

Index

S No.	Topic	Date	Teacher's Signature
Exp - 0	Introduction To IBM SPSS		
Exp - 1	Transportation of Dataset to SPSS Editor		
Exp - 2	Merging of Dataset and providing missing values		
Exp - 3	Graphical representation of imported dataset using Bar and Pie charts		
Exp - 4	Drawing of Histogram and Distribution Curve		
Exp - 5	Descriptive statistics like mean, mode, median, variance, skewness, kurtosis, etc		
Exp - 6	Correlation between two Random Variables		
Exp - 7	Regression Analysis		
Exp - 8	Hypothesis Testing 'T'- test		
Exp - 9	Chi Square Test		
Exp - 10	ANOVA One-Way		

Introduction to IBM SPSS

Statistical Package for Social Sciences (**SPSS**) is a proprietary software of IBM.

Data Editor: It displays the contents of the data file. This is the window that opens automatically when you start an SPSS session. In this window, you can create new data files or modify existing ones.

1. The data editor provides two views of your data:

- **Data View** - It displays the data view. Each variable is a column & each row is a case.

The screenshot shows the IBM SPSS Data Editor window titled "car_sales.sav [DataSet1] - IBM SPSS Statistics Data Editor". The menu bar includes File, Edit, View, Data, Transform, Analyze, Graphs, Utilities, Extensions, Window, and Help. The toolbar contains icons for opening, saving, printing, and various data manipulation functions. The main area displays a table with 11 rows of data and 13 columns. The columns are labeled: manufact, model, sales, resale, type, price, engine_s, horsepower, wheelbas, width. The first few rows of data are as follows:

	manufact	model	sales	resale	type	price	engine_s	horsepow	wheelbas	width
1	Acura	Integra	16.919	16.360	0	21.500	1.8	140	101.2	67.1
2	Acura	TL	39.384	19.875	0	28.400	3.2	225	108.1	70.1
3	Acura	CL	14.114	18.225	0	.	3.2	225	106.9	70.1
4	Acura	RL	8.588	29.725	0	42.000	3.5	210	114.6	71.1
5	Audi	A4	20.397	22.255	0	23.990	1.8	150	102.6	68.1
6	Audi	A6	18.780	23.555	0	33.950	2.8	200	108.7	76.1
7	Audi	A8	1.380	39.000	0	62.000	4.2	310	113.0	74.1
8	BMW	323i	19.747	.	0	26.990	2.5	170	107.3	68.1
9	BMW	328i	9.231	28.675	0	33.400	2.8	193	107.3	68.1
10	BMW	528i	17.527	36.125	0	38.900	2.8	193	111.4	70.1
11	Ruick	Century	91.561	12.475	0	21.975	3.1	175	109.0	72.1

- **Variable View** – It displays variable definition information, including defined variables and value labels, data type (for e.g. string, date, and numeric), measurement level (nominal, ordinal, scale) and user-defined missing values. Some of the points are discussed below.

The screenshot shows the IBM SPSS Data Editor window titled "car_sales.sav [DataSet1] - IBM SPSS Statistics Data Editor". The menu bar includes File, Edit, View, Data, Transform, Analyze, Graphs, Utilities, Extensions, Window, and Help. The toolbar contains icons for opening, saving, printing, and various data manipulation functions. The main area displays a table with 13 rows of variable definitions. The columns are labeled: Name, Type, Width, Decimals, Label, Values, Missing, Columns, Align, Measure, Role. The first few rows of data are as follows:

	Name	Type	Width	Decimals	Label	Values	Missing	Columns	Align	Measure	Role
1	manufact	String	13	0	Manufacturer	None	None	13	Left	Nominal	Input
2	model	String	17	0	Model	None	None	17	Left	Nominal	Input
3	sales	Numeric	11	3	Sales in thousand...	None	None	8	Right	Scale	Input
4	resale	Numeric	11	3	4-year resale va...	None	None	8	Right	Scale	Input
5	type	Numeric	11	0	Vehicle type	{0, Automobi...	None	8	Right	Ordinal	Input
6	price	Numeric	11	3	Price in thousands	None	None	8	Right	Scale	Input
7	engine_s	Numeric	11	1	Engine size	None	None	8	Right	Scale	Input
8	horsepow	Numeric	11	0	Horsepower	None	None	8	Right	Scale	Input
9	wheelbas	Numeric	11	1	Wheelbase	None	None	8	Right	Scale	Input
10	width	Numeric	11	1	Width	None	None	8	Right	Scale	Input
11	length	Numeric	11	1	Length	None	None	8	Right	Scale	Input
12	curb_wgt	Numeric	11	3	Curb weight	None	None	8	Right	Scale	Input
13	fuel_cap	Numeric	11	1	Fuel capacity	None	None	8	Right	Scale	Input

The following variable attributes (including others) can be entered and modified in the variable view:

- I. **Variable Name:** It is the name given to a variable and each name should be unique: duplication is not allowed. *Variable names cannot contain spaces.* Variable name can be defined with any mixture of upper case and lower case characters and are preserved for display purposes.
- II. **Type:** It specifies the type of data for each variable. It can be one of the following;
 - a. **Numeric:** A variable whose values are numbers. Values are displayed in standard numeric format.
 - b. **Comma:** A numeric value whose values are displayed with commas delimiting every three places and displayed with the period as a decimal delimiter(eg. 1,234 or 1,234,567)
 - c. **Dot:** A numeric variable whose values are displayed with periods delimiting every three places and with a comma as a decimal limiter (eg. 1.234 or 1.234.567).
 - d. **Scientific Notation:** A numeric variable whose values are displayed with an embedded E and a signed power of 10 exponents. For eg 123, 1.23E2, 1.23E+2 and 1.23+2.
 - e. **Date:** A numerical variable whose values are displayed in one of the calendar-date or clock-time formats.
 - f. **Dollar:** a numerical variable displayed with a leading dollar sign (\$), commas delimiting every three places, and a period as a decimal delimiter.
 - g. **Custom Currency:** Numeric variable whose values are displayed in one of the custom currency formats that you have defined on the currency tab of the options dialog box.
 - h. **String:** A variable whose values are not numeric and therefore are not used in calculations.
 - i. **Restricted Numeric:** A variable whose values are restricted to a nonnegative integer.
- III. **Width:** It is used to specify the number of digits or characters that can be entered.
- IV. **Decimals:** It is used to specify the number of decimal places in data representation.
- V. **Label:** The name of variable can be described in detail here.
- VI. **Values:** Nominal or categorical variable can be assigned numerical values. eg. If variable is gn (Gender) and entered as male or female then here assigned values 1 or 2 respectively
- VII. **Missing values:** These are user-defined missing values.
- VIII. **Columns:** It is used to specify the width of columns (can be changed by dragging also).
- IX. **Measurement Level:** It can be specified as follows:

- a) **Ordinal**- A categorical variable where sense of ordering (eg. low, medium high) holds. Such type of variables has several ordered categories. Examples of ordinal variables include attitude scores representing the degree of satisfaction or confidence and preference rating scores.
- b) **Scale**- A variable can be treated as a scale (continuous) when its value represents ordered categories with a meaningful metric, so that distance comparisons between values are appropriate. Examples of scale variables include age in years and income in thousands of dollars.
- c) **Nominal**- A categorical variable used without any ordering. Such type of variable has several unordered categories. Examples of nominal variables include gender, region, zip code, and religious affiliation.

Variable types - IBM Documentation

Table 1. Measurement level icons

	Numeric	String	Date	Time
Scale (Continuous)		n/a		
Ordinal				
Nominal				

2. Data: Enter data in data view as per the property of variable specified in variable view various formats, or data (if already stored in other format) can be imported to SPSS editor and the properties of variables can be changed according to need in variable view.

3. Operations: Any operation on data value will be done & observed in data view.

Descriptives

Sales in thousands [sales]
Vehicle type [type]
Price in thousands [price]
Engine size [engine_s]
Horsepower [horsepow]
Wheelbase [wheelbas]
Width [width]
Length [length]
Curb weight [curb_wgt]
Fuel capacity [fuel_cap]
Fuel efficiency [mpg]

Save standardized values as variables

OK | Paste | Reset | Cancel | Help

Variable(s): 4-year resale value [resale]

Options...
Style...
Bootstrap...

