
2 t_quantile_low = stats.t.ppf(q=0.025, df=49)
3 t_quantile_high = stats.t.ppf(q=0.975, df=49)
4 print("The t-quantiles for a 95% confidence level with 49 degrees of freedom are {} and {}".format(t_quantile_low, t_quantile_high))
[24]   ✓ 0.1s
...
The t-quantiles for a 95% confidence level with 49 degrees of freedom are -2.0095752344892093 and 2.009575234489209.
Python
```

```

1 # Calculate the p-value with the t-statistic and 49 degrees of freedom
2 p_val_calculated = stats.t.cdf(x=t_stat, df=49) * 2
3 print("The calculated p-value is {}".format(p_val_calculated))
[25]   ✓ 0.0s
...
The calculated p-value is 0.00012990129848405887.
Python
```

```

1 # Calculate the standard error of the mean for the sample
2 sigma = minnesota_ages.std() / math.sqrt(50)
[26]   ✓ 0.0s
...
Python
```

```

1 # Construct a 95% confidence interval for the sample mean
2 conf_int_95 = stats.t.interval(alpha=0.95, df=49, loc=minnesota_ages.mean(), scale=sigma)
3 print("The 95% confidence interval for the sample mean is {}".format(conf_int_95))
[27]   ✓ 0.0s
...
The 95% confidence interval for the sample mean is (39.29748102093582, 46.542518979064184).
C:\Users\mainp\AppData\Local\Temp\ipykernel_14524\4103696102.py:2: DeprecationWarning: Use of keyword argument 'alpha' for method 'interval' is deprecated and will be removed in SciPy 1.0.0.
conf_int_95 = stats.t.interval(alpha=0.95, df=49, loc=minnesota_ages.mean(), scale=sigma)
Python
```

```

1 # Construct a 99% confidence interval for the sample mean
2 conf_int_99 = stats.t.interval(alpha=0.99, df=49, loc=minnesota_ages.mean(), scale=sigma)
3 print("The 99% confidence interval for the sample mean is {}".format(conf_int_99))
[28]   ✓ 0.2s
...
The 99% confidence interval for the sample mean is (38.08904034203606, 47.750959657963946).
C:\Users\mainp\AppData\Local\Temp\ipykernel_14524\176364369.py:2: DeprecationWarning: Use of keyword argument 'alpha' for method 'interval' is deprecated and will be removed in SciPy 1.0.0.
conf_int_99 = stats.t.interval(alpha=0.99, df=49, loc=minnesota_ages.mean(), scale=sigma)
Python
```

## Two Sample T-Test

```

1 # Generate a sample of voter age data for Wisconsin
2 np.random.seed(10)
3 wisconsin_ages1 = stats.poisson.rvs(loc=18, mu=30, size=30)
4 wisconsin_ages2 = stats.poisson.rvs(loc=18, mu=15, size=20)
5 wisconsin_ages = np.concatenate((wisconsin_ages1, wisconsin_ages2))
[29]   ✓ 0.1s
...
Python
```

```

1 # Print the mean of the Wisconsin sample
2 print("Wisconsin Sample Mean:", wisconsin_ages.mean())
[30]   ✓ 0.1s
...
Wisconsin Sample Mean: 41.8
Python
```

```

1 # Perform a two-sample t-test to compare the Wisconsin sample to the Minnesota sample
2 t_stat, p_val = stats.ttest_ind(a=minnesota_ages, b=wisconsin_ages, equal_var=False)
[31]   ✓ 0.0s
...
Python
```

```

1 # Print the p-value
2 print("P-Value:", p_val)
[32] ✓ 0.0s
... P-Value: 0.6175514813720644

```

Python

---

```

1 # Check if the p-value is less than the significance level of 0.05
2 if p_val < 0.05:
3     print("Reject null hypothesis, the means are different.")
4 else:
5     print("Fail to reject null hypothesis, the means are the same.")
[33] ✓ 0.0s
... Fail to reject null hypothesis, the means are the same.

```

Python

### Paired T-Test

```

1 # set the random seed for reproducibility
2 np.random.seed(11)
3
4 # generate weight data before and after treatment
5 before = stats.norm.rvs(scale=30, loc=250, size=100)
6 after = before + stats.norm.rvs(scale=5, loc=-1.25, size=100)
7
8 # create a dataframe to store the weight data
9 weight_df = pd.DataFrame({"weight_before": before,
10                           "weight_after": after,
11                           "weight_change": after - before})
12
13 # print a summary of the weight data
14 print(weight_df.describe())
[34] ✓ 0.1s
...    weight_before  weight_after  weight_change
count      100.000000   100.000000   100.000000
mean       250.345546  249.115171   -1.230375
std        28.132539   28.422183   4.783696
min        170.400443  165.913930  -11.495286
25%       230.421042  229.148236   -4.046211
50%       250.830805  251.134089   -1.413463
75%       270.637145  268.927258   1.738673
max       314.700233  316.720357   9.759282

```

---

