

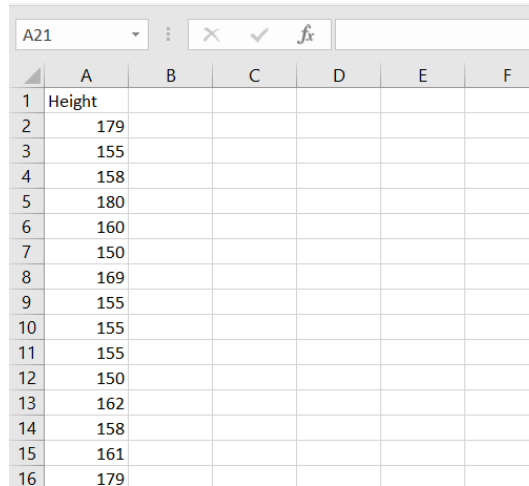
INDEX				
	Topic	Date	Page No.	Sign
1.	Practical 1	17/09/22	3	
a.	Write a program for obtaining descriptive statistics of data		3	
b.	Import data from different data sources (from Excel, csv, mysql, sql server, oracle to R /Python /Excel)		5	
2.	Practical 2	23/09/22	9	
a.	Design a survey form for a given case study, collect the primary data and analyse it.		9	
b.	Perform analysis of given secondary data.		11	
3.	Practical 3	24/09/22	14	
a.	Perform testing of hypothesis using one sample t-test.		14	
b.	Write a program for t-test comparing two means for independent samples.		15	
c.	Perform testing of Hypothesis using paired t-test.		17	
4.	Practical 4	01/10/22	19	
a.	Perform testing of hypothesis using chi-squared goodness-of-fit test.		19	
b.	Perform testing of hypothesis using chi-squared test of independence.		20	
5.	Practical 5	08/10/22	22	
a.	Perform testing of hypothesis using Z-test.		22	
6.	Practical 6	15/10/22	24	
a.	Perform testing of hypothesis using One-way ANOVA.		24	
b.	Perform testing of hypothesis using Two-way ANOVA.		26	
c.	Perform testing of hypothesis using MANOVA.		28	
7.	Practical 7	22/10/22	31	
a.	Perform the Random sampling for the given data and analyse it.		31	

b.	Perform the Stratified sampling for the given data and analyse it		33	
8.	Practical 8	12/11/22	36	
a.	Write a program for computing different correlation		36	
9.	Practical 9	19/11/22	36	
a.	Write a program to Perform linear regression for prediction.		36	
b.	Polynomial Regression		41	
10.	Practical 10	26/11/22	44	
a.	Multiple Linear Regression		44	
b.	Logistic Regression		45	

Practical 1

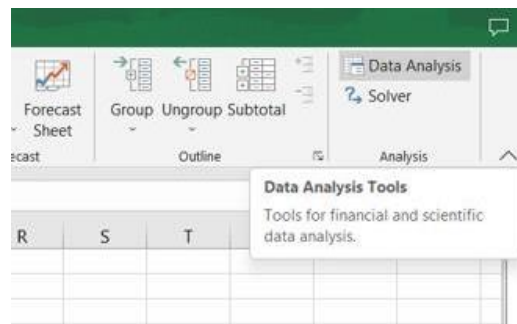
A. Write a program for obtaining descriptive statistics of data.

Step 1: Open your data in excel

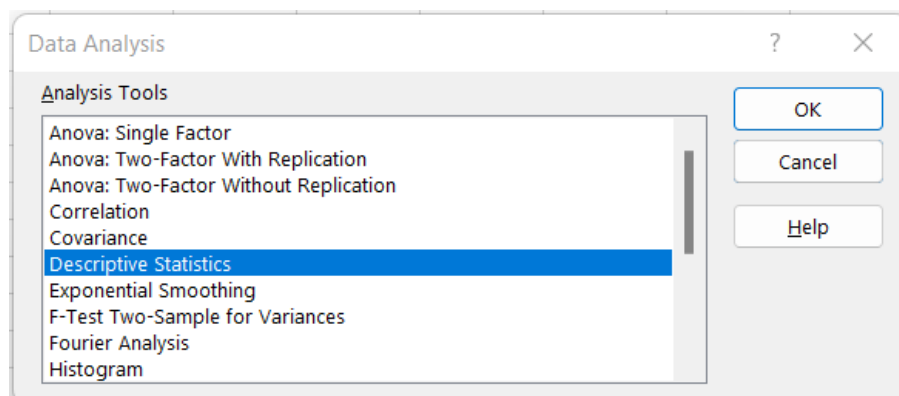


	A	B	C	D	E	F
1	Height					
2	179					
3	155					
4	158					
5	180					
6	160					
7	150					
8	169					
9	155					
10	155					
11	155					
12	150					
13	162					
14	158					
15	161					
16	179					

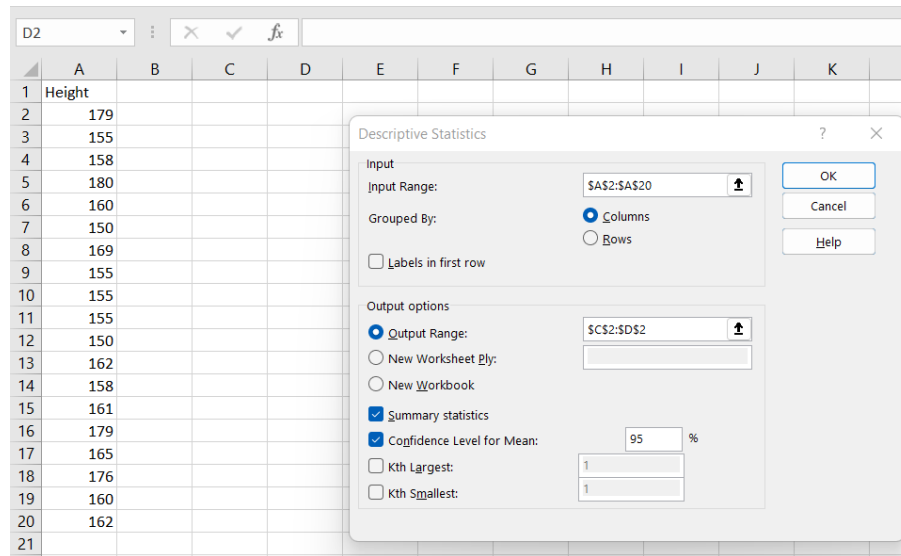
Step 2: From the Data Tool in the ribbon choose Data Analysis.



Step 3: Select the option of descriptive statistics



Step 4: Select an input range, output range, and check summary statistics and confidence level.



Step 5: The output may appear as follows.

1	Height		
2	179		
3	155		
4	158	Mean	162.5789474
5	180	Standard Error	2.212822533
6	160	Median	160
7	150	Mode	155
8	169	Standard Deviation	9.645469803
9	155	Sample Variance	93.03508772
10	155	Kurtosis	-0.561004927
11	155	Skewness	0.753286117
12	150	Range	30
13	162	Minimum	150
14	158	Maximum	180
15	161	Sum	3089
16	179	Count	19
17	165	Confidence Level(95.0%)	4.648967631
18	176		
19	160		
20	162		

Conclusion: Descriptive statistics is a form of data analysis that involves the use of numerical measures to describe the characteristics of a given set of data. The goal of descriptive statistics is to summarize the data into a meaningful and useful form. The most common forms of descriptive statistics include the mean, median, mode, and standard deviation. These measures allow researchers to understand the characteristics of the data, such as its spread and distribution. Descriptive statistics can also be used to make predictions about future data. For example, if the mean of a set of data is consistently

higher than the median, it may be possible to predict that future data will also have a higher mean.

B. Import data from different data sources (from Excel, csv, mysql, sql server, oracle to R /Python /Excel)

Code:

1 - SQLite3

```
import sqlite3 as sq

import pandas as pd

#####

Base='C:/VKHCG'

sDatabaseName=Base + '/01-Vermeulen/04-
Transform//SQLite/vermeulen.db'

conn = sq.connect(sDatabaseName)

#####

sFileName='C:/VKHCG/01-Vermeulen/01-Retrieve/01-EDS/02-
Python/Retrieve_IP_DATA.csv'

print('Loading :',sFileName)

IP_DATA_ALL_FIX=pd.read_csv(sFileName,header=0,low_memory=False)

IP_DATA_ALL_FIX.index.names = ['RowIDCSV']

sTable='IP_DATA_ALL'

print('Storing :',sDatabaseName,' Table:',sTable)

IP_DATA_ALL_FIX.to_sql(sTable, conn, if_exists="replace")

print('Loading :',sDatabaseName,' Table:',sTable)

TestData=pd.read_sql_query("select * from IP_DATA_ALL;", conn)
```

```

print('## Data Values')

print(TestData)

print('## Data Profile')

print('Rows :', TestData.shape[0])

print('Columns :', TestData.shape[1])

print('Successful')

```

Output:

```

Loading : C:/VKHCG/01-Vermeulen/01-Retrieve/01-EDS/02-Python/
Retrieve_IP_DATA.csv
Storing : C:/VKHCG/01-Vermeulen/04-Transform//SQLite/
vermeulen.db Table: IP_DATA_ALL
Loading : C:/VKHCG/01-Vermeulen/04-Transform//SQLite/
vermeulen.db Table: IP_DATA_ALL
## Data Values
      RowIDCSV  RowID  ...  First.IP.Number
Last.IP.Number
0              0      0  ...      692781056
692781567
1              1      1  ...      692781824
692783103
2              2      2  ...      692909056
692909311
3              3      3  ...      692909568
692910079
4              4      4  ...      693051392
693052415
...          ...    ...  ...          ...
...
1247497  1247497  1247497  ...      1068157850
1068157850
1247498  1247498  1247498  ...      1334409600
1334409607
1247499  1247499  1247499  ...      1596886528
1596886783
1247500  1247500  1247500  ...      1742189568
1742190591
1247501  1247501  1247501  ...      1905782573
1905782573

[1247502 rows x 11 columns]
## Data Profile
Rows : 1247502

```

2 - Excel

```

import os

import pandas as pd

```

```
Base='C:/VKHCG'

sFileDir=Base + '/01-Vermeulen/01-Retrieve/01-EDS/02-Python'

CurrencyRawData = pd.read_excel('C:/VKHCG/01-Vermeulen/00-
RawData/Country_Currency.xlsx')

sColumns = ['Country or territory', 'Currency', 'ISO-4217']

CurrencyData = CurrencyRawData[sColumns]

CurrencyData.rename(columns={'Country or territory': 'Country',
'ISO-4217':
'CurrencyCode'}, inplace=True)

CurrencyData.dropna(subset=['Currency'],inplace=True)

CurrencyData['Country'] = CurrencyData['Country'].map(lambda x:
x.strip())

CurrencyData['Currency'] = CurrencyData['Currency'].map(lambda
x:
x.strip())

CurrencyData['CurrencyCode'] =
CurrencyData['CurrencyCode'].map(lambda x:
x.strip())

print(CurrencyData)

print('~~~~~ Data from Excel Sheet Retrived Successfully
~~~~~')

sFileName=sFileDir + '/Retrieve-Country-Currency.csv'

CurrencyData.to_csv(sFileName, index = False)
```

```

IPython console
Console 1/A
255         Tuvalu      ...      AUD
257         Uganda     ...      UGX
258         Ukraine     ...      UAH
259         United Arab Emirates ...  AED
260         United Kingdom ...      GBP
261         United States of America ...  USD
262         Uruguay     ...      UYU
263         US Virgin Islands (USA) ...  USD
264         Uzbekistan  ...      UZS
266         Vanuatu     ...      VUV
267         Vatican City (Holy See) ...  EUR
268         Venezuela   ...      VEF
269         Vietnam     ...      VND
271         Wake Island (USA) ...      USD
272         Wallis and Futuna (France) ...  XPF
274         Yemen       ...      YER
276         Zambia      ...      ZMW
277         Zimbabwe    ...      USD