Mean Sum

Dispersion
 Std. deviation Minimum
 Variance Maximum
 Range S.E. mean

Distribution
 Kurtosis Skewness

Display Order
 Variable list
 Alphabetic
 Ascending means
 Descending means

Continue | Cancel | Help

And statistical operation's information will be created in another window

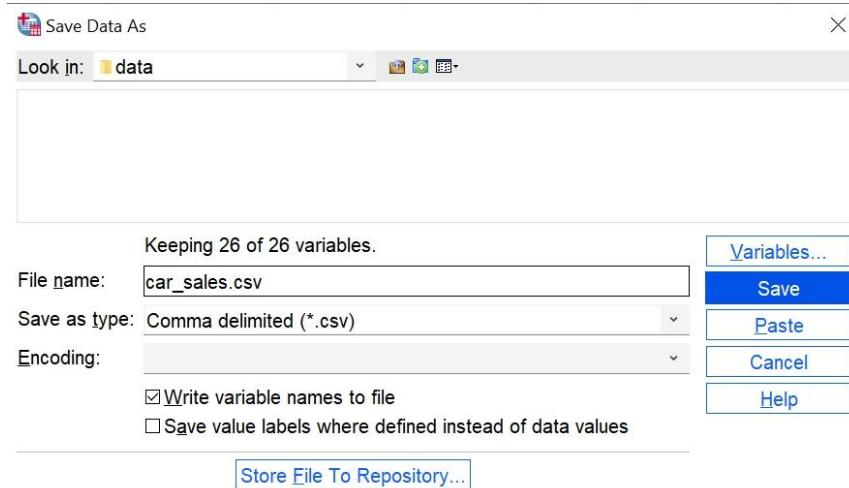
→ Descriptives

[DataSet1] C:\Users\Administrator\OneDrive\Desktop\Lab_Semester_3_Data\Probability And Statistics\data\car_sales.sav

Descriptive Statistics

	N	Minimum	Maximum	Mean	Std. Deviation
4-year resale value	121	5.160	67.550	18.07298	11.453384
Valid N (listwise)	121				

4. Save: The last step is to 'Save' the files which can be done in various format using 'Export' option under 'File'



5. Saved: File will be saved at desired location in desired format.



Experiment - 1

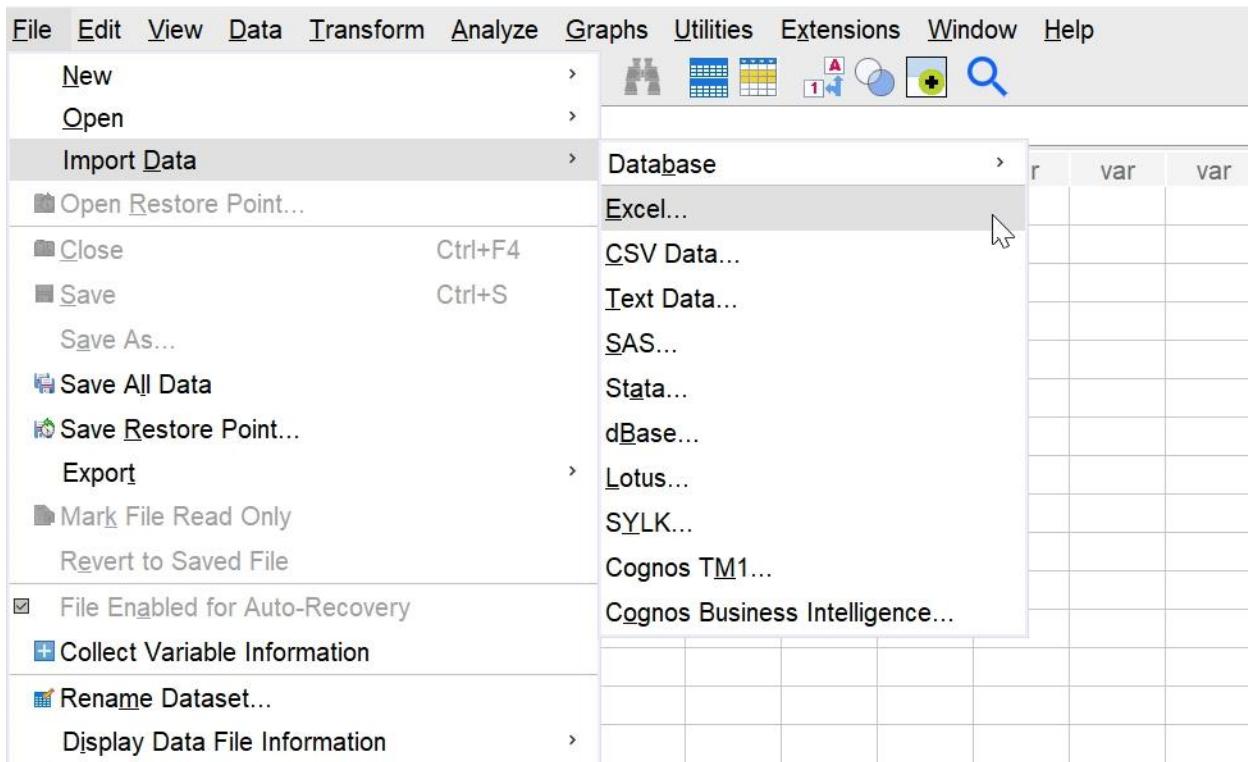
Transportation of Dataset to SPSS Editor

STEPS

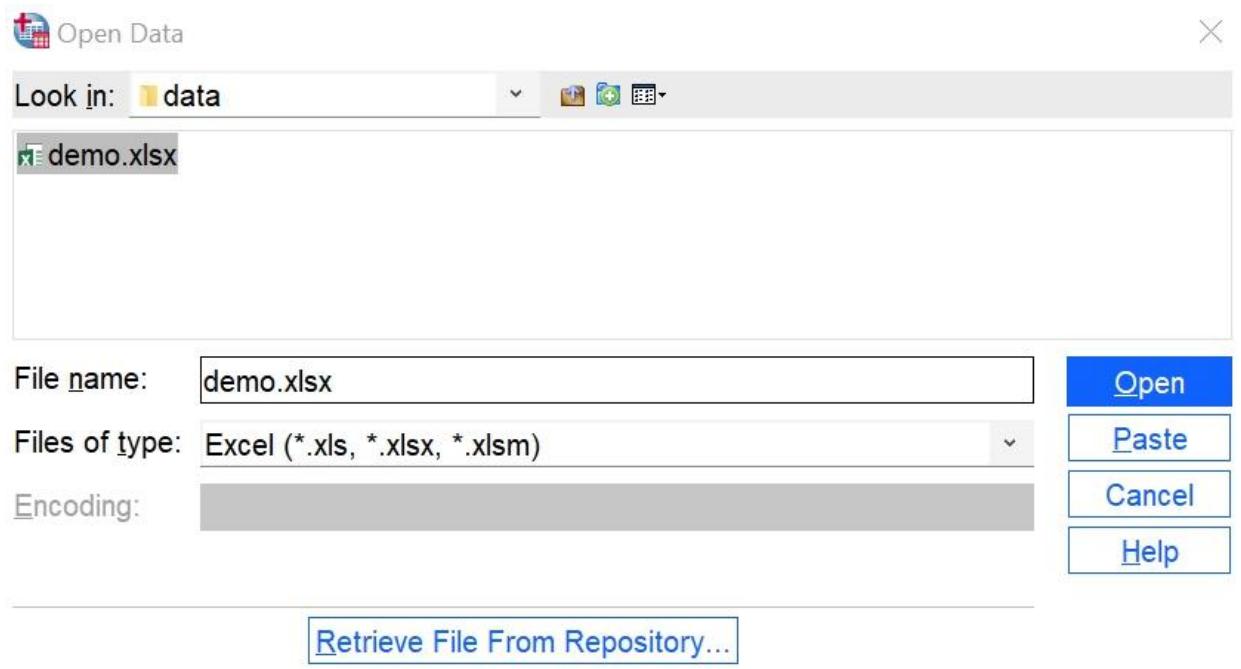
CASE – 1

- Data from an excel file

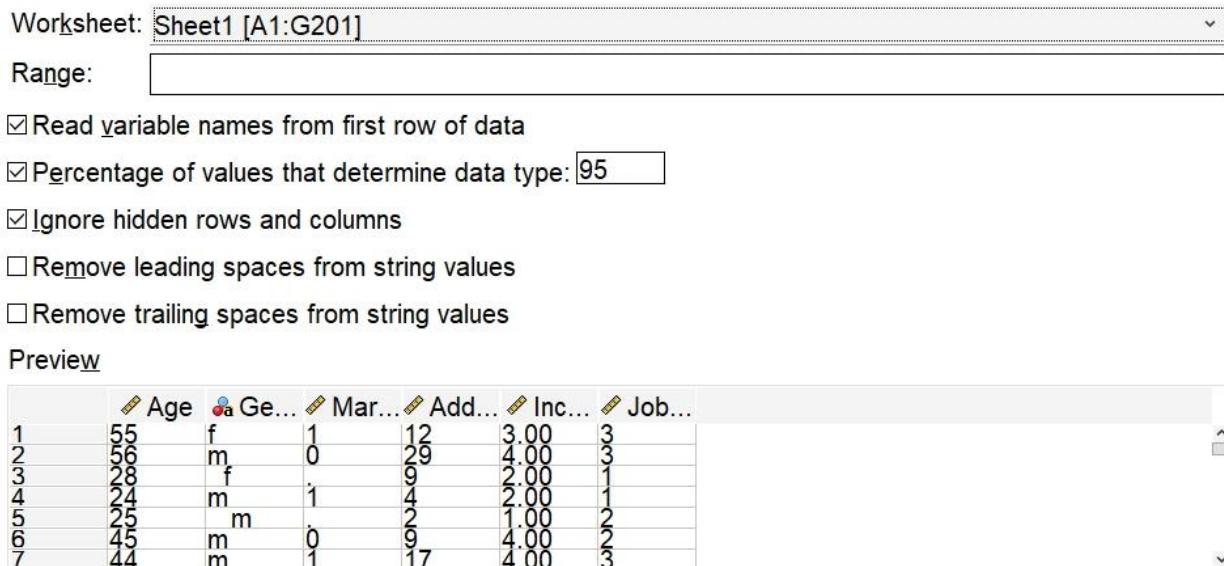
Commands → **File>Open>Import Data>Excel**



Commands -> click <file-name> -> Open



Commands -> Press OK



(i) Final data type is based on all data and can be different from the preview, which is based on the first 200 data rows. The preview displays only the first 500 columns.