```

1 # conduct a paired t-test to check if the means of the samples differ
2 t_statistic, p_value = stats.ttest_rel(a=before, b=after)
[35] ✓ 0.1s

```

Python

---

```

1 # print the p-value
2 print("Paired t-test p-value:", p_value)
[36] ✓ 0.1s
... Paired t-test p-value: 0.011596444318439859

```

Python

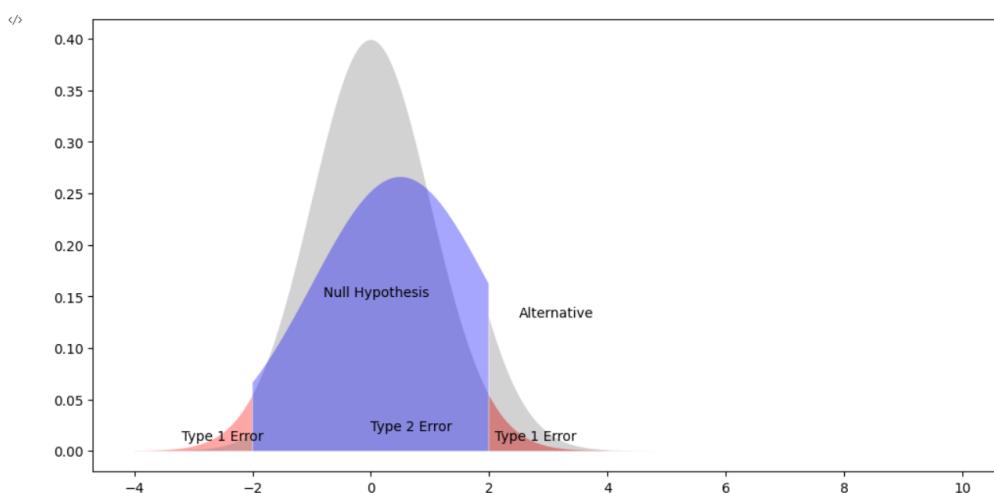
## Type I and Type II Errors

```

1 # Generate three normal distributions with different means and standard deviations
2 null_distribution = stats.norm(loc=0, scale=1)
3 alternative_distribution = stats.norm(loc=0.5, scale=1)
4 observed_distribution = stats.norm(loc=0.5, scale=1.5)
5
6 # Set the plot size
7 plt.figure(figsize=(12, 10))
8
9 # Fill in the areas corresponding to the null distribution
10 plt.fill_between(x=np.arange(-4, -2, 0.01),
11                   y1= null_distribution.pdf(np.arange(-4, -2, 0.01)),
12                   facecolor='red', alpha=0.35)
13 plt.fill_between(x=np.arange(-2, 2, 0.01),
14                   y1= null_distribution.pdf(np.arange(-2, 2, 0.01)),
15                   facecolor='grey', alpha=0.35)
16 plt.fill_between(x=np.arange(2, 4, 0.01),
17                   y1= null_distribution.pdf(np.arange(2, 4, 0.01)),
18                   facecolor='red', alpha=0.5)
19
20 # Fill in the areas corresponding to the alternative and observed distributions
21 plt.fill_between(x=np.arange(-4, -2, 0.01),
22                   y1= alternative_distribution.pdf(np.arange(-4, -2, 0.01)),
23                   facecolor='grey', alpha=0.35)
24 plt.fill_between(x=np.arange(-2, 2, 0.01),
25                   y1= observed_distribution.pdf(np.arange(-2, 2, 0.01)),
26                   facecolor='blue', alpha=0.35)
27 plt.fill_between(x=np.arange(2, 10, 0.01),
28                   y1= alternative_distribution.pdf(np.arange(2, 10, 0.01)),
29                   facecolor='grey', alpha=0.35)
30
31 # Add labels to the plot
32 plt.text(x=-0.8, y=0.15, s="Null Hypothesis")
33 plt.text(x=2.5, y=0.13, s="Alternative")
34 plt.text(x=2.1, y=0.01, s="Type 1 Error")
35 plt.text(x=-3.2, y=0.01, s="Type 1 Error")
36 plt.text(x=0, y=0.02, s="Type 2 Error")
[41] ✓ 0.7s
... Text(0, 0.02, 'Type 2 Error')

```

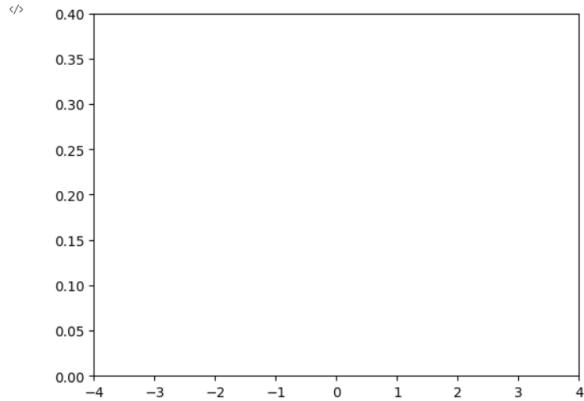
Python