[253 rows x 3 columns]
Data from Excel Sheet Retrived Successfully

In [8]:
IPython console  History log

```

Country	Currency	CurrencyCode
1 Afghanistan	Afghan a/AfN	
2 Albania	Albanian a/ALL	
3 Algeria	Algerian d/DZD	
4 American Samoa	United St/USD	
5 Andorra	European EUR	
6 Angola	Angolan k/AOA	
7 Anguilla	East Carib XCD	
8 Antigua and Barbuda	East Carib XCD	
9 Argentina	Argentine ARS	
10 Armenia	Armenian AMD	
11 Aruba	Aruban f/AWG	
12 Ascension and Saint Helena	GBP	
13 Australia	Australian AUD	
14 Austria	European EUR	
15 Azerbaijan	Azerbaijani AZN	
16 Bahamas	Bahamian BSD	
17 Bahrain	Bahraini d/BHD	
18 Bangladesh	Bangladeshi BDT	
19 Barbados	Barbadian BBD	
20 Belarus	Belarusian BYN	
21 Belgium	European EUR	
22 Belize	Belize dol/BZD	
23 Benin	West Africa XOF	
24 Bermuda	Bermudian BMD	
25 Bhutan	Bhutanese BTN	
26 Bolivia	Bolivian b/BOB	
27 Botswana	United St/Pula	

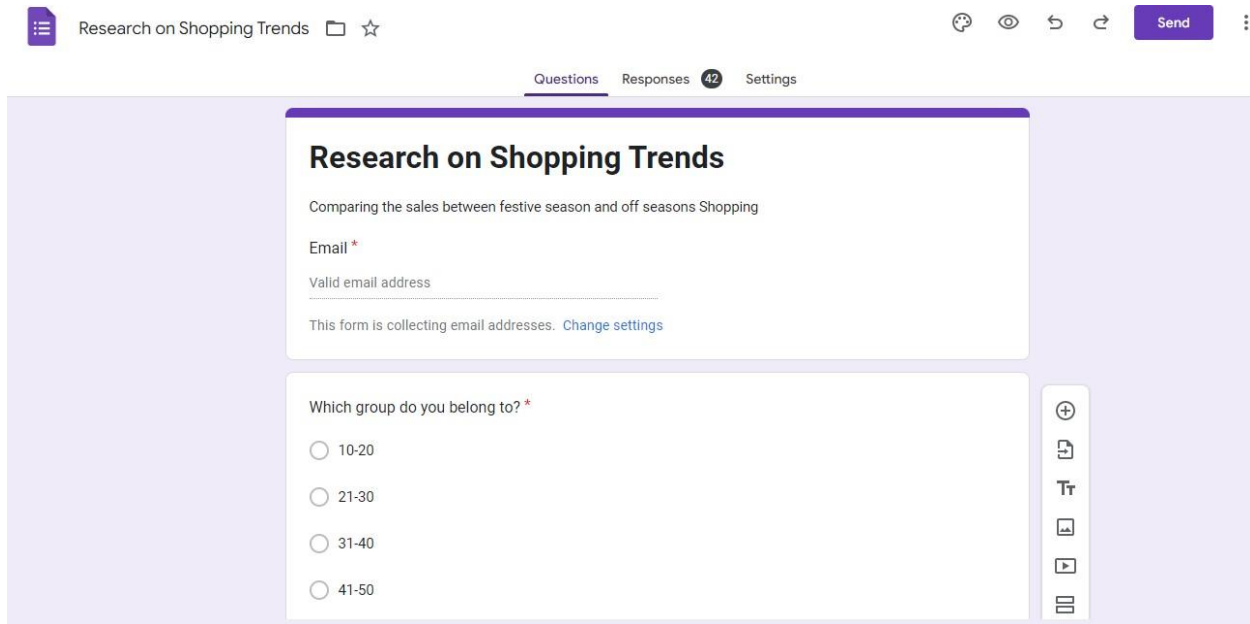
Conclusion: SQLite/Python/Excel

This type of coding are very much helpful to get the correct exact output we can compare all the outputs for getting the more exact value for the accurate answering. Using Python and Excel helps us to do any formulation for finding accuracy getting the answers/the outputs very quickly and easily. SQLite is a database to store various data easily and easy to access.

Practical 2

A. Design a survey form for a given case study, collect the primary data and analyse it.

Step 1: In order to make a survey we used Google forms. Our case study's aim was to find out which age group did Shopping in festive seasons.

The image shows a Google Form titled "Research on Shopping Trends" with the subtitle "Comparing the sales between festive season and off seasons Shopping". The form includes an email field with a red asterisk and a "Valid email address" error message. Below this is a question "Which group do you belong to? *" with four radio button options: "10-20", "21-30", "31-40", and "41-50". The form is displayed in a preview mode with a purple header and a sidebar on the right containing icons for adding questions, duplicating, deleting, and other actions. The top navigation bar shows "Questions", "Responses" (with a count of 42), and "Settings".

Step 2: The questions and options entered were as follows.

Q1.Which group do you belong to?

(a)10-20 (b)21-30 (c)31-40 (d)41-50 (e)51 Above

Q2.What is your Gender?

(a)Male (b)Female

Q3.What is your Occupation?

(a)Business (b)Student (c)Employee (d)Retired (e)Self Employed (f)Household

Q4.What's your Monthly Income?

(a)10,000 - 20,000 (b)21,000 - 30,000 (c)31,000 - 40,000 (d)41,000 and above (e)N/A

Q5.How often do you shop?

(a)Monthly (b)Occasionally (c)Festive Seasons (d)Rarely

Q6.Purpose of Shopping?

(a)Personal Use (b)Gifting

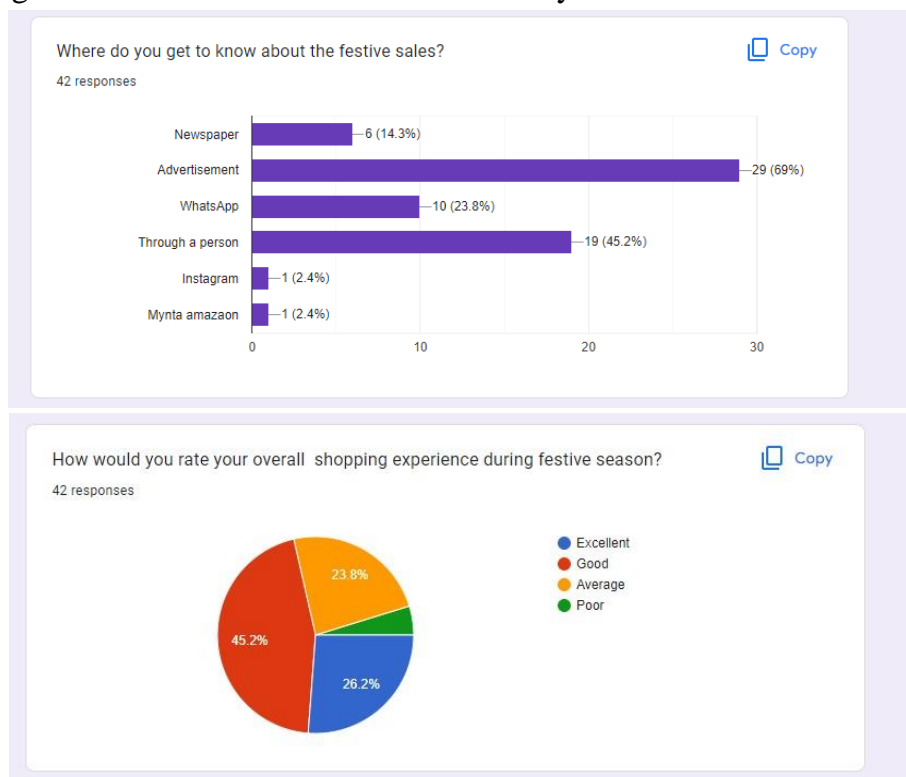
Q7.Where do you get to know about the festive sales?

(a)Newspaper (b)Advertisement (c)WhatsApp (d)Through a person

Q8.How would you rate your overall shopping experience during festive season?

(a)Excellent (b)Good (c)Average (d)Poor

Step 3 :Google forms tool was chosen because analysis becomes easier on Google forms.



Conclusion: we have successfully made the survey form for the given case study.

In the above practical, we have successfully demonstrated Designing of survey form to collect the shopping trend data of people from different age group. From the above survey we can observe that 42.9% people shop on monthly basis and in this 57.1% women are involved shopping. 45.2% people have Good shopping experience whereas 26.2% people have Excellent shopping experience. Advertisement (69%) is the main source for people to know about festive shopping.

B. Perform analysis of given secondary data.

Step 1: Open Data in excel

1	age	male	female	total	
2	0 to 4	328759	307079	635838	
3	4 to 9	315119	293664	608783	
4	10 to 14	311456	290598	602054	
5	15 to 19	312831	293313	606144	
6	20 to 24	311077	295739	606816	
7	25 to 29	284258	273379	557637	
8	30 to 34	255596	247383	502979	
9	35 to 39	248575	241938	490513	
10	40 to 44	232217	226914	459131	
11	45 to 49	202633	201142	403775	
12	50 to 54	176241	176440	352681	
13	55 to 59	153494	156283	309777	
14	60 to 64	114194	121200	235394	
15	65 to 69	83129	92071	175200	
16	70 to 74	65266	77990	143256	
17	75 to 79	43761	56895	100656	
18	80 to 84	25060	37873	62933	
19	85+	14164	28156	42320	
20					
21		3477830	3418057	6895887	

Step 2: Calculate the total sum of each column. Select the cell for Sum→ add formula SUM in formula bar→ select the range.

1	age	male	female	total	
2	0 to 4	328759	307079	635838	
3	4 to 9	315119	293664	608783	
4	10 to 14	311456	290598	602054	
5	15 to 19	312831	293313	606144	
6	20 to 24	311077	295739	606816	
7	25 to 29	284258	273379	557637	
8	30 to 34	255596	247383	502979	
9	35 to 39	248575	241938	490513	
10	40 to 44	232217	226914	459131	
11	45 to 49	202633	201142	403775	
12	50 to 54	176241	176440	352681	
13	55 to 59	153494	156283	309777	
14	60 to 64	114194	121200	235394	
15	65 to 69	83129	92071	175200	
16	70 to 74	65266	77990	143256	
17	75 to 79	43761	56895	100656	
18	80 to 84	25060	37873	62933	
19	85+	14164	28156	42320	
20					
21		3477830	3418057	6895887	

Step 3: Calculate the percentage of male in cell E. Use formula $-1*100*B2/(\$D\21

SUM						
	A	B	C	D	E	F
1	age	male	female	total	Male(%)	
2	0 to 4	328759	307079	635838	=-1*100*B2/\$D\$21	
3	4 to 9	315119	293664	608783		
4	10 to 14	311456	290598	602054		
5	15 to 19	312831	293313	606144		
6	20 to 24	311077	295739	606816		
7	25 to 29	284258	273379	557637		
8	30 to 34	255596	247383	502979		
9	35 to 39	248575	241938	490513		
10	40 to 44	232217	226914	459131		
11	45 to 49	202633	201142	403775		
12	50 to 54	176241	176440	352681		

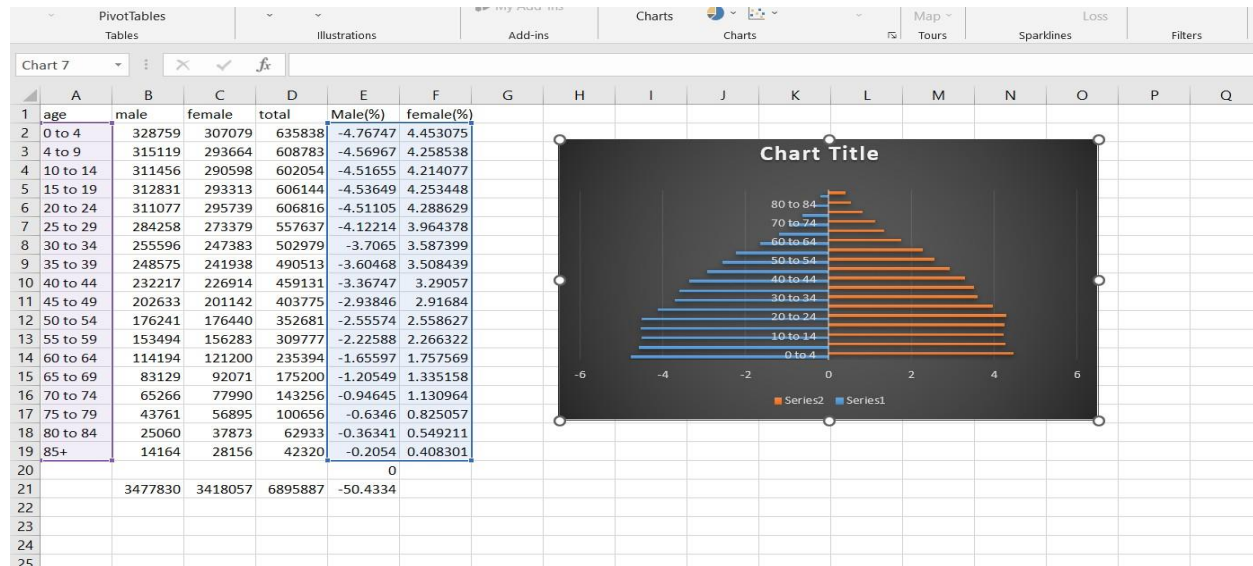
Step 4: Calculate the percentage female in cell F. Use formula $100 \times C2 / \$D\21

F2						
	A	B	C	D	E	F
1	age	male	female	total	Male(%)	female(%)
2	0 to 4	328759	307079	635838	-4.76747	4.453075
3	4 to 9	315119	293664	608783	-4.56967	4.258538
4	10 to 14	311456	290598	602054	-4.51655	4.214077
5	15 to 19	312831	293313	606144	-4.53649	4.253448
6	20 to 24	311077	295739	606816	-4.51105	4.288629
7	25 to 29	284258	273379	557637	-4.12214	3.964378
8	30 to 34	255596	247383	502979	-3.7065	3.587399
9	35 to 39	248575	241938	490513	-3.60468	3.508439
10	40 to 44	232217	226914	459131	-3.36747	3.29057
11	45 to 49	202633	201142	403775	-2.93846	2.91684
12	50 to 54	176241	176440	352681	-2.55574	2.558627
13	55 to 59	153494	156283	309777	-2.22588	2.266322
14	60 to 64	114194	121200	235394	-1.65597	1.757569
15	65 to 69	83129	92071	175200	-1.20549	1.335158

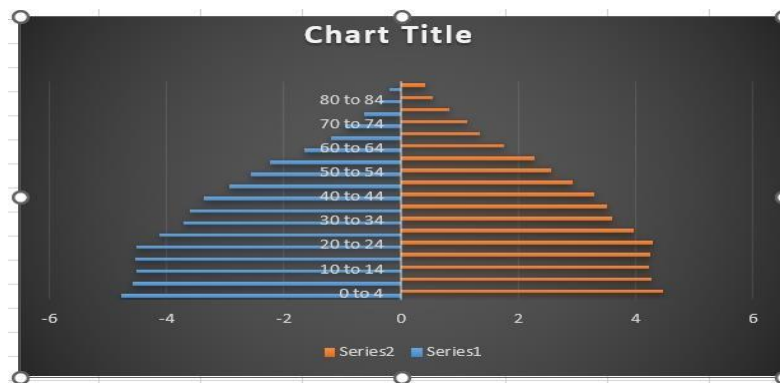
Step 5: The modified data may seem as follows.

	A	B	C	D	E	F
1	age	male	female	total	Male(%)	female(%)
2	0 to 4	328759	307079	635838	-4.76747	4.453075
3	4 to 9	315119	293664	608783	-4.56967	4.258538
4	10 to 14	311456	290598	602054	-4.51655	4.214077
5	15 to 19	312831	293313	606144	-4.53649	4.253448
6	20 to 24	311077	295739	606816	-4.51105	4.288629
7	25 to 29	284258	273379	557637	-4.12214	3.964378
8	30 to 34	255596	247383	502979	-3.7065	3.587399
9	35 to 39	248575	241938	490513	-3.60468	3.508439
10	40 to 44	232217	226914	459131	-3.36747	3.29057
11	45 to 49	202633	201142	403775	-2.93846	2.91684
12	50 to 54	176241	176440	352681	-2.55574	2.558627
13	55 to 59	153494	156283	309777	-2.22588	2.266322

Step 6: For analysis go to Insert → Bar → 2D



Step 7 : Drag the data and set the graph for analysis.



Conclusion: We have successfully performed the analysis for the given secondary data. We have opened the given data in excel and calculated the total sum of each column and applied formula for the respectively. Then we selected the graph to be a bar graph in 2D for the modified data and the output is obtained.

Practical 3

A.Perform testing of hypothesis using one sample t-test.

Code:

```
from scipy.stats import ttest_1samp

import numpy as np

ages=np.genfromtxt('H:/ages.csv')

print(ages)

ages_mean=np.mean(ages)

print(ages_mean)

tset,pval=ttest_1samp(ages,30)

print('p-values-',pval)

if pval<0.05:#alpha value is 0.05

    print("we are rejecting null hypothesis")

else:

    print("we are accepeting null hypothesis")
```

Output:

```
In [3]: runfile('H:/3a.py', wdir='H:')
[20. 30. 25. 13. 16. 17. 34. 35. 38. 42. 43. 45. 48. 49. 50. 51. 54. 55.
 56. 59. 61. 62. 18. 22. 29. 30. 31. 39. 52. 53. 67. 36. 47. 54. 40. 40.
 35. 22. 59. 58. 30. 43. 22. 45. 21. 59. 51. 47. 25. 58. 50. 23. 24. 45.
 37. 59. 28. 28. 48. 42. 54. 36. 36. 24. 26. 24. 50. 48. 34. 44. 56. 55.
 35. 33. 39. 53. 34. 28. 56. 24. 21. 29. 28. 58. 35. 57. 26. 25. 59. 56.
 22. 57. 48. 33. 23. 26. 57. 32. 53. 31. 35. 44. 54. 25. 31. 58. 26. 32.
 26. 50. 41. 49. 26. 33. 34. 24. 43. 42. 51. 36. 38. 38. 40. 38. 56. 39.
 23. 33. 53. 30. 38.]
39.47328244274809
p-values- 5.362905195437013e-14
we are rejecting null hypothesis
```

Conclusion: we have successfully performed testing of hypothesis using one sample t-test. The one-sample t-test is a statistical hypothesis test used to determine whether an unknown population mean is different from a specific value.

B. Write a program for t-test comparing two means for independent samples.

Step 1: Calculate the mean of the samples (Total /Count)

SUM				
	A	B	C	D
1			Female	Male
2			26	23
3			25	30
4			43	18
5			34	25
6			18	28
7			52	0
8		Total	198	124
9		Mean	=C8/7	
10				

D10				
	A	B	C	D
1			Female	Male
2			26	23
3			25	30
4			43	18
5			34	25
6			18	28
7			52	0
8		Total	198	124
9		Mean	28.28571	17.71429
10				

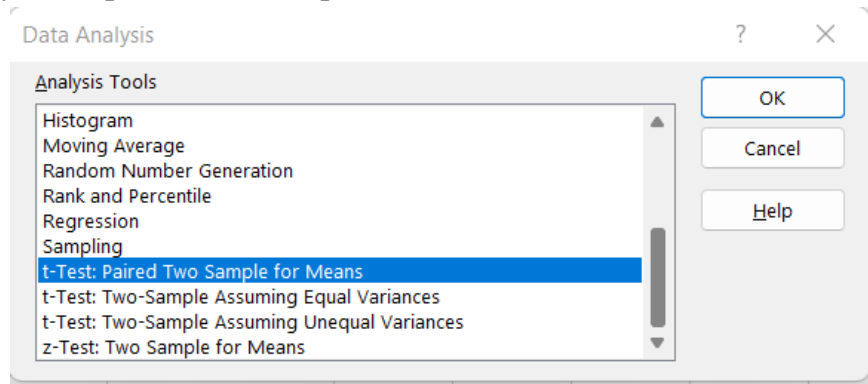
Step 2: Calculate the SD of samples STDEV(range)

C2			
	A	B	C
1			Female
2			26
3			25
4			43
5			34
6			18
7			52
8		Total	198
9		Mean	28.28571
10		SD	=STDEV(C2:C7)

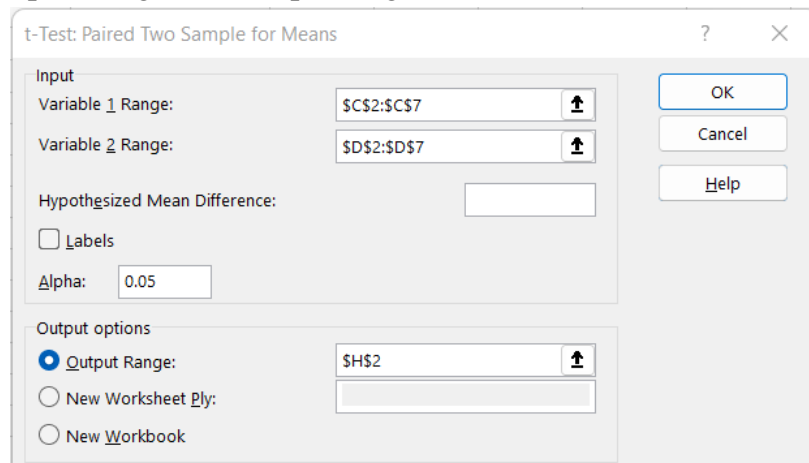
Step 3: Calculate the t-test

=(C9-D9)/SQRT((C10*C10)/COUNT(C2:C9)+(D10*D10)/COUNT(D2:D6))										
	A	B	C	D	E	F	G	H	I	J
1			Female	Male						
2			26	23						
3			25	30						
4			43	18						
5			34	25						
6			18	28						
7			52	0						
8		Total	198	124						
9		Mean	28.28571	17.71429						
10		SD	12.64911	10.94836						
11										

Step 4: Apply t-test paired two Samples for Means.



Step 5: Apply input Range and Output range.



Step 6: Hence we reject null hypothesis.

H15	t Critical two-tail										
	A	B	C	D	E	F	G	H	I	J	K
1			Female	Male							
2			26	23				t-Test: Paired Two Sample for Means			
3			25	30		Calculated T test					
4			43	18		1.594185981					
5			34	25				Mean	33	20.66666667	
6			18	28				Variance	160	119.8666667	
7			52	0				Observations	6		6
8		Total	198	124				Pearson Correlation	-0.889613782		
9		Mean	28.28571	17.71429				Hypothesized Mean Difference	0		
10		SD	12.64911	10.94836				df	5		
11								t Stat	1.316901108		
12								P(T<=t) one-tail	0.122498157		
13								t Critical one-tail	2.015048373		
14								P(T<=t) two-tail	0.244996315		
15								t Critical two-tail	2.570581836		

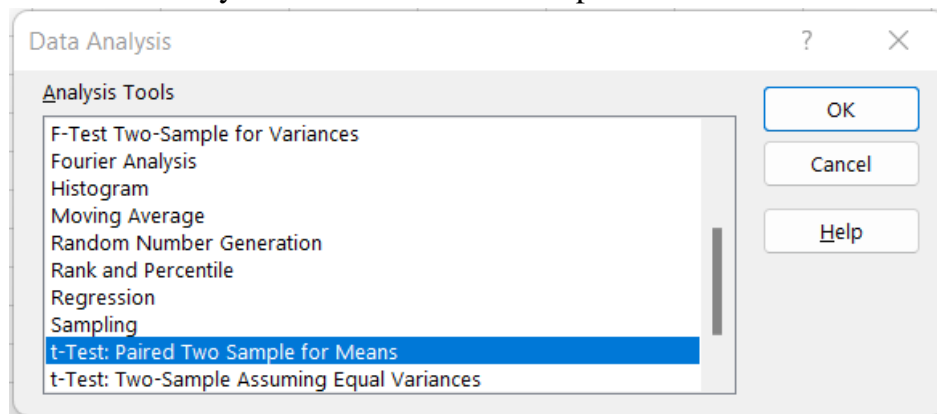
Conclusion: we have successfully compared the two means of independent samples for a t-test. The Independent Samples t Test compares the means of two independent groups in order to determine whether there is statistical evidence that the associated population means are significantly different.

C.Perform testing of Hypothesis using paired t-test.

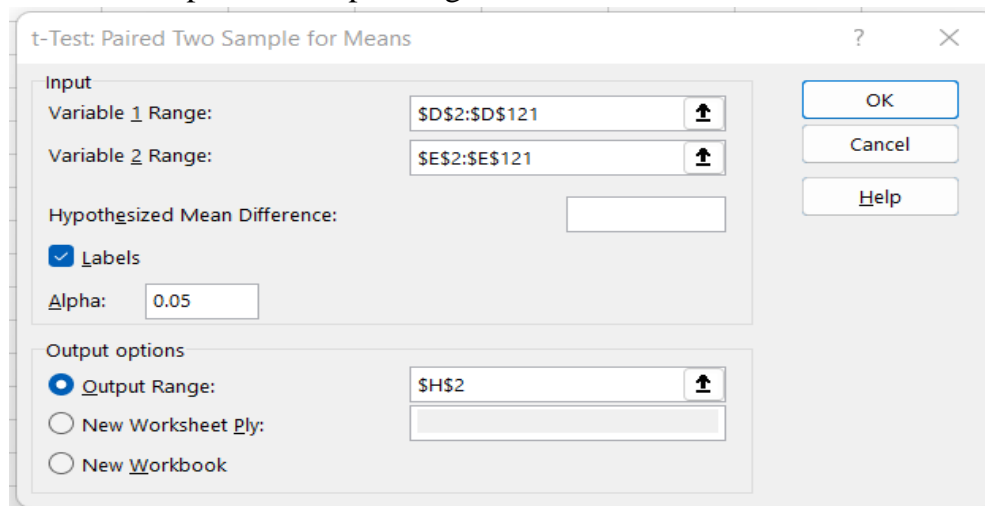
Step 1: Load the Data

	A	B	C	D	E	F
1	patient	gender	agegrp	bp_before	bp_after	
2	1	Male	30-45	143	153	
3	2	Male	30-45	163	170	
4	3	Male	30-45	153	168	
5	4	Male	30-45	153	142	
6	5	Male	30-45	146	141	
7	6	Male	30-45	150	147	
8	7	Male	30-45	148	133	
9	8	Male	30-45	153	141	
10	9	Male	30-45	153	131	
11	10	Male	30-45	158	125	
12	11	Male	30-45	149	164	
13	12	Male	30-45	173	159	

Step 2: Data -> Data Analysis -> t-test between samples for means.



Step 3: Choose the Input and Output range.



Step 4: Since the samples Means of the Data samples are not equal we reject the null hypothesis that they might be equal.

H21											
	A	B	C	D	E	F	G	H	I	J	K
1	patient	gender	agegrp	bp_before	bp_after						
2	1	Male	30-45	143	153			t-Test: Paired Two Sample for Means			
3	2	Male	30-45	163	170						
4	3	Male	30-45	153	168				143	153	
5	4	Male	30-45	153	142			Mean	156.5630252	151.3445378	
6	5	Male	30-45	146	141			Variance	129.2820111	202.6853725	
7	6	Male	30-45	150	147			Observations	119	119	
8	7	Male	30-45	148	133			Pearson Correlation	0.161241417		
9	8	Male	30-45	153	141			Hypothesized Mean Difference	0		
10	9	Male	30-45	153	131			df	118		
11	10	Male	30-45	158	125			t Stat	3.403463555		
12	11	Male	30-45	149	164			P(T<=t) one-tail	0.000454744		
13	12	Male	30-45	173	159			t Critical one-tail	1.657869522		
14	13	Male	30-45	165	135			P(T<=t) two-tail	0.000909488		
15	14	Male	30-45	145	159			t Critical two-tail	1.980272249		
16	15	Male	30-45	143	153						
17	16	Male	30-45	152	126						
18	17	Male	30-45	141	162						
19	18	Male	30-45	176	134						
20	19	Male	30-45	143	136						
21	20	Male	30-45	162	150						

Conclusion: We have performed the testing of hypothesis using paired t-test. A paired t-test is used when we are interested in the difference between two variables for the same subject. The Data is loaded and Data analysis for paired t-test is selected. Input and Output Range of the variables are entered (Labels=tick, Alpha =0.5). In Step 4, the output for the following is displayed.

Practical 4

A. Perform testing of hypothesis using chi-squared goodness-of-fit test.

Step 1: Load the data.

D2				
	A	B	C	D
1	O	E		
2	29	21.33		
3	24	21.33		
4	22	21.33		
5	19	21.33		
6	21	21.33		
7	18	21.33		
8	19	21.33		
9	20	21.33		
10	21	21.33		
11	18	21.33		
12	20	21.33		
13	23	21.33		
14				

Step 2: Type CHITEST and select Actual(observed value range) and Expected value range.

F5									
	A	B	C	D	E	F	G	H	
1	O	E	O-E	(O-E) ²					
2	29	21.33	7.67	58.8289					
3	24	21.33	2.67	7.1289					
4	22	21.33	0.67	0.4489					
5	19	21.33	-2.33	5.4289		=CHITEST(A2:A13,B2:B13)			
6	21	21.33	-0.33	0.1089					
7	18	21.33	-3.33	11.0889					