OK **Paste** **Reset** **Cancel** **Help**

Data View

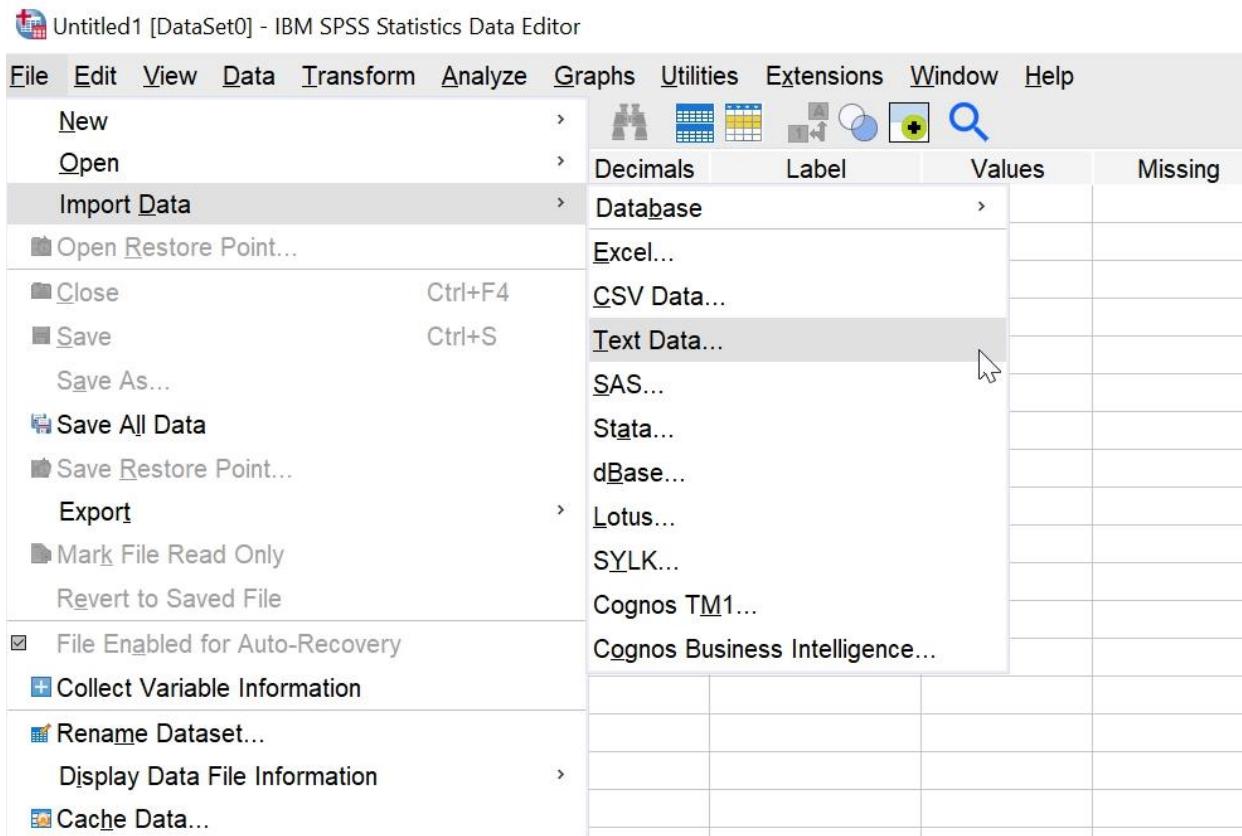
	Age	Gender	MaritalStatus	Address	IncomeCategory	JobCategory	var	var
1	55 f		1	12	3.00	3		
2	56 m		0	29	4.00	3		
3	28 f		.	9	2.00	1		
4	24 m		1	4	2.00	1		
5	25 m		.	2	1.00	2		
6	45 m		0	9	4.00	2		
7	44 m		1	17	4.00	3		
8	46 m		.	20	4.00	3		
9	41 m		.	10	2.00	2		
10	29 f		.	4	1.00	2		
11	34 m		0	0	4.00	2		
12	55 f		0	17	3.00	1		
13	28 m		0	9	3.00	1		
14	21 f		1	2	1.00	1		
15	55 f		0	8	4.00	2		
16	35 m		0	8	3.00	2		
17	45 f		0	4	2.00	2		
18	21 m		0	1	2.00	1		
19	32 f		0	0	2.00	1		
20	42 f		0	9	4.00	3		
21	40 f		1	12	4.00	3		
22	36 f		0	6	2.00	1		
23	42 m		1	13	3.00	2		
24	65 m		1	17	2.00	3		
25	52 m		1	5	4.00	3		
26	51 m		1	17	4.00	2		
27	44 m		1	1	2.00	2		
28	26 f		0	6	2.00	1		
29	41 f		1	19	3.00	3		
30	49 m		0	11	2.00	2		
31	64 f		0	36	4.00	3		
32	39 m		0	8	3.00	2		
33	53 m		0	8	4.00	3		
34	38 f		1	17	2.00	2		
35	46 f		1	6	2.00	3		
36	58 m		0	2	3.00	1		
37	25 f		0	0	3.00	1		
38	57 f		0	28	4.00	3		
39	47 m		0	21	2.00	3		
40	21 f		0	0	1.00	1		
41	45 f		1	21	3.00	3		
42	56 f		0	7	4.00	3		
43	24 m		0	2	1.00	1		

Variable View

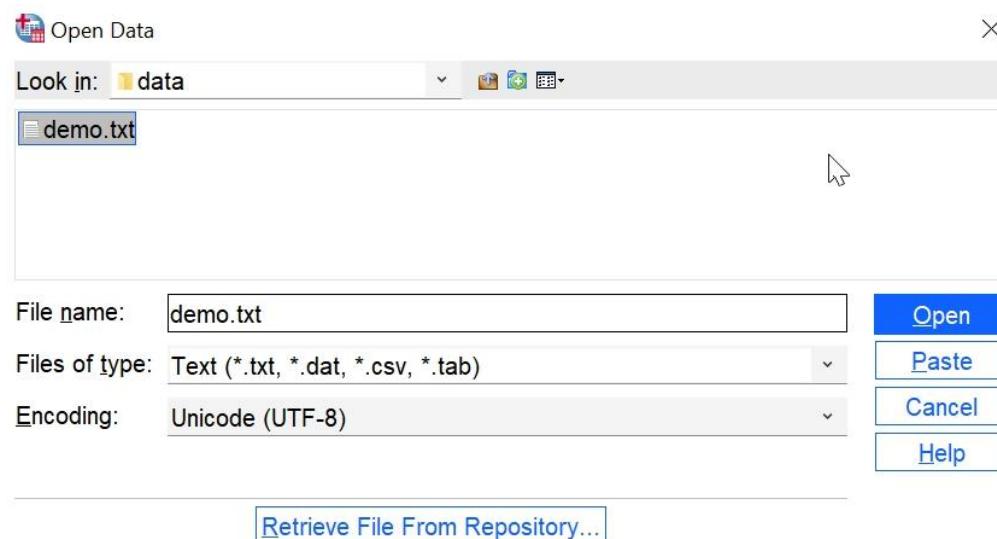
CASE – 2

- Data from an text file

Commands → File>Open>Import Data>Text Data



Commands → click <file-name> → Open



Commands -> Press Next (Till all the steps end)

Text Import Wizard - Step 1 of 6

Welcome to the text import wizard!

This wizard will help you read data from your text file and specify information about the variables.

Does your text file match a predefined format?

Yes

No

Text file: C:\Users\Administrator\OneDrive\Desktop\Lab_Semester_3_Data\Probability And Statistics\data\demo...