```
1 # Set the x and y limits of the plot
2 plt.xlim([-4, 4])
3 plt.ylim([0, 0.4])
4
● 5 # Calculate the Type II error rate for the observed distribution compared to the null distribution
6 alpha = 0.05 # Significance level
7 z_alpha = stats.norm.ppf(1 - alpha/2) # Two-tailed z-score corresponding to alpha
8 z_beta = (0.5 - 0) / (1.5 / np.sqrt(100)) + z_alpha # z-score corresponding to beta
9 beta = 1 - stats.norm.cdf(z_beta, loc=0, scale=1) # Calculate beta
10 power = 1 - beta # Calculate power
11
12 # Print the Type II error rate and power
13 print(f"Type II error rate: {beta:.3f}")
14 print(f"Power: {power:.3f}")
15
16 # Show the plot
17 plt.show()
```

[38] ✓ 0.3s

... Type II error rate: 0.000



Python

## EXPERIMENT NO. 13

<b>Student Name and Roll Number:</b> Piyush Gambhir – 21CSU349
<b>Semester /Section:</b> Semester-IV – AIML-B (AL-3)
<b>Link to Code:</b> <a href="https://github.com/Piyush-Gambhir/NCU-CSL227-ACS-Lab_Manual">Piyush-Gambhir/NCU-CSL227-ACS-Lab Manual (github.com)</a>
<b>Date:</b>
<b>Faculty Signature:</b>
<b>Marks/Grade:</b>

<b>Objective(s):</b> To understand the steps of hypothesis testing for two populations
<b>Outcome:</b> Students will analysis the steps to test hypothesis for two populations
<b>Problem Statement:</b> Write a python program to conduct hypothesis testing for two populations
<b>Background Study:</b> Hypothesis testing allows the statistical evaluation of the claim made
<b>Question Bank:</b>
<b>1. What is meant by T-test? Explain with example.</b> A t-test is a statistical test that compares the means of two samples. It is used in hypothesis testing, with a null hypothesis that the difference in group means is zero and an alternate hypothesis that the difference in group means is different from zero.
<b>2. What is difference between T-test and paired t-test?</b> Two-sample t-test is used when the data of two samples are statistically independent, while the paired t-test is used when data is in the form of matched pairs.
<b>3. What is significance of T-Test?</b> The One Sample t Test examines whether the mean of a population is statistically different from a known or hypothesized value. The One Sample t Test is a parametric test. This test is also known as: Single Sample t Test.

## Student Work Area

### Algorithm/Flowchart/Code/Sample Outputs

#### Experiment 13

Write a python program to conduct hypothesis testing for two populations

```

1 # importing required libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from scipy import stats
[5] ✓ 0.0s Python

1 # Generate two random samples
2 sample1 = np.random.normal(loc=10, scale=2, size=50)
3 sample2 = np.random.normal(loc=12, scale=2, size=50)
4
5 # Calculate the sample statistics
6 mean1 = np.mean(sample1)
7 mean2 = np.mean(sample2)
8 var1 = np.var(sample1, ddof=1)
9 var2 = np.var(sample2, ddof=1)
10 n1 = len(sample1)
11 n2 = len(sample2)
[6] ✓ 0.0s Python

1 # Calculate the pooled standard deviation
2 s_p = np.sqrt((n1 - 1) * var1 + (n2 - 1) * var2) / (n1 + n2 - 2)
3
4 # Calculate the test statistic
5 t = (mean1 - mean2) / (s_p * np.sqrt(1/n1 + 1/n2))
6
7 # Calculate the degrees of freedom
8 df = n1 + n2 - 2
9
10 # Calculate the p-value
11 p_value = stats.t.sf(abs(t), df) * 2
12
13 # Set the significance level
14 alpha = 0.05
[7] ✓ 0.1s Python

1 # Compare the p-value to the significance level
2 if p_value < alpha:
3     print("We reject the null hypothesis that the two populations have equal means.")
4 else:
5     print("We fail to reject the null hypothesis that the two populations have equal means.")
[8] ✓ 0.1s Python
...
We reject the null hypothesis that the two populations have equal means.

```