8	19	21.33	-2.33	5.4289					
9	20	21.33	-1.33	1.7689					
10	21	21.33	-0.33	0.1089					
11	18	21.33	-3.33	11.0889					
12	20	21.33	-1.33	1.7689					
13	23	21.33	1.67	2.7889					
14									

Step 3: This calculated value is less than table value which is 19.68. Hence we accept null hypothesis.

i.e $0.932663 < 19.68$. H_0 accepted.

	A	B	C	D	E	F
1	O	E	O-E	(O-E)^2		
2	29	21.33	7.67	58.8289		
3	24	21.33	2.67	7.1289		
4	22	21.33	0.67	0.4489		
5	19	21.33	-2.33	5.4289		0.932663
6	21	21.33	-0.33	0.1089		
7	18	21.33	-3.33	11.0889		
8	19	21.33	-2.33	5.4289		
9	20	21.33	-1.33	1.7689		
10	21	21.33	-0.33	0.1089		
11	18	21.33	-3.33	11.0889		
12	20	21.33	-1.33	1.7689		
13	23	21.33	1.67	2.7889		
14						

Conclusion: In the above practical, we have successfully performed Testing of hypothesis using Chi squared goodness of fit test using Excel. The calculated value of Chi-Square goodness of fit test is compared with the table value. If the calculated value of Chi-Square goodness of fit test is greater than the table value, we will reject the null hypothesis and conclude that there is a significant difference between the observed and the expected frequency.

B.Perform testing of hypothesis using chi-squared test of independence.

Step 1: Load the data

	A	B	C	D
1	System	O	E	
2	Windows	20	33.33	
3	Mac	60	33.33	
4	Linux	20	33.33	
5				

Step 2: Calculate the value of $(O-E)^2 / E$

SUM					
	A	B	C	D	E
1	System	O	E	$(O-E)^2 / E$	
2	Windows	20	33.33	$=(B2-C2)^2/C2$	
3	Mac	60	33.33		
4	Linux	20	33.33		
5					

Step 3: Calculate Chi square at 5% confidence and degree of freedom $n-1=2$ (in our case)

F2						
	A	B	C	D	E	F
1	System	O	E	$(O-E)^2 / E$		
2	Windows	20	33.33	5.33120012		5.991465
3	Mac	60	33.33	21.34080108		
4	Linux	20	33.33	5.33120012		
5			Total	32.00320132		

Conclusion: Using a chi square test, you've just determined that there is, in fact, a statistically significant relationship between our two categorical variables, s2q10 and s1truan. In addition, we've run another chi square, determining that there is a statistically significant relationship between s2q10 and s1q62a, a measure of whether or not a respondent's father had obtained a degree. Remember that you are simply able to say now that paternal degree and Year 11 truancy both have relationships with respondent enrolment in full time education after secondary school. We cannot say, for example, that a paternal degree causes enrolment in full time education.

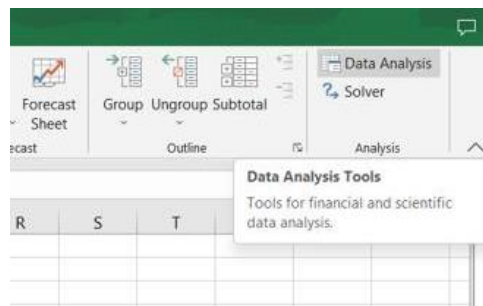
Practical 5

Perform testing of hypothesis using Z-test.

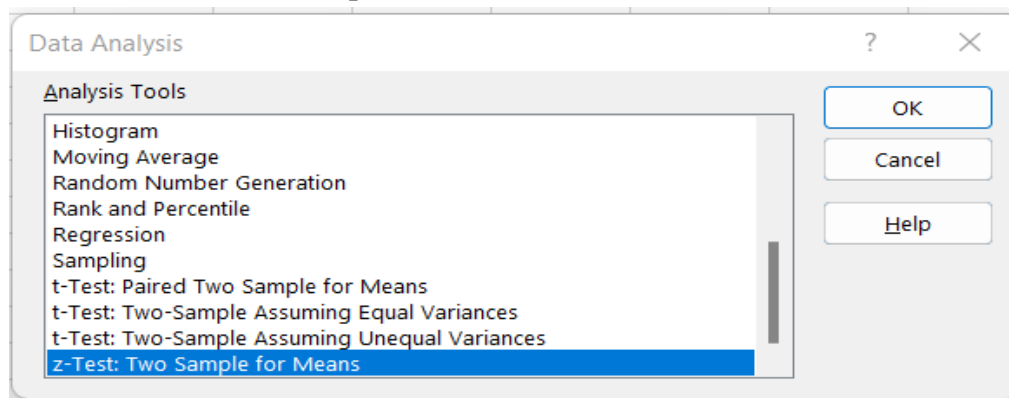
Step 1: Load the data

	A	B	C	D	E	F
1	patient	gender	agegrp	bp_before	bp_after	
2	1	Male	30-45	143	153	
3	2	Male	30-45	163	170	
4	3	Male	30-45	153	168	
5	4	Male	30-45	153	142	
6	5	Male	30-45	146	141	
7	6	Male	30-45	150	147	
8	7	Male	30-45	148	133	
9	8	Male	30-45	153	141	
10	9	Male	30-45	153	131	
11	10	Male	30-45	158	125	

Step 2: To apply Z test we need a sample size over 30. Here our sample size is 120 data points , so to apply Z-test go to Data -> Data analysis.



Step 3: Select Z-Test: Two Sample for means.



Step 5: Set the Variable 1 and 2 range.

	A	B	C	D	E	F	G	H
1	patient	gender	agegrp	bp_before	bp_after			
2	1	Male	30-45	143	153		129.7286	
3	2	Male	30-45	163	170			
4	3	Male	30-45	153	168		201.005	
5	4	Male	30-45	153	142			
6	5	Male	30-45	146	141			
7	6	Male	30-45	150	147			
8	7	Male	30-45	148	133			
9	8	Male	30-45	153	141			

Step 5: Set the Variable 1 and 2 range.

z-Test: Two Sample for Means

Input

Variable 1 Range:

Variable 2 Range:

Hypothesized Mean Difference:

Variable 1 Variance (known):

Variable 2 Variance (known):

☒ Labels

Alpha:

Output options

☒ Output Range:

☐ New Worksheet Ply:

☐ New Workbook

OK Cancel Help

Step 6: Output

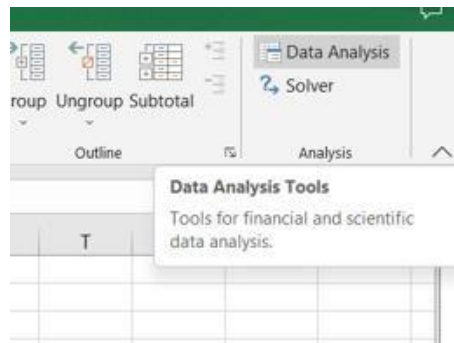
	A	B	C	D	E	F	G	H	I	J	K	L
1	patient	gender	agegrp	bp_before	bp_after							
2		1 Male	30-45	143	153		129.7286			z-Test: Two Sample for Means		
3		2 Male	30-45	163	170							
4		3 Male	30-45	153	168		201.005				143	153
5		4 Male	30-45	153	142					Mean	156.5630252	151.3445378
6		5 Male	30-45	146	141					Known Variance	129.7286	201.005
7		6 Male	30-45	150	147					Observations	119	119
8		7 Male	30-45	148	133					Hypothesized Mean Difference	0	
9		8 Male	30-45	153	141					z	3.13024954	
10		9 Male	30-45	153	131					P(Z<=z) one-tail	0.000873289	
11		10 Male	30-45	158	125					z Critical one-tail	1.644853627	
12		11 Male	30-45	149	164					P(Z<=z) two-tail	0.001746579	
13		12 Male	30-45	173	159					z Critical two-tail	1.959963985	
14		13 Male	30-45	165	135							

Conclusion: The pre-requisite for using one-sample Z-test for means is that data must belong to normal distribution and the population standard deviation is known. In case, the population standard deviation is unknown, t-test gets used. The z-test for hypothesis testing will help you make confident decisions about your data.

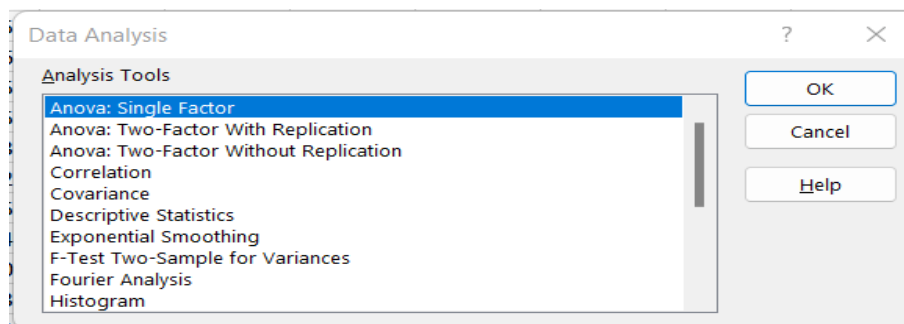
A. Perform testing of hypothesis using One-way ANOVA.

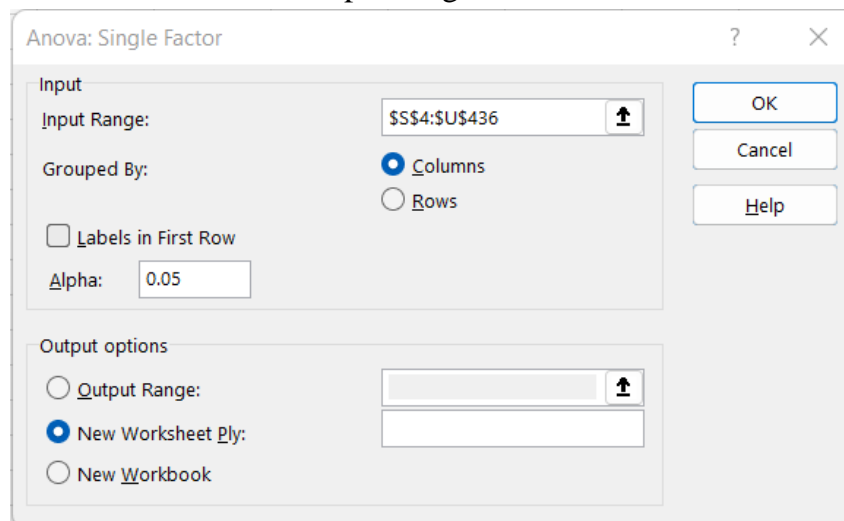
345 East 15th Street																						
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	
1	School IC	School Name	Borough	Buildin Str	Address	City	State	Zip Code	Latitude	Longitude	Phone Nu	Start Time	End Time	Student En	Prct W	Percent BI	Percent HI	Percent As	Percent S	Percent T	Percent U	Percent V
2	02M261	Clinton School W/Manhattan	M933	425 West 33rd St	Manhattan	NY	10001	40.75321	-73.9979	212-695-9114												
3	02M262	Lincoln School W/Manhattan	M932	425 West 33rd St	Manhattan	NY	10001	40.86605	-73.9249	718-935-3	8:30 AM	3:00 PM	87	3.40%	21.80%	67.80%	4.60%					
4	01M529	Univ Early Coll Manhattan	M022	111 Columbia Str	Manhattan	NY	10002	40.71873	-73.9794	212-677-5	8:15 AM	4:00 PM	1735	28.60%	13.30%	18.00%	38.50%	657	601	601	91.00%	
5	02M324	Essex Street Acad Manhattan	M445	350 Grand Street	Manhattan	NY	10002	40.71867	-73.9895	212-475-4	8:00 AM	2:45 PM	358	11.70%	38.20%	41.30%	5.90%	395	411	387	78.90%	
6	02M308	Lower Manhattan Manhattan	M445	350 Grand Street	Manhattan	NY	10002	40.71867	-73.9895	212-505-0	8:30 AM	3:00 PM	383	3.10%	28.00%	56.90%	8.60%	418	428	415	65.10%	