Name	Age	Gender	Marital Status	Address	Income	Income Category	Job Category
2	55	f	1	12	72.00	3.00	3
3	56	m	0	29	153.00	4.00	3
4	28	f	no answer	9	28.00	2.00	1
5	24	m	1	4	26.00	2.00	1
6	25	m	no answer	2	23.00	1.00	2
7	45	m	0	9	76.00	4.00	2
8	44	m	1	17	144.00	4.00	3
9	46	m	no answer	20	75.00	4.00	3
10	41	m	no answer	10	26.00	2.00	2
11	29	f	no answer	4	19.00	1.00	2
12	34	m	0	0	89.00	4.00	2
13	55	f	0	17	72.00	3.00	1
14	28	m	0	9	55.00	3.00	1
15	21	f	1	2	20.00	1.00	1
16	55	f	0	8	283.00	4.00	2
17	35	m	0	8	70.00	3.00	2

< Back Finish Cancel Help

Text Import Wizard - Step 2 of 6

How are your variables arranged?

Delimited - Variables are delimited by a specific character (i.e., comma, tab).
 Fixed width - Variables are aligned in fixed width columns.

Are variable names included at the top of your file?

Yes
 Line number that contains variable names:

No

What is the decimal symbol?

Period
 Comma

Text file: C:\Users\Administrator\OneDrive\Desktop\Lab_Semester_3_Data\Probability And Statistics\data\demo...

Name	Age	Gender	Marital Status	Address	Income	Income Category	Job Category
2	55	f	1	12	72.00	3.00	3
3	56	m	0	29	153.00	4.00	3
4	28	f	no answer	9	28.00	2.00	1
5	24	m	1	4	26.00	2.00	1
6	25	m	no answer	2	23.00	1.00	2
7	45	m	0	9	76.00	4.00	2
8	44	m	1	17	144.00	4.00	3
9	46	m	no answer	20	75.00	4.00	3
10	41	m	no answer	10	26.00	2.00	2

< Back Finish Cancel Help

Text Import Wizard - Delimited Step 3 of 6

X

The first case of data begins on which line number?

How are your cases represented?

- Each line represents a case
- A specific number of variables represents a case:

How many cases do you want to import?

- All of the cases
- The first cases.
- A random percentage of the cases (approximate): %

Data preview

Name	Age	Gender	Marital Status	Address	Income	Income Category	Job Category
2	55	f	1	12	72.00	3.00	3
3	56	m	0	29	153.00	4.00	3
4	28	f	no answer	9	28.00	2.00	1
5	24	m	1	4	26.00	2.00	1
6	25	m	no answer	2	23.00	1.00	2
7	45	m	0	9	76.00	4.00	2
8	44	m	1	17	144.00	4.00	3
9	46	m	no answer	20	75.00	4.00	3
10	41	m	no answer	10	26.00	2.00	2
11	29	f	no answer	4	19.00	1.00	2
12	34	m	0	0	89.00	4.00	2
13	55	f	0	17	72.00	3.00	1
14	28	m	0	9	55.00	3.00	1

< Back **Next >** Finish Cancel Help

Text Import Wizard - Delimited Step 4 of 6

X

Which delimiters appear between variables?

- Tab Space
- Comma Semicolon
- Other:

What is the text qualifier?

- None
- Single quote
- Double quote
- Other:

Leading and Trailing Spaces

- Remove leading spaces from string values
- Remove trailing spaces from string values

Data preview

Age	Gender	Marital	Status	Addr...	Income	Income	Cate...	Job	Cate...	V11
55	f	1	12	72.00	3.00	3				
56	m	0	29	153.00	4.00	3				
28	f	no answer	9	28.00	2.00	1				
24	m	1	4	26.00	2.00	1				
25	m	0	9	76.00	4.00	2				
44	m	1	17	144.00	4.00	3				
46	m	no answer	20	75.00	4.00	3				
41	m	no answer	10	26.00	2.00	2				
29	f	no answer	4	19.00	1.00	2				
34	m	0	0	89.00	4.00	2				
55	f	0	17	72.00	3.00	1				
28	m	0	9	55.00	3.00	1				
21	f	1	2	20.00	1.00	1				
55	f	0	8	283.00	4.00	2				
35	m	0	8	70.00	3.00	2				
25	f	0	4	48.00	3.00	2				
21	m	0	1	37.00	2.00	1				
32	f	0	0	28.00	2.00	1				
42	f	0	9	109.00	4.00	3				

< Back **Next >** Finish Cancel Help

 Text Import Wizard - Step 6 of 6

X

You have successfully defined the format of your text file.

Would you like to save this file format for future use?

Yes

[Save As...](#)

◎No

Would you like to paste the syntax?

Cache data locally

Yes

• No

Press the Finish button to complete the text import wizard.

Data preview

< Back Next > Finish Cancel Help

Commands → Press **Finish**

Variable View

Data View

*Untitled2 [DataSet1] - IBM SPSS Statistics Data Editor

File Edit View Data Transform Analyze Graphs Utilities Extensions Window Help

The screenshot shows the IBM SPSS Statistics Data Editor window. The menu bar includes File, Edit, View, Data, Transform, Analyze, Graphs, Utilities, Extensions, Window, and Help. Below the menu is a toolbar with various icons for file operations like Open, Save, Print, and data manipulation. The main area displays a data grid with 43 rows and 13 columns. The columns are labeled: Age, Gender, Marital, Status, Address, Income_A, Income, Category_A, Job, Category, and V11. The data includes demographic information such as age, gender, marital status, address, income levels, and job categories. A cursor is visible over the data grid.

	Age	Gender	Marital	Status	Address	Income_A	Income	Category_A	Job	Category	V11
1	55 f		1	12	72.00	3.00	3.00
2	56 m		0	29	153.00	4.00	3.00
3	28		9.00	28.00	2.00	1.00	.
4	24 m		1	4	26.00	2.00	1.00
5	25		2.00	23.00	1.00	2
6	45 m		0	9	76.00	4.00	2.00
7	44 m		1	17	144.00	4.00	3.00
8	46 m		.	.	20.00	75.00	4.00	3.00	.	.	.
9	41 m		.	.	10.00	26.00	2.00	2.00	.	.	.
10	29 f		.	.	4.00	19.00	1.00	2.00	.	.	.
11	34 m		0	0	89.00	4.00	2.00
12	55 f		0	17	72.00	3.00	1.00
13	28 m		0	9	55.00	3.00	1.00
14	21 f		1	2	20.00	1.00	1.00
15	55 f		0	8	283.00	4.00	2.00
16	35 m		0	8	70.00	3.00	2.00
17	45 f		0	4	48.00	2.00	2.00
18	21 m		0	1	37.00	2.00	1.00
19	32 f		0	0	28.00	2.00	1.00
20	42 f		0	9	109.00	4.00	3.00
21	40 f		1	12	117.00	4.00	3.00
22	36 f		0	6	39.00	2.00	1.00
23	42 m		1	13	53.00	3.00	2.00
24	65 m		1	17	42.00	2.00	3.00
25	52 m		1	5	83.00	4.00	3.00
26	51 m		1	17	148.00	4.00	2.00
27	44 m		1	1	29.00	2.00	2.00
28	26 f		0	6	28.00	2.00	1.00
29	41 f		1	19	70.00	3.00	3.00
30	49 m		0	11	40.00	2.00	2.00
31	64 f		0	36	102.00	4.00	3.00
32	39 m		0	8	60.00	3.00	2.00
33	53 m		0	8	78.00	4.00	3.00
34	38 f		1	17	43.00	2.00	2.00
35	46 f		1	6	31.00	2.00	3.00
36	58 m		0	2	60.00	3.00	1.00
37	25 f		0	0	58.00	3.00	1.00
38	57 f		0	28	92.00	4.00	3.00
39	47 m		0	21	48.00	2.00	3.00
40	21 f		0	0	13.00	1.00	1.00
41	45 f		1	21	67.00	3.00	3.00
42	56 f		0	7	213.00	4.00	3.00
43	24 m		0	2	19.00	1.00	1.00

CONCLUSION

- Data from files in excel or text format can be transported to SPSS editor window

PRECAUTIONS

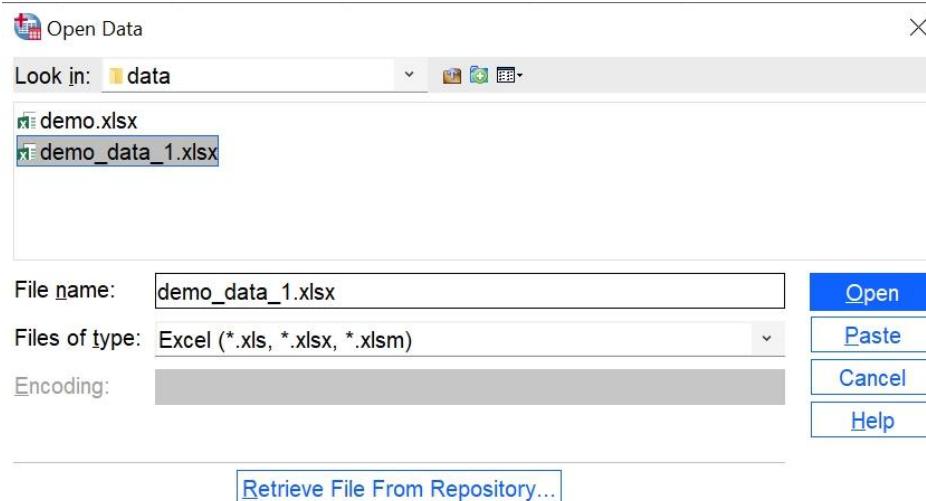
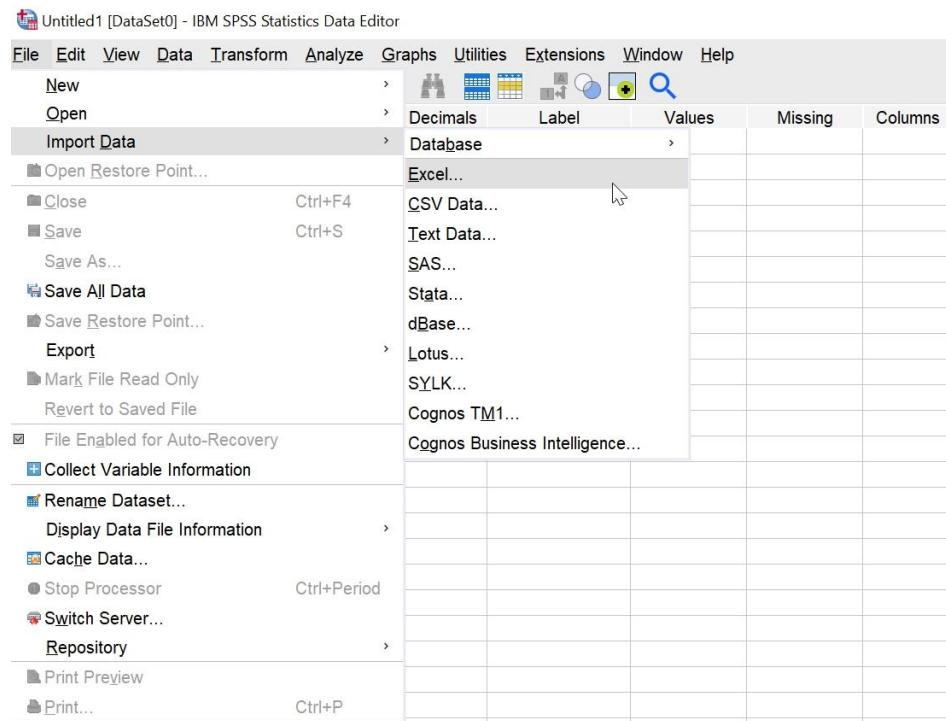
- There should be proper delimiter between data entry in text file
- Extensions of the files should be strictly taken care of

Experiment - 2

Merging of Dataset and providing missing values

STEPS

Command → File>Import Data>Excel



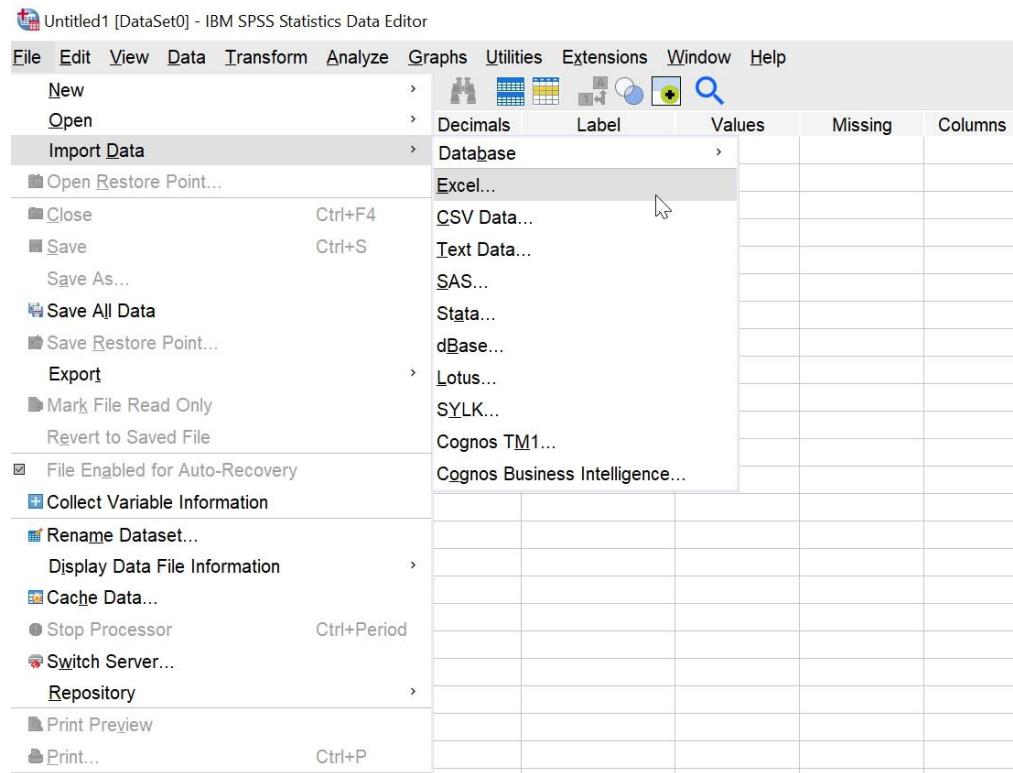
Data View

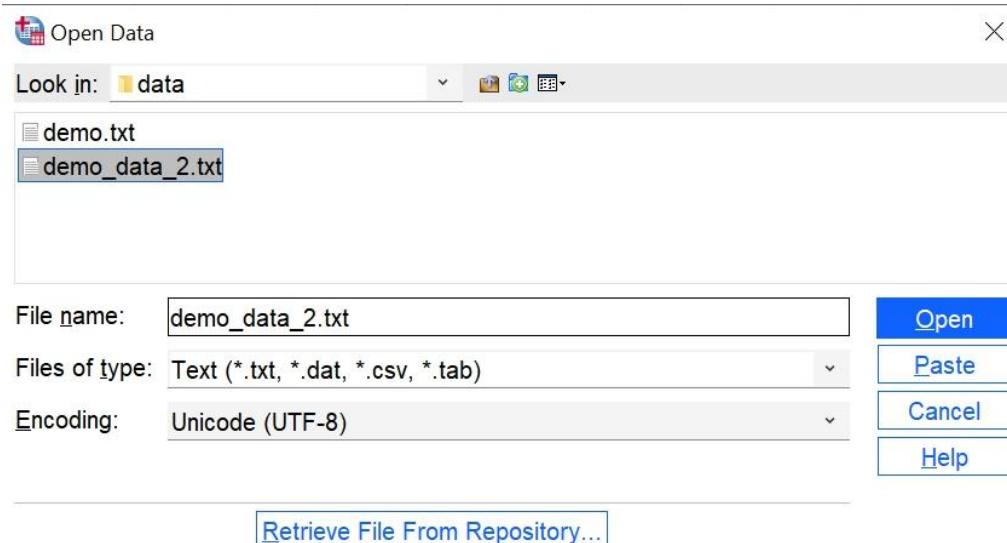
*Untitled2 [DataSet1] - IBM SPSS Statistics Data Editor

The screenshot shows the IBM SPSS Statistics Data Editor window. The menu bar includes File, Edit, View, Data, Transform, Analyze, Graphs, Utilities, Extensions, Window, and Help. Below the menu is a toolbar with various icons. The main area displays a data table with the following structure:

	name	maths	physics	english	var	var	var	
1	Gauss	48	49	43				
2	Alex	.	35	45				
3	Miyazaki	48	39	30				
4	Maya	40	.	42				
5								
6								
7								
8								
q								