7	02M455	High School for D-Manhattan	M445	350 Grand Street	Manhattan	NY	10002	40.71867	-73.9895	212-475-4	8:00 AM	3:25 PM	541	1.70%	31.0%	58.90%	613	453	463	95.90%		
8	01M292	Henry Street Schc Manhattan	M056	220 Henry Street	Manhattan	NY	10002	41.71376	-73.9853	212-496-8	8:30 AM	3:30 PM	255	3.90%	24.40%	56.60%	13.20%	410	406	381	59.70%	
9	01M696	Bard High School Manhattan	M097	525 East Houston	Manhattan	NY	10002	40.71896	-73.9761	212-905-9	8:00 AM	3:50 PM	545	45.30%	17.20%	18.70%	17.10%	634	641	639	70.80%	
10	01M305	Urban Assembly f Manhattan	M445	350 Grand Street	Manhattan	NY	10002	40.71867	-73.9895	212-505-0	8:32 AM	3:45 PM	329	2.70%	41.90%	49.20%	5.80%	389	395	381	80.80%	
11	01M509	Marta Valle High Manhattan	M025	145 Stanton Street	Manhattan	NY	10002	40.72057	-73.9857	212-473-8	8:00 AM	3:30 PM	363	2.50%	39.00%	51.20%	5.80%	438	413	394	35.60%	
12	01M448	University Neighs Manhattan	M446	200 Monroe Street	Manhattan	NY	10002	40.71233	-73.9848	212-962-4	8:15 AM	3:15 PM	304	3.30%	25.00%	41.10%	29.90%	437	355	352	69.90%	
13	01M543	New Design High Manhattan	M445	350 Grand Street	Manhattan	NY	10002	40.71867	-73.9895	212-475-4	8:00 AM	2:56 PM	441	3.90%	30.80%	56.90%	5.90%	381	396	372	73.70%	
14	02M298	Pace High School Manhattan	M131	100 Hester Street	Manhattan	NY	10002	40.71641	-73.9927	212-344-4	9:00 AM	3:15 PM	423	1.90%	28.10%	45.40%	13.70%	430	435	427	87.90%	
15	02M420	High School for H-Manhattan	M475	345 East 15th Str	Manhattan	NY	10003	40.73249	-73.9831	212-780-9	9:00 AM	3:45 PM	1664	7.30%	18.90%	50.90%	22.40%	452	445	430	86.00%	
16	02M399	High School for L-Manhattan	M460	40 Irving Place	Manhattan	NY	10003	40.73552	-73.9876	212-253-2	8:00 AM	3:30 PM	437	5.70%	20.40%	40.30%	31.10%	446	433	411	70.20%	
17	02M546	Academy for Soft Manhattan	M460	40 Irving Place	Manhattan	NY	10003	40.73552	-73.9876	212-253-3	8:45 AM	3:36 PM	344	9.00%	28.80%	45.90%	11.00%					
18	02M533	Union Square Aca Manhattan	M460	40 Irving Place	Manhattan	NY	10003	40.73552	-73.9876	212-253-3	8:00 AM	3:49 PM	319	4.70%	20.10%	56.70%	16.00%					
19	02M438	International Higl Manhattan	M460	40 Irving Place	Manhattan	NY	10003	40.73552	-73.9876	212-533-2	8:45 AM	3:05 PM	353	9.90%	14.20%	45.30%	30.60%	403	330	316	53.20%	
20	02M407	Institute for Colla Manhattan	M475	345 East 15th Str	Manhattan	NY	10003	40.7324														

Data -> Data Analysis.



Step 3: Select Anova: Single Factor



Step 4: Select the cells of S-T-U for input range


The image shows the 'Anova: Single Factor' dialog box in Excel. The 'Input' section has 'Input Range' set to '\$S\$4:\$U\$436' and 'Grouped By' set to 'Columns'. The 'Labels in First Row' checkbox is unchecked, and 'Alpha' is set to '0.05'. The 'Output options' section has 'New Worksheet Ply' selected. Buttons for 'OK', 'Cancel', and 'Help' are on the right.

Step 5: Output

	A	B	C	D	E	F	G
1	Anova: Single Factor						
2							
3	SUMMARY						
4	Groups	Count	Sum	Average	Variance		
5	Column 1	375	162354	432.944	5177.144		
6	Column 2	375	159189	424.504	3829.267		
7	Column 3	375	156922	418.4587	4166.522		
8							
9							
10	ANOVA						
11	Source of Variation	SS	df	MS	F	P-value	F crit
12	Between Groups	39700.57	2	19850.28	4.520698	0.01108	3.003745
13	Within Groups	4926677	1122	4390.977			
14							
15	Total	4966377	1124				
16							

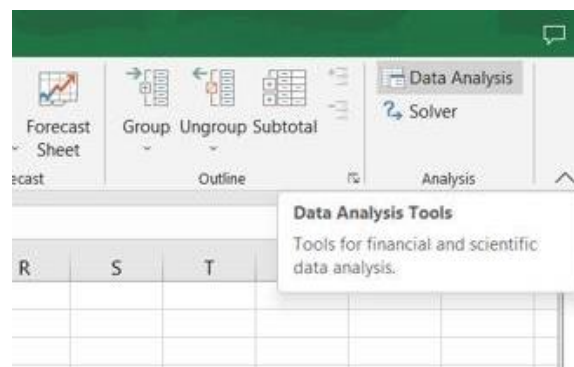
Conclusion: In the above practical, we have successfully performed One-way ANOVA on 'scores.csv' dataset. Since, in the resulting output the p-value is less than 0.05 so we are rejecting Null Hypothesis and we conclude that there is a significant difference between the SAT scores for each group mean.

B. Perform testing of hypothesis using Two-way ANOVA.

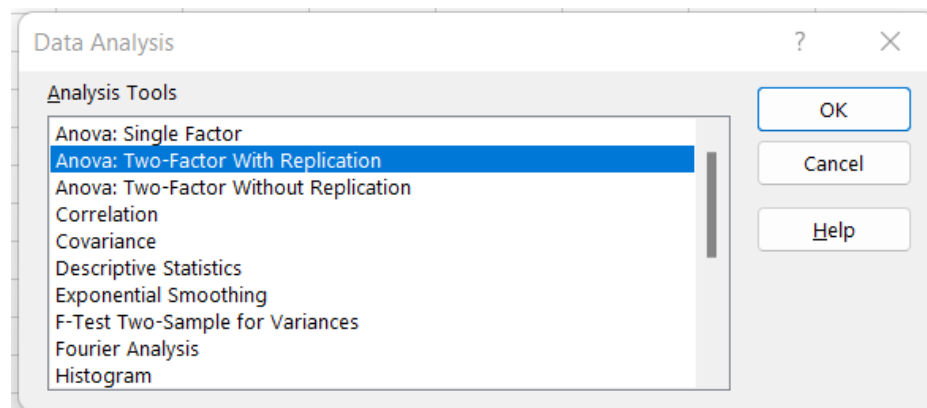
Step 1: Load the data.

	A	B	C	D	E	F
1		len	supp	dose		
2	1	4.2	VC	0.5		
3	2	11.5	VC	0.5		
4	3	7.3	VC	0.5		
5	4	5.8	VC	0.5		
6	5	6.4	VC	0.5		
7	6	10	VC	0.5		
8	7	11.2	VC	0.5		
9	8	11.2	VC	0.5		
10	9	5.2	VC	0.5		
11	10	7	VC	0.5		
12	11	16.5	VC	1		
13	12	16.5	VC	1		
14	13	15.2	VC	1		

Step 2: Data -> Data Analysis



Step 3: Select Anova: Two Factor with Replication



Step 4: Select input and output range

Anova: Two-Factor With Replication

Input

Input Range:

Rows per sample:

Alpha:

Output options

☒ Output Range:

☐ New Worksheet Ply:

☐ New Workbook

OK Cancel Help

Step 5: Output

Anova: Two-Factor With Replication						
SUMMARY	len	dose	Total			
1						
Count	30	30	60			
Sum	508.9	35	543.9			
Average	16.9633	1.16667	9.065			
Variance	68.3272	0.4023	97.2233			
31						
Count	30	30	60			
Sum	619.9	35	654.9			
Average	20.6633	1.16667	10.915			
Variance	43.6334	0.4023	118.285			
Total						
Count	60	60				
Sum	1128.8	70				
Average	18.8133	1.16667				
Variance	58.512	0.39548				
ANOVA						
Source of Variation	SS	df	MS	F	P-value	F crit
Sample	102.675	1	102.675	3.64208	0.05881	3.92288
Columns	9342.15	1	9342.15	331.384	8.5E-36	3.92288
Interaction	102.675	1	102.675	3.64208	0.05881	3.92288
Within	3270.19	116	28.1913			
Total	12817.7	119				

Conclusion: In the above practical, we have successfully performed Two-way ANOVA. Since, P-value = 0.0588 column in the ANOVA Source of Variation table at the bottom of the output. Because the p-values for both medicine dose and interaction are less than our significance level, these factors are statistically significant. On the other hand, the interaction effect is not significant because its p-value (0.05881) is greater than our

significance level. Because the interaction effect is not significant, we can focus on only the main effects and not consider the interaction effect of the dose.

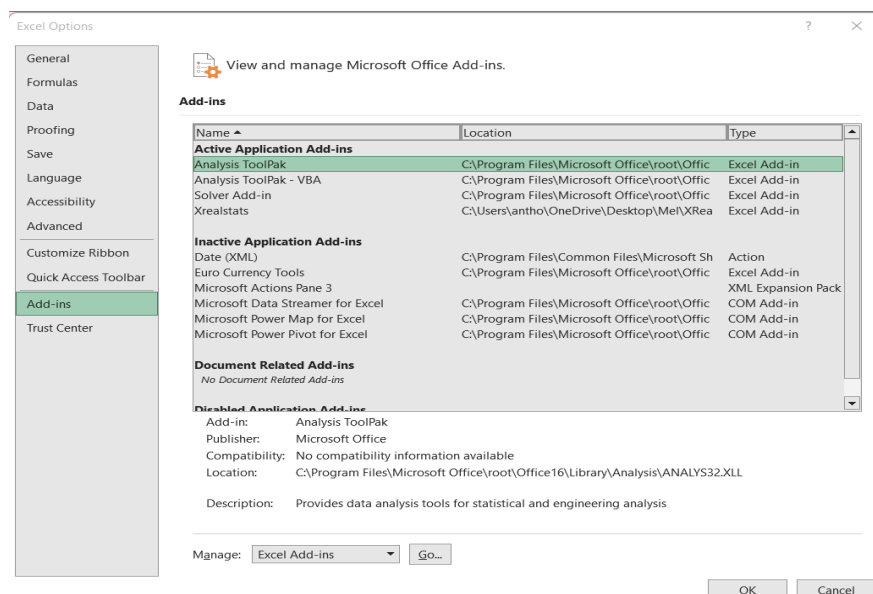
C.Perform testing of hypothesis using MANOVA.

Step 1: Load the Data.

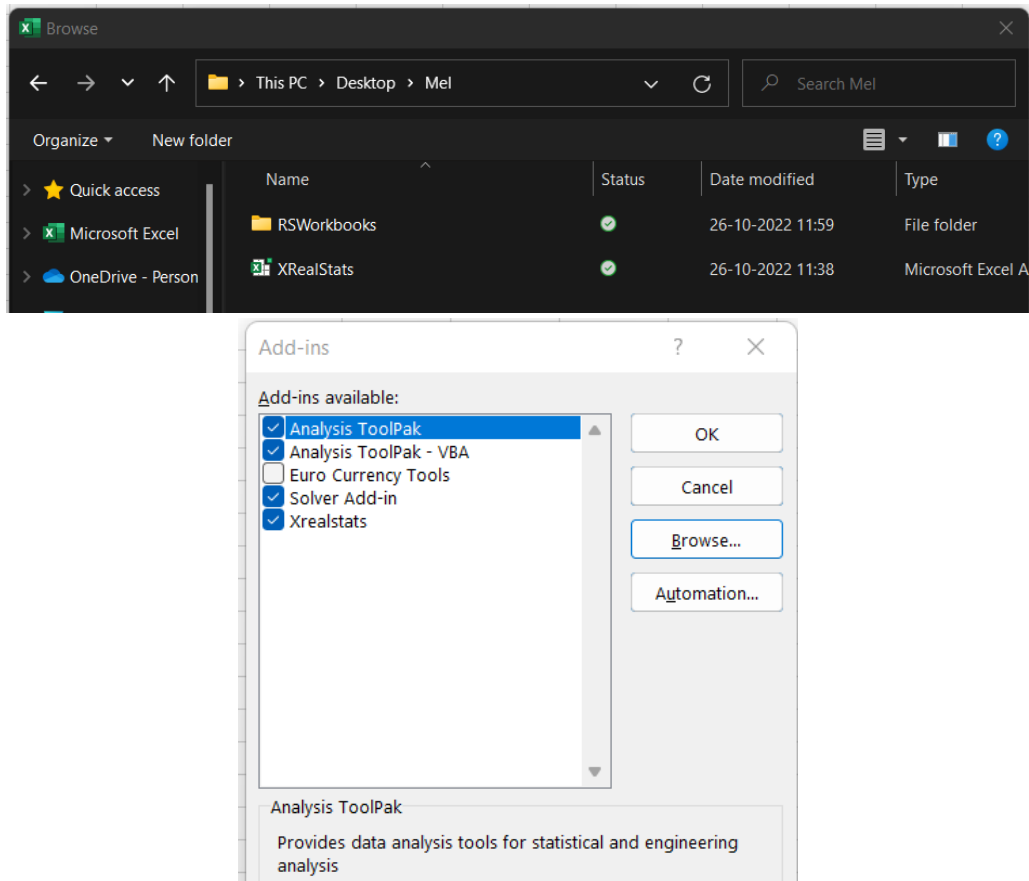
	A	B	C	D	E
1	Gender	Economics	Kindness	Optimism	
2	male	wealthy	5	3	
3	male	wealthy	4	6	
4	male	wealthy	3	4	
5	male	wealthy	2	4	
6	male	wealthy	4	6	
7	male	wealthy	3	6	
8	male	middle	5	4	
9	male	middle	5	5	
10	male	middle	7	5	

Step 2: Install Add-in in excel.

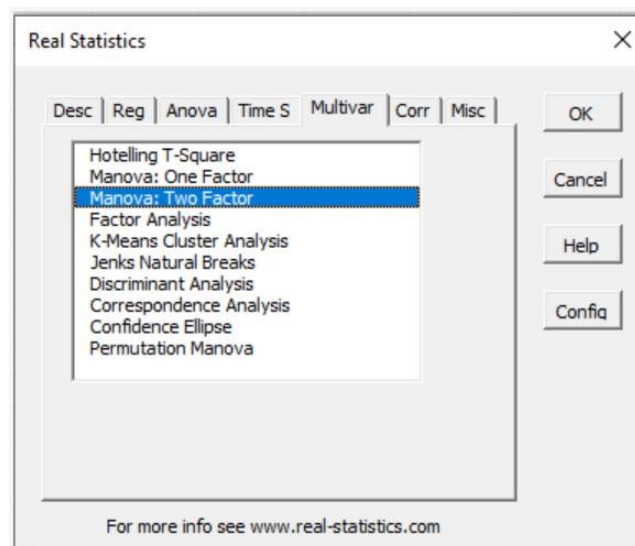
Select **File -> Help | Options -> Add-Ins** and click on the **Go** button at the bottom of the window.