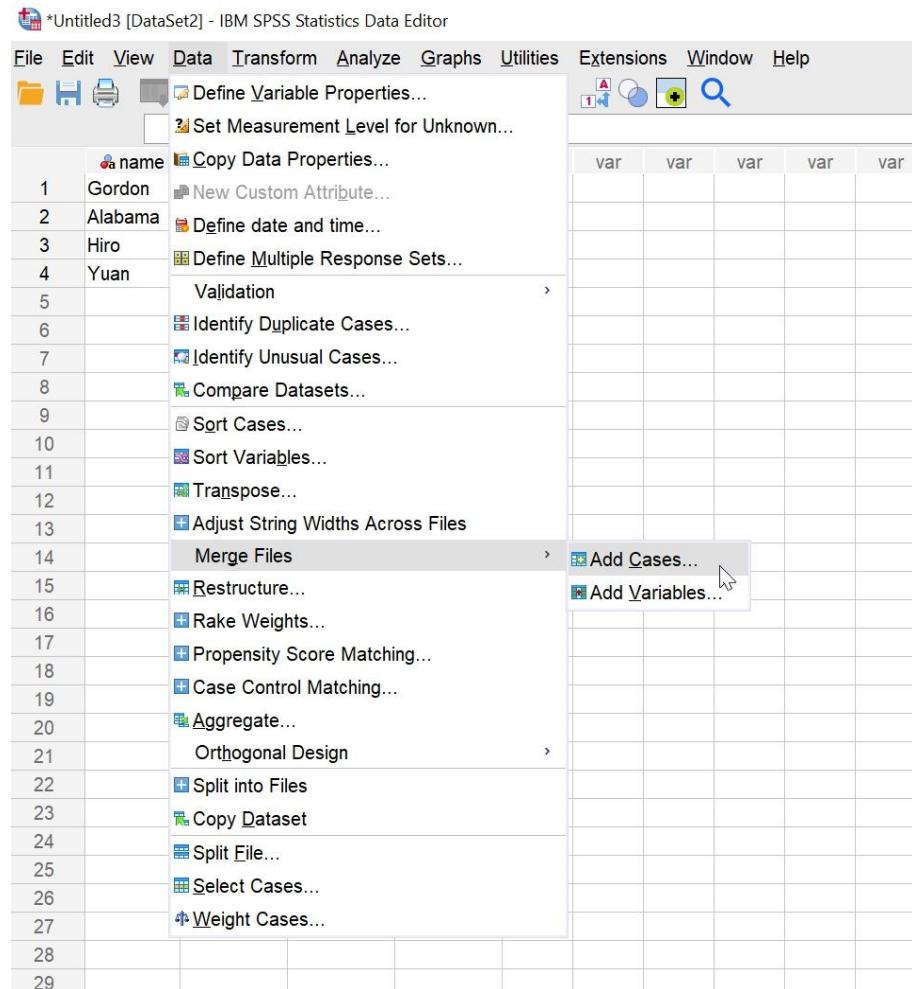
Command → File>Import Data>Text Data





Data View

Command -> Data>Merge Files>Add Cases



Command -> Click Continue

Add Cases to Untitled3[DataSet2] X

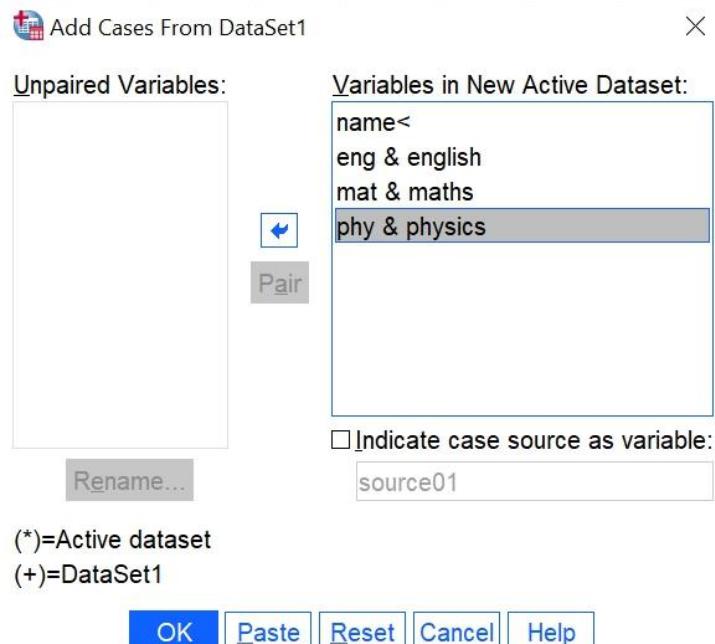
Select a dataset from the list of open datasets or from a file to merge with the active dataset

An open dataset An external SPSS Statistics data file

Untitled2[DataSet1] Browse...

Non-SPSS Statistics data files must be opened in SPSS Statistics before they can be used as part of a merge.

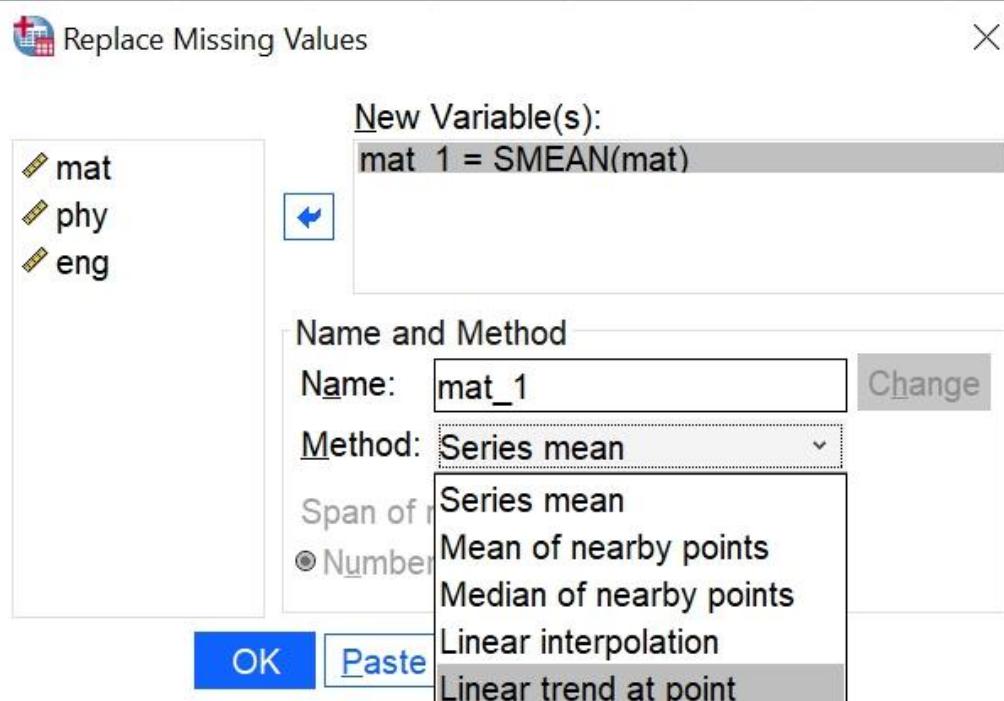
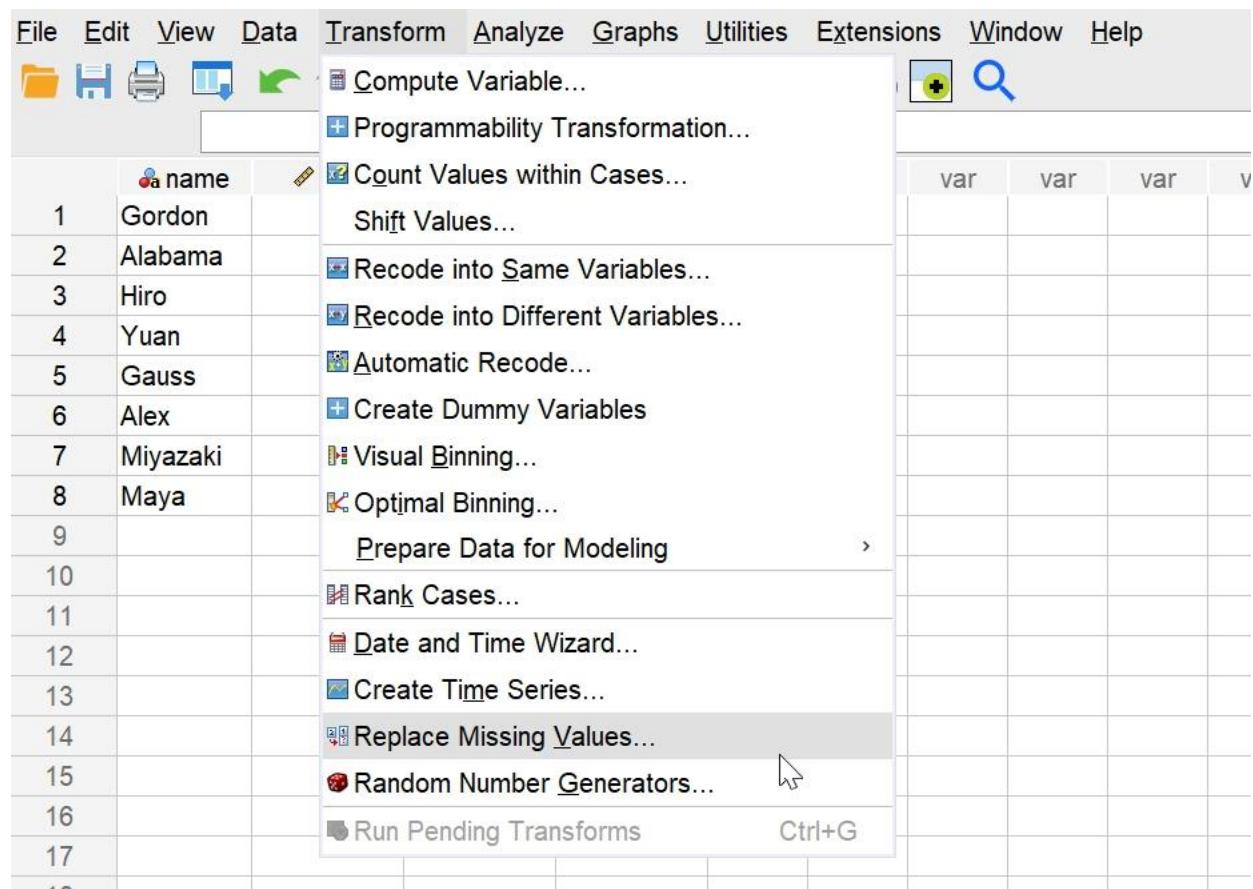
Continue Cancel Help

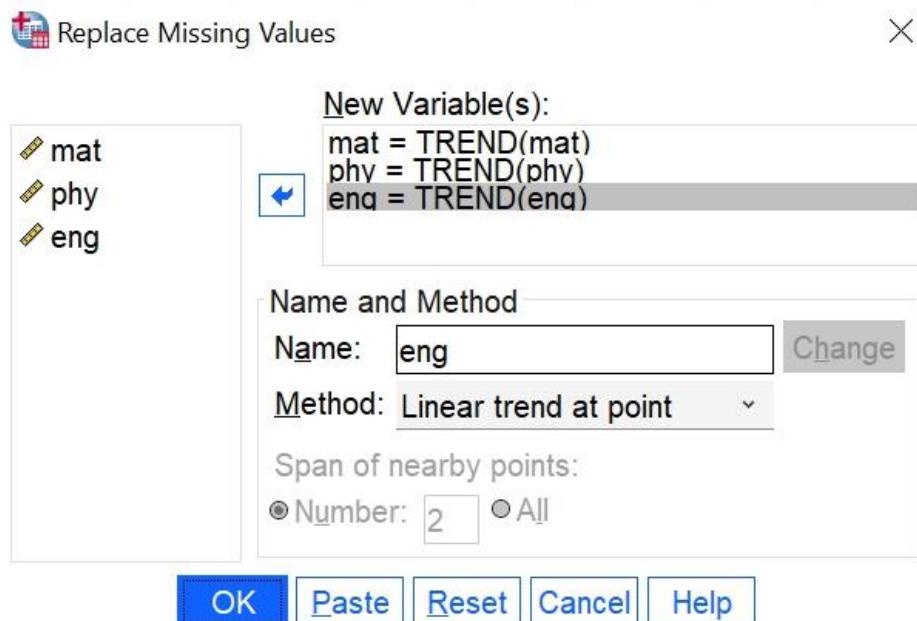


Data View (Some Missing Values)

	name	mat	phy	eng	var
1	Gordon	45	41	.	
2	Alabama	49	45	35	
3	Hiro	38	.	29	
4	Yuan	.	40	32	
5	Gauss	48	49	43	
6	Alex	.	35	45	
7	Miyazaki	48	39	30	
8	Maya	40	.	42	
9					
10					

Command -> Transform>Replace Missing Values





Data View (All Missing Values are replaced)

name	mat	phy	eng
Gordon	45.0	41.0	31.4
Alabama	49.0	45.0	35.0
Hiro	38.0	42.3	29.0
Yuan	44.7	40.0	32.0
Gauss	48.0	49.0	43.0
Alex	44.3	35.0	45.0
Miyazaki	48.0	39.0	30.0
Maya	40.0	38.9	42.0

CONCLUSION

- In this practical we have successfully merged two different types of datasets (excel file and text file) using the “Merge Dataset” feature
- Also, we have filled the missing values in the merged dataset using “Replace Missing Value” feature by using ‘Linear trend at point’ method

Experiment - 3

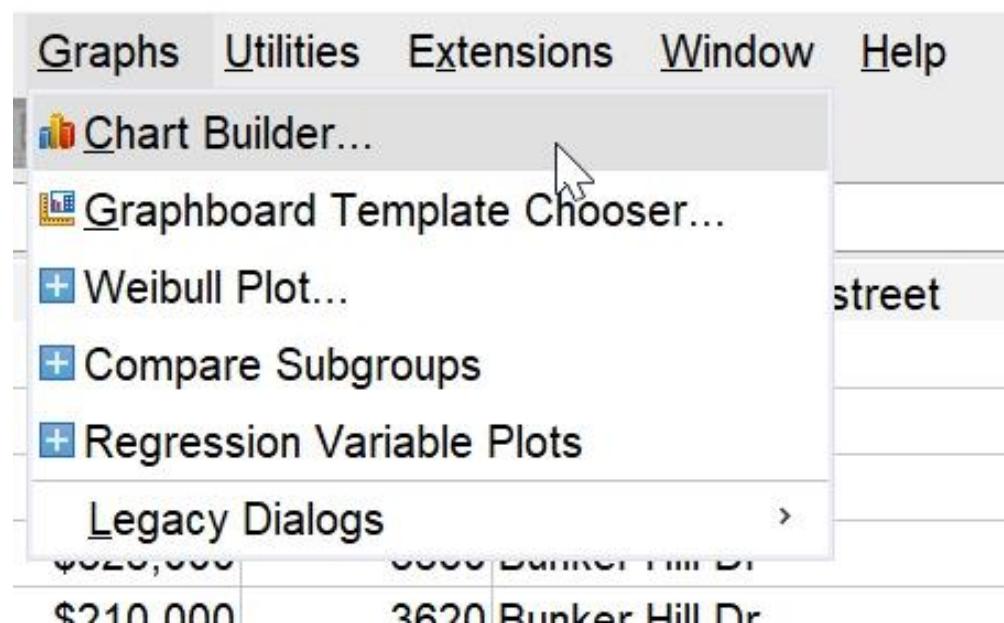
Graphical representation of imported dataset using Bar and Pie charts

STEPS

CASE – 1

- BAR Chart

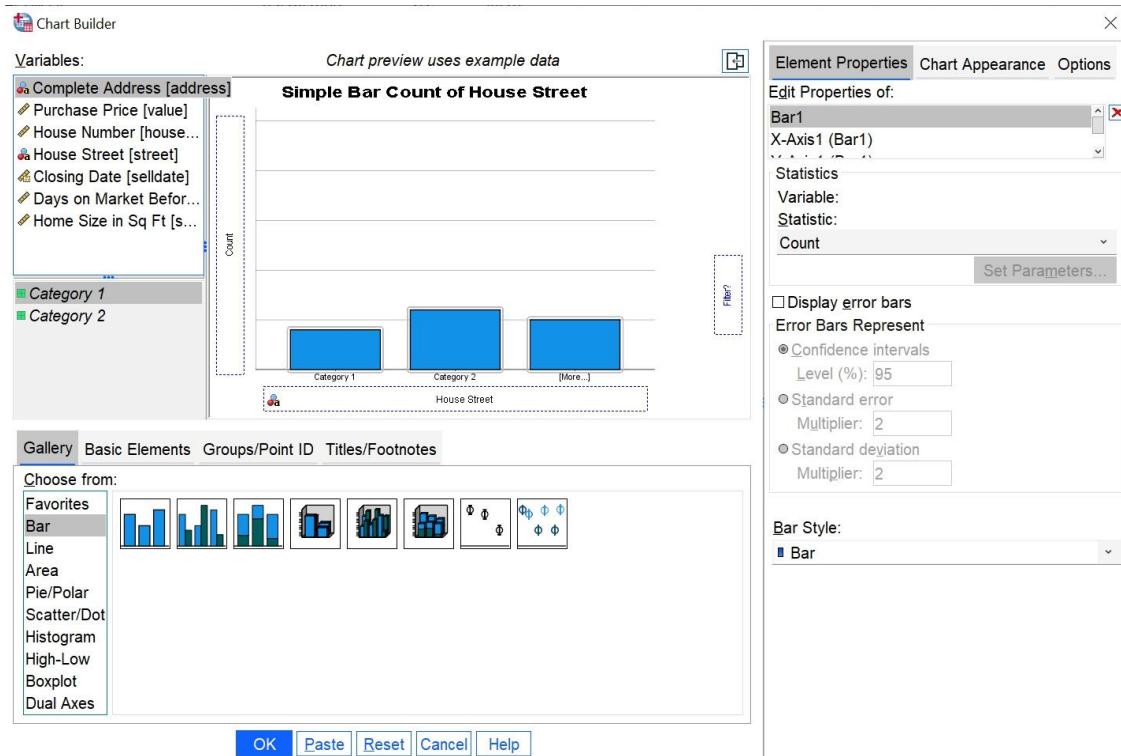
Command → [Graphs>Chart Builder](#)



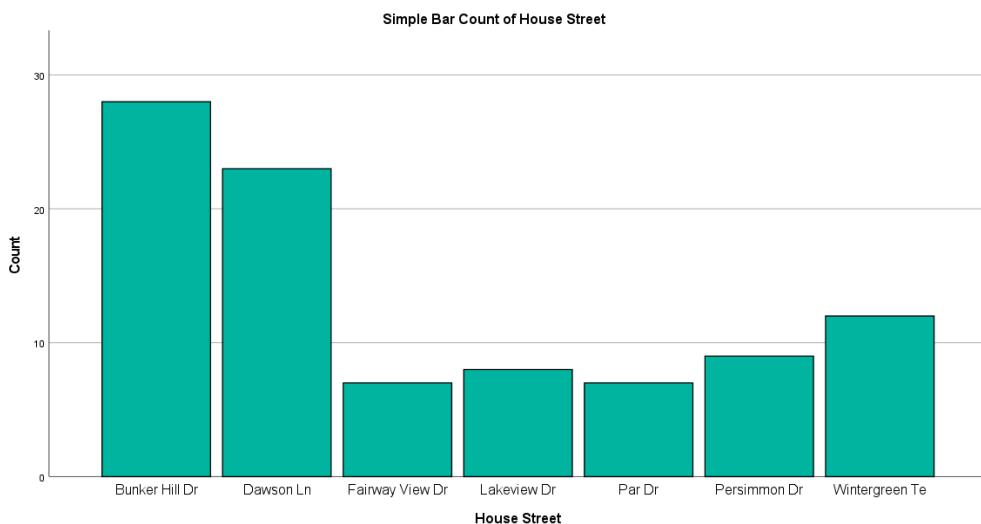
Command -> Drag Bar Chart to Chart Preview Section

Command -> Select Variables>House Street [street] to the X-axis

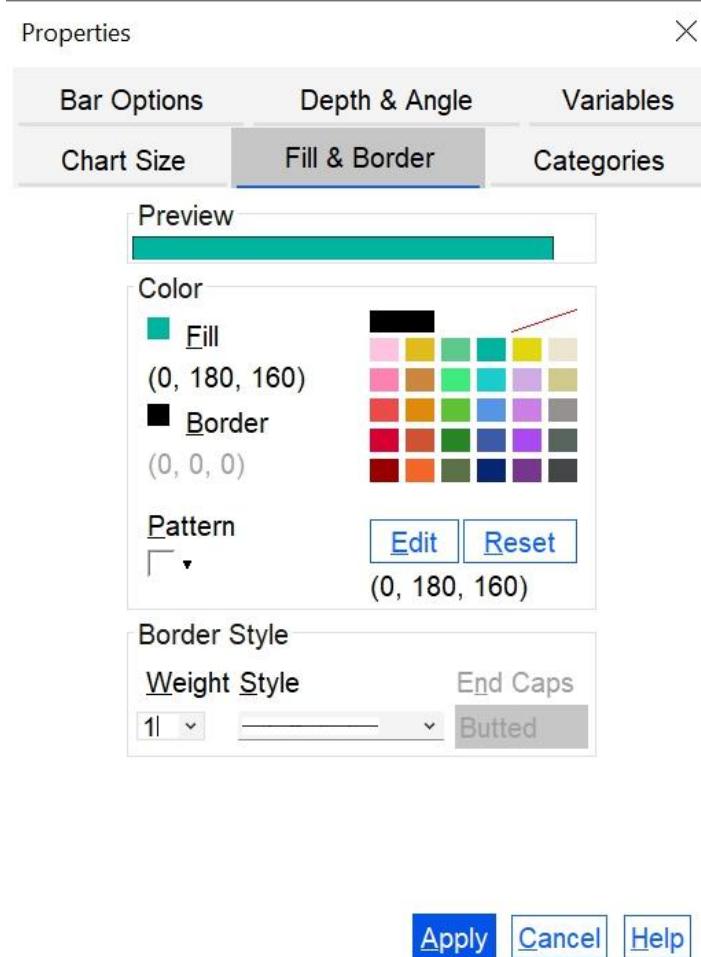
Command -> Select Ok



Bar-Graph



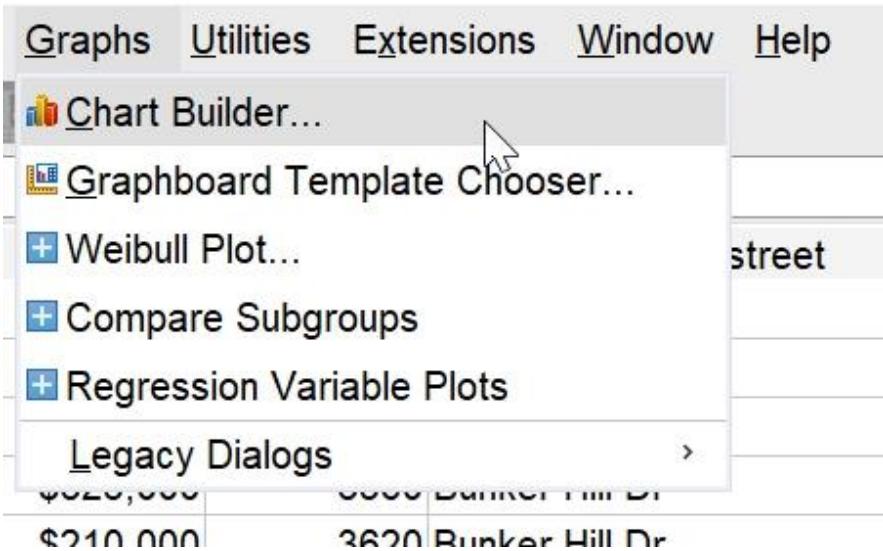
Command -> Double Click Graph to open Properties
***settings**



CASE – 2

- PIE Chart

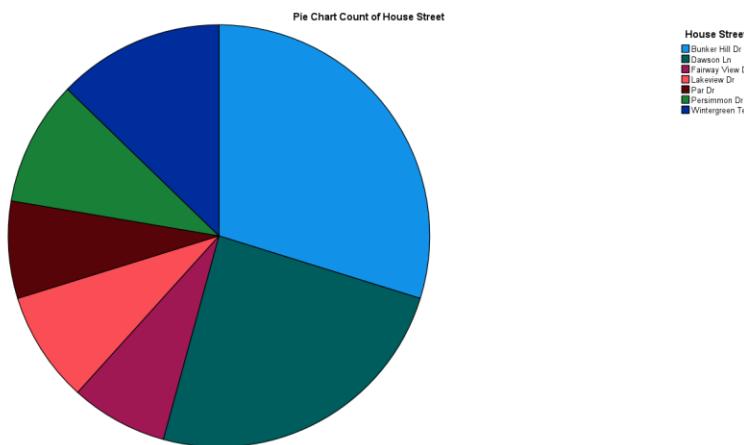
Command -> Graphs>Chart Builder



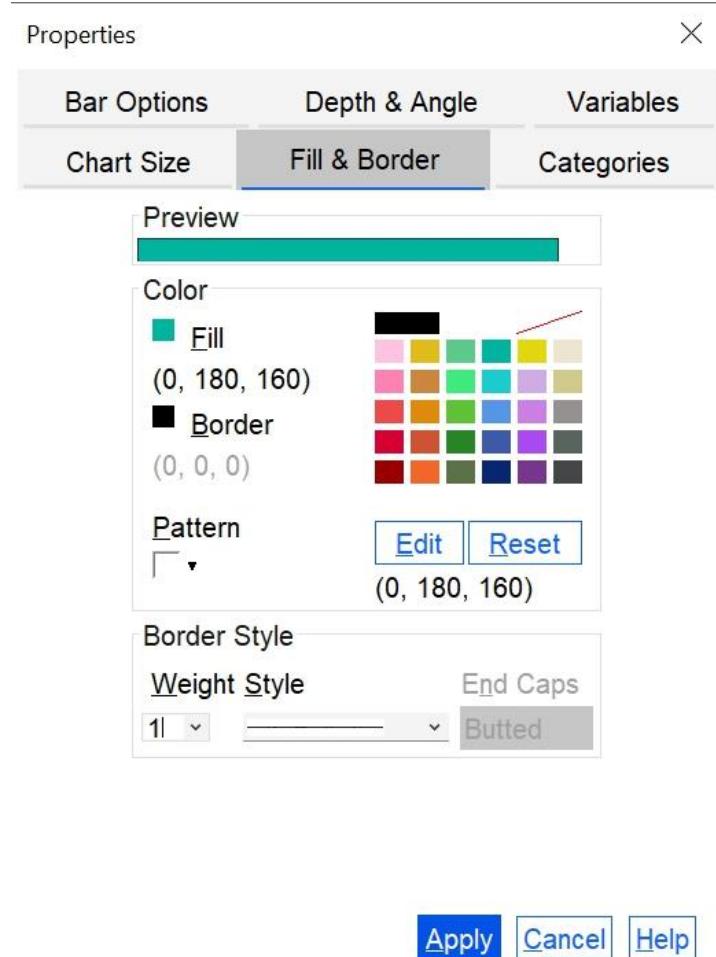
Command -> Drag Pie Chart to Chart Preview Section

Command -> Select Variables>House Street [street] to the X-axis

Command -> Select Ok



**Command -> Double Click Graph to open Properties
*settings**



CONCLUSION