Step 3: Click on browse and select XrealStats file. Check the following check boxes.



Step 4: Once loaded , CTRL+M to perform MANOVA. Click Multivar -> Manova: two factor -> OK.



Step 5: Select the data excluding column names. Select a cell for output

Manova: Two Factors

Input Range: Sheet1!\$A\$2:\$D\$27

Analysis type: ☒ Regular ☐ Repeated Measures

Options:

- ☒ Significance Analysis
- ☒ Sum of Squares and Cross Product Matrices
- ☒ Covariance Matrices
- ☒ Outliers
- ☒ Box's Test
- ☒ Group Means
- ☐ Contrast

Alpha: 0.05

Output Range: H6

Step 6: Output

Two-Way MANOVA							SSCP Matrices	
fact A	stat	df1	df2	F	p-value	part eta-sq	Tot	
Pillai Trace	0.214135	2	18	2.452347	0.114323	0.214135	130	64
Wilk's Lam	0.785865	2	18	2.452347	0.114323	0.214135	64	142
Hotelling	0.272483	2	18	2.452347	0.114323	0.214135		
Roy's Lg R	0.272483							
							Row (A)	
							27.08333	17.5
							17.5	11.30769
fact B	stat	df1	df2	F	p-value	part eta-sq	Column (B)	
Pillai Trace	0.339044	4	38	1.939195	0.123777	0.169522		
Wilk's Lam	0.820428	4	36	1.969888	0.120034	0.179572		
Hotelling	0.467615	4	34	1.987365	0.118648	0.189501	31.90278	28.18056
Roy's Lg R	0.425176						28.18056	30.36111
fact AB	stat	df1	df2	F	p-value	part eta-sq	Interaction (AB)	
Pillai Trace	0.249126	4	38	1.35172	0.268831	0.124563	5.963889	0.019444
Wilk's Lam	0.871554	4	36	1.326381	0.278864	0.128446	0.019444	21.0312
Hotelling	0.304976	4	34	1.296146	0.291053	0.132312		
Roy's Lg R	0.260917							
							Res	
							65.05	18.3
							18.3	79.3

Conclusion: In the above practical, we have successfully performed two-way MANOVA. Since, the p-value of fact A is 0.114323 and fact B is 0.118648 as we can see they are different so we are rejecting null Hypothesis as there is significant difference between mean SAT scores of two groups.

Practical 7

A. Perform the Random sampling for the given data and analyze it.

Step 1: Load the data.

	A	B	C	D	E	F	G	H	I	J	K
1	Sr. No.	Roll No.	Name	Gender	Grade		Sr. No.	Roll No.	Name	Gender	Grade
2	1	1	xyz	Male	O		24	3	add	Female	A
3	2	2	mno	Male	A		25	7	sd	Female	A
4	3	5	pqr	Male	O		26	9	hg	Female	C
5	4	13	abc	Male	B		27	11	dg	Female	D
6	5	16	xyz	Male	O		28	14	pgu	Female	F
7	6	17	mno	Male	O		29	25	yt	Female	O
8	7	34	pqr	Male	D		30	36	f	Female	C
9	8	35	abc	Male	O		31	40	gh	Female	C
10	9	38	xyz	Male	O		32	41	t	Female	B
11	10	42	mno	Male	C		33	46	es	Female	A
12	11	43	pqr	Male	O		34	47	gg	Female	B
13	12	45	abc	Male	O		35	10	gfh	Female	D
14	13	48	xyz	Male	B		36	20	ghy	Female	C
15	14	49	mno	Male	O		37	21	tg	Female	B
16	15	50	pqr	Male	D		38	72	sfe	Female	O
17	16	51	abc	Male	O		39	73	sfg	Female	D
18	17	54	xyz	Male	A		40	75	dgt	Female	A
19	18	55	mno	Male	B		41	77	dgt	Female	B
20	19	56	pqr	Male	D		42	82	jgt	Female	D
21	20	59	abc	Male	O		43	84	ghh	Female	D
22	21	61	xyz	Male	F		44	91	fer	Female	D
23	22	62	mno	Male	A		45	95	gyd	Female	D
24	23	63	pqr	Male	C		46	4	gtt	Female	C

Step 2: Set Cell O1= Male and Cell P2= Female

To generate a random sample for male students from given population go to Cell O1 and type =INDEX(E\$2:E\$62,RANK(B2,B\$2:B\$62))

	H	I	J	K	L	M	N	O	P	Q	R
1	Roll No.	Name	Gender	Grade				=INDEX(E\$2:E\$62,RANK(B2,B\$2:B\$62))			
2	3	add	Female	A							
3	7	sd	Female	A							
4	9	hg	Female	C							
5	11	dg	Female	D							
6	14	pgu	Female	F							
7	25	yt	Female	O							
8	36	f	Female	C							
9	40	gh	Female	C							
10	41	t	Female	B							
11	46	es	Female	A							
12	47	gg	Female	B							
13	10	gfh	Female	D							
14	30	abc	Female	C							

B.Perform the Stratified sampling for the given data and analyse it

Step 1: Load the data.

	A	B	C	D	E
1	Cost	Product	Status	Random	
2	485	book	new		
3	697	bag	old		
4	225	bag	new		
5	256	bag	old		
6	562	book	old		
7	743	book	old		
8	985	bag	old		
9	321	bag	new		
10	954	book	new		

Step 2: Assign Random values using RAND() function.

	A	B	C	D	E
1	Cost	Product	Status	Random	
2	485	book	new	0.250817	
3	697	bag	old		

Step 3: Copy the entire column D and paste only values

C	D	E	F	G	H	I	J	K
Status	Random							
new	0.27945							
old	0.296495							
new	0.471366							
old	0.185464							
old	0.851319							
old	0.312694							
old	0.823398							
new	0.969015							
new	0.969434							
old	0.194029							
old	0.495866							
new	0.060707							
new	0.132045							
new	0.910547							
new	0.017235							
old	0.710536							
old	0.851558							
old	0.941389							

Paste Special

Paste

☐ All
 ☐ Formulas
 ☒ Values
 ☐ Formats
 ☐ Comments
 ☐ Validation

☐ All using Source theme
 ☐ All except borders
 ☐ Column widths
 ☐ Formulas and number formats
 ☐ Values and number formats
 ☐ All merging conditional formats

Operation

☒ None
 ☐ Add
 ☐ Subtract

☐ Multiply
 ☐ Divide

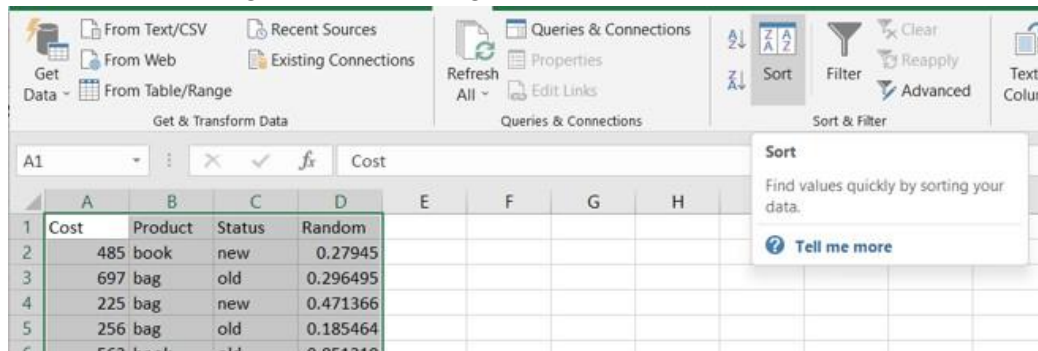
☐ Skip blanks
 ☐ Transpose

Paste Link

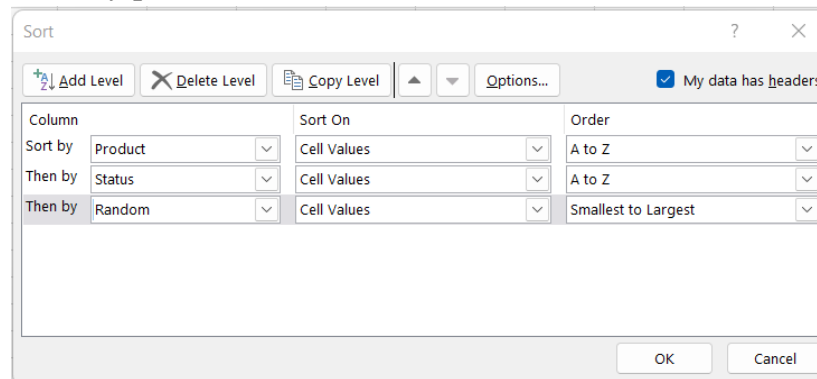
OK

Cancel

Step 4: To perform stratified sampling and obtain data where we get 2 new books cost, 2 old books cost, 2 new bags cost, 2 old bags cost. Select the data and click on sort.



Step 5: Sort the data by product , status and random.



Step 6: the output obtained will be as follows.

	A	B	C	D	E
1	Cost	Product	Status	Random	
2	125	bag	new	0.017235	
3	225	bag	new	0.471366	
4	321	bag	new	0.969015	
5	256	bag	old	0.185464	
6	870	bag	old	0.194029	
7	697	bag	old	0.296495	
8	985	bag	old	0.823398	
9	985	bag	old	0.851558	
10	415	bag	old	0.941389	
11	632	book	new	0.060707	
12	416	book	new	0.132045	
13	485	book	new	0.27945	
14	857	book	new	0.910547	
15	954	book	new	0.969434	
16	743	book	old	0.312694	

Conclusion: Conclusion: In stratified sampling, researchers divide subjects into subgroups called strata based on characteristics that they share (e.g., race, gender, educational attainment). Once divided, each subgroup is randomly sampled using another probability sampling method. We have successfully performed stratified sampling for the given data and its analysis. We have opened the data in excel and assigned random values using RAND() function followed by data sorting based on the columns and the final output is obtained.

Practical 8

Write a program for computing different correlation

A. Positive Correlation.

Code:

```
import matplotlib

import numpy as np

import matplotlib.pyplot as plt

np.random.seed(1)

# 1000 random integers between 0 and 50

x = np.random.randint(0, 50, 1000)

# Positive Correlation with some noise

y = x + np.random.normal(0, 10, 1000)

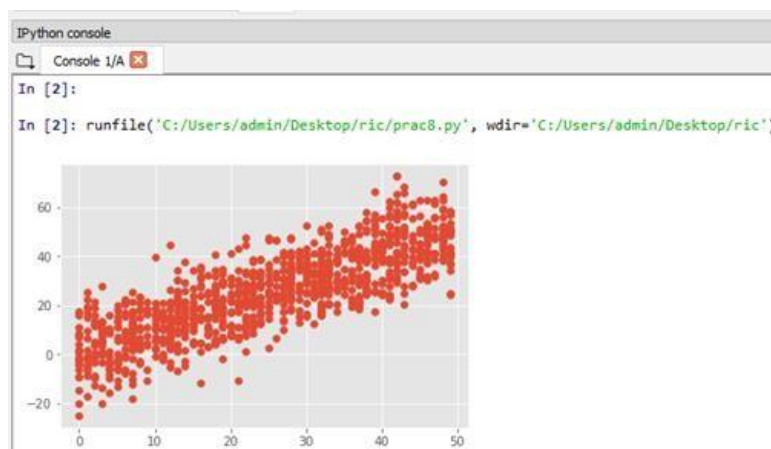
np.corrcoef(x, y)

matplotlib.style.use('ggplot')

plt.scatter(x, y)

plt.show()
```

Output:



Conclusion: A positive correlation is a relationship between two variables that move in tandem—that is, in the same direction. A positive correlation exists when one variable decreases as the other variable decreases, or one variable increases while the other increases. we have successfully computed positive correlation.

B. Negative Correlation.

Code:

```
import matplotlib

import numpy as np

import matplotlib.pyplot as plt

np.random.seed(1)

# 1000 random integers between 0 and 50

x = np.random.randint(0, 50, 1000)

# Negative Correlation with some noise

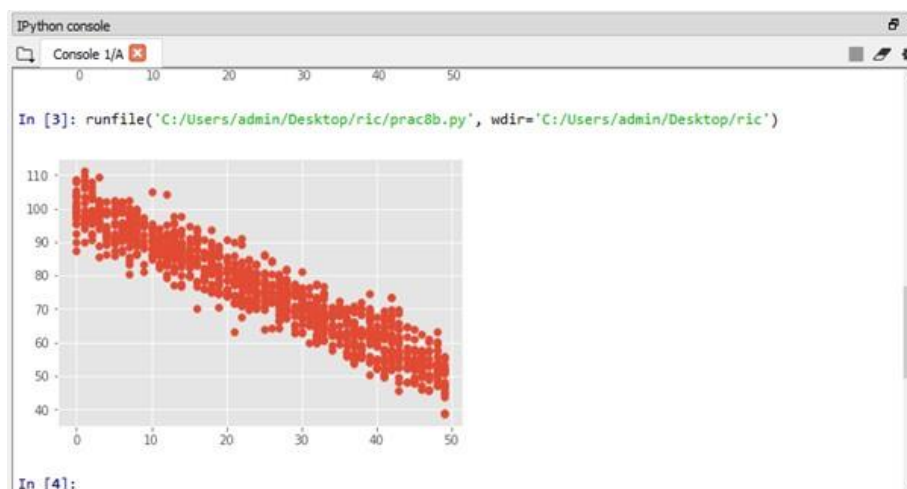
y = 100 - x + np.random.normal(0, 5, 1000)

np.corrcoef(x, y)

plt.scatter(x, y)

plt.show()
```

Output:



Conclusion: we have successfully computed negative correlation. A negative correlation is a relationship between two variables such that as the value of one variable increases, the other decreases. Correlation is expressed on a range from +1 to -1, known as the correlation coefficient. Values below zero express negative correlation. A perfect negative correlation has a coefficient of -1, indicating that an increase in one variable reliably predicts a decrease in the other one.

C.No/Weak Correlation.

Code:

```
import numpy as np

import matplotlib.pyplot as plt

np.random.seed(1)

x = np.random.randint(0, 50, 1000)

y = np.random.randint(0, 50, 1000)

np.corrcoef(x, y)

plt.scatter(x, y)

plt.show()
```

Output:



Conclusion: we have successfully computed no/weak correlation. A weak positive correlation indicates that, although both variables tend to go up in response to one another, the relationship is not very strong. A strong negative correlation, on the other hand, indicates a strong connection between the two variables, but that one goes up whenever the other one goes down.

Practical 9

A. Write a program to Perform linear regression for prediction.

Code:

```
> #Perform linear regression

> m<-c(1,2,3,4,5,6)

> t<-c(25,22,30,34,45,52)

> #Label the chart

> png(file="Linear Regression")

> plot(t,m,col="red",main="Months and
Temperature",abline(lm(m~t)),cex=1.6,pch=10,xlab="Months",ylab="
Temperature")

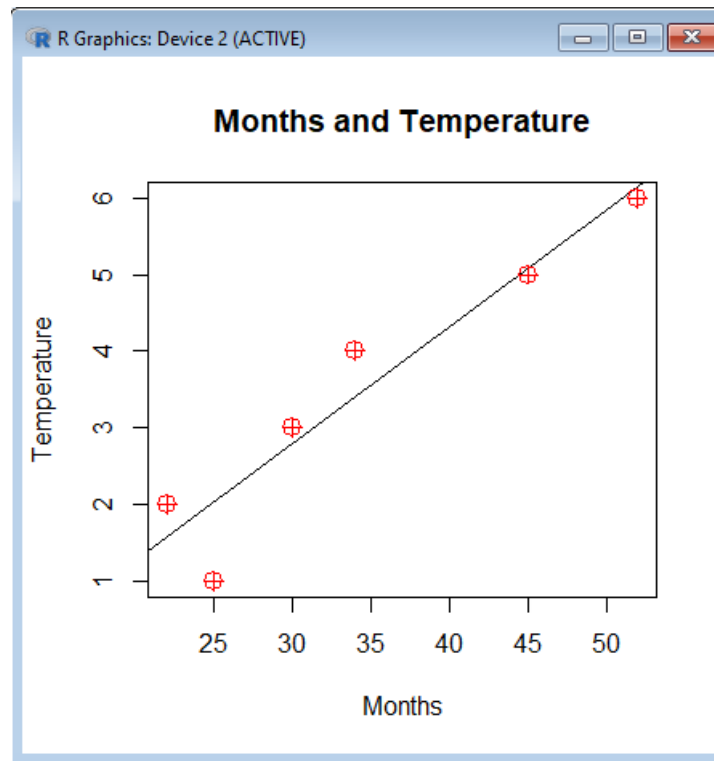
> dev.off()

null device

      1

> plot(t,m,col="red",main="Months and
Temperature",abline(lm(m~t)),cex=1.6,pch=10,xlab="Months",ylab="
Temperature")
```

```
> #Perform linear regression
> m<-c(1,2,3,4,5,6)
> t<-c(25,22,30,34,45,52)
> #Label the chart
> png(file="Linear Regression")
> plot(t,m,col="red",main="Months and Temperature",abline(lm(m~t)),cex=1.6,pch=10)
> dev.off()
null device
      1
> plot(t,m,col="red",main="Months and Temperature",abline(lm(m~t)),cex=1.6,pch=10)
> |
```

Conclusion: we have successfully performed linear regression for predicted. Linear regression analysis is used to predict the value of a variable based on the value of another variable. The variable you want to predict is called the dependent variable. The variable you are using to predict the other variable's value is called the independent variable. We have successfully performed linear regression for prediction in R. We have achieved this by defining dataset values for the variables month "m" and temperature "t". We have also labelled the graph as "Months and Temperature" and a precisely jotted straight line graph is obtained.

B.Polynomial Regression.

Code:

```
> #Polynomial Regression
> set.seed(16)
> x<-0:50
> y<-2.3-15.1*x+1.2*x^2+rnorm(length(x),20,50)
```

```
> plot(x,y)

> fit <- lm(y ~ 1 + x + I(x^2))

> points(x, predict(fit), type="l")

> summary(fit)
```

Output:

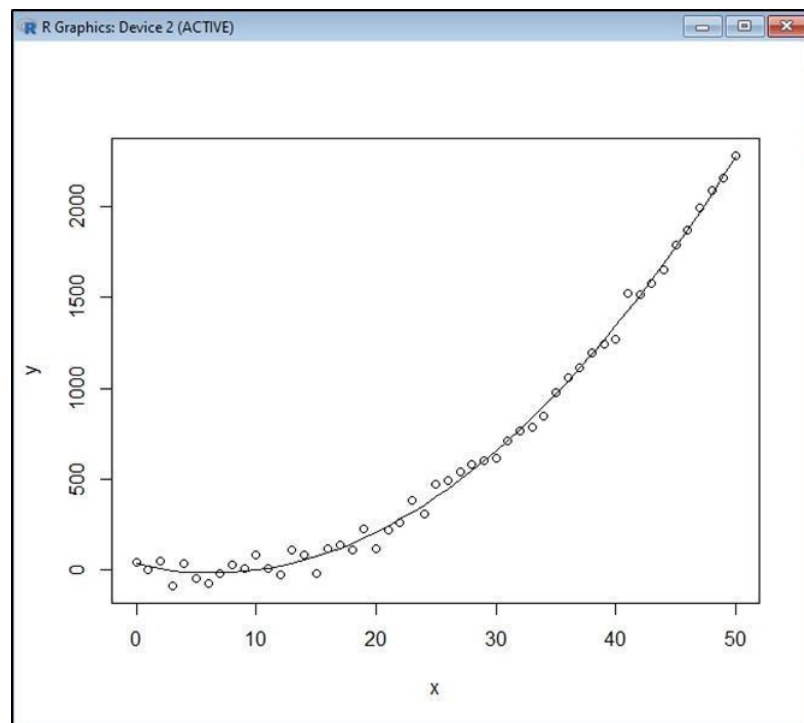
```
Call:
lm(formula = y ~ 1 + x + I(x^2))

Residuals:
    Min       1Q   Median       3Q      Max
-92.173 -28.968   3.673  24.953  97.269

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  33.84216   18.36178   1.843  0.0715 .
x          -15.28705    1.69836  -9.001 7.07e-12 ***
I(x^2)       1.20126    0.03285  36.569 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 45.44 on 48 degrees of freedom
Multiple R-squared:  0.996,    Adjusted R-squared:  0.9959
F-statistic: 6034 on 2 and 48 DF,  p-value: < 2.2e-16

> |
```



Conclusion:

Polynomial regression is a form of regression analysis in which the relationship between the independent variable x and the dependent variable y is modelled as an n th degree polynomial in x . We have successfully performed polynomial regression in R. We have achieved this by defining dataset values for the variables “ x ” and “ y ”. Followed by which we have fit a polynomial regression model to get predicted values and a curved scatterplot graph is obtained.

Practical 10

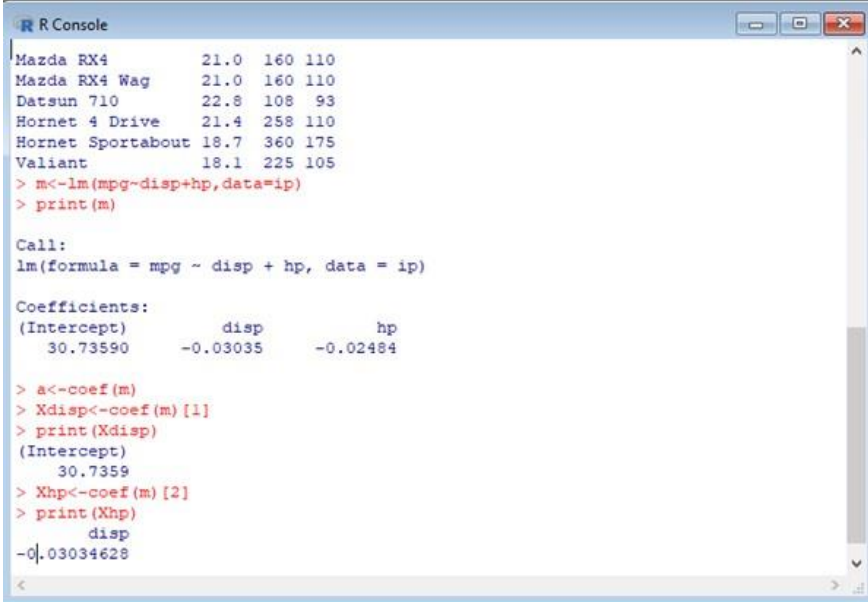
A. Multiple linear regression.

Code:

```
> #Multiple linear regression
> ip<-mtcars[,c("mpg", "disp", "hp")]
> print(head(ip))
> m<-lm(mpg~disp+hp, data=ip)
> print(m)
> a<-coef(m)
> Xdisp<-coef(m)[1]
> print(Xdisp)
> Xhp<-coef(m)[2]
> print(Xhp)
```

Output:

```
> print(head(ip))
      am    wt  mpg
Mazda RX4      1 2.620 21.0
Mazda RX4 Wag  1 2.875 21.0
Datsun 710      1 2.320 22.8
Hornet 4 Drive  0 3.215 21.4
Hornet Sportabout 0 3.440 18.7
Valiant        0 3.460 18.1
```



```

R Console
Mazda RX4           21.0  160 110
Mazda RX4 Wag       21.0  160 110
Datsun 710           22.8  108  93
Hornet 4 Drive       21.4  258 110
Hornet Sportabout    18.7  360 175
Valiant              18.1  225 105
> m<-lm(mpg~disp+hp,data=ip)
> print(m)

Call:
lm(formula = mpg ~ disp + hp, data = ip)

Coefficients:
(Intercept)          disp           hp
   30.73590      -0.03035     -0.02484

> a<-coef(m)
> Xdisp<-coef(m)[1]
> print(Xdisp)
(Intercept)
   30.7359
> Xhp<-coef(m)[2]
> print(Xhp)
          disp
-0.03034628

```

Conclusion: Multiple linear regression is a generalization of simple linear regression, in the sense that this approach makes it possible to relate one variable with several variables through a linear function in its parameters. We have successfully performed multiple linear regression.

B. Logistic Regression.

Code:

```

> #Logistic regression

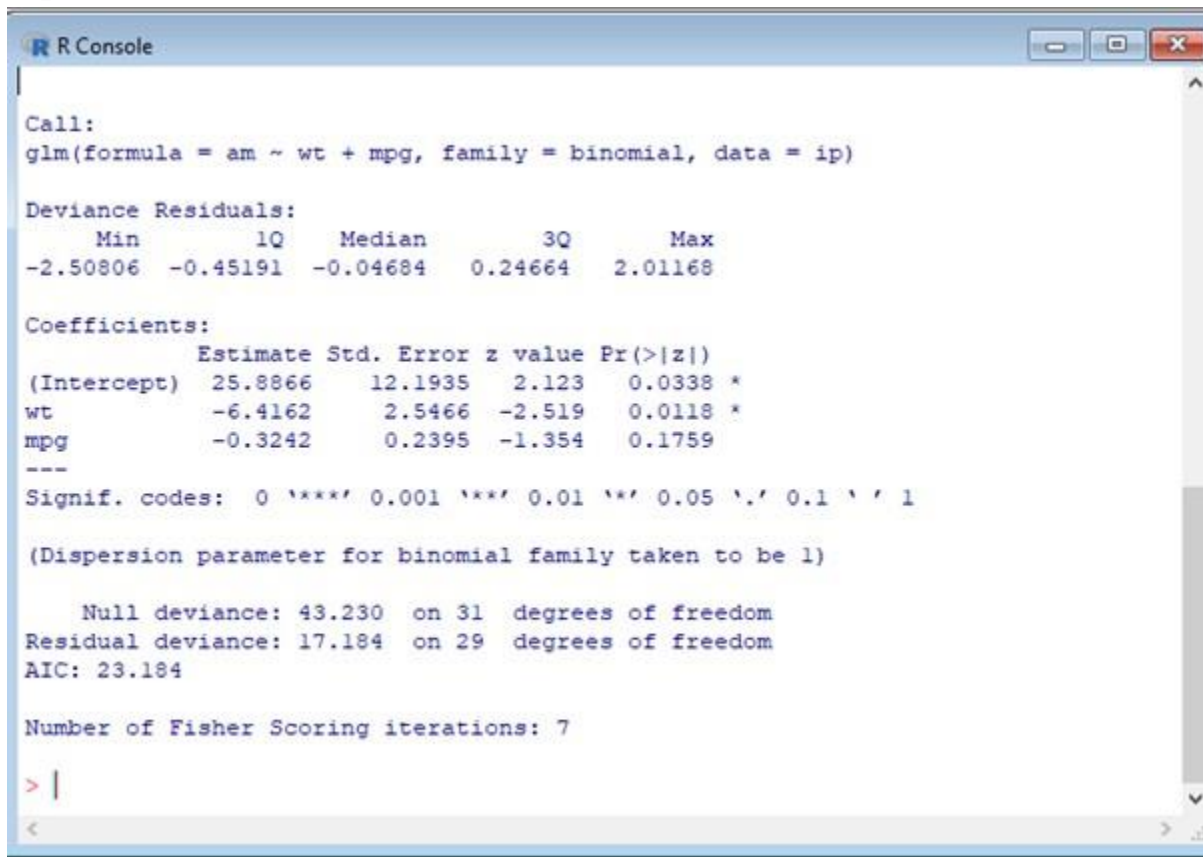
> ip<-mtcars[,c("am", "wt", "mpg")]

> print(head(ip))

> am.data<-glm(formula=am~wt+mpg,data=ip,family=binomial)

> summary(am.data)

```

Output:The image shows a screenshot of an R Console window. The title bar says "R Console". The output text is as follows:

```
Call:
glm(formula = am ~ wt + mpg, family = binomial, data = ip)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.50806  -0.45191  -0.04684   0.24664   2.01168

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  25.8866    12.1935   2.123  0.0338 *
wt          -6.4162     2.5466  -2.519  0.0118 *
mpg         -0.3242     0.2395  -1.354  0.1759
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 43.230  on 31  degrees of freedom
Residual deviance: 17.184  on 29  degrees of freedom
AIC: 23.184

Number of Fisher Scoring iterations: 7

> |
```

Conclusion: Logistic regression in R Programming is a classification algorithm used to find the probability of event success and event failure. Logistic regression is used when the dependent variable is binary (0/1, True/False, Yes/No) in nature. Logit function is used as a link function in a binomial distribution. we have successfully performed logistic regression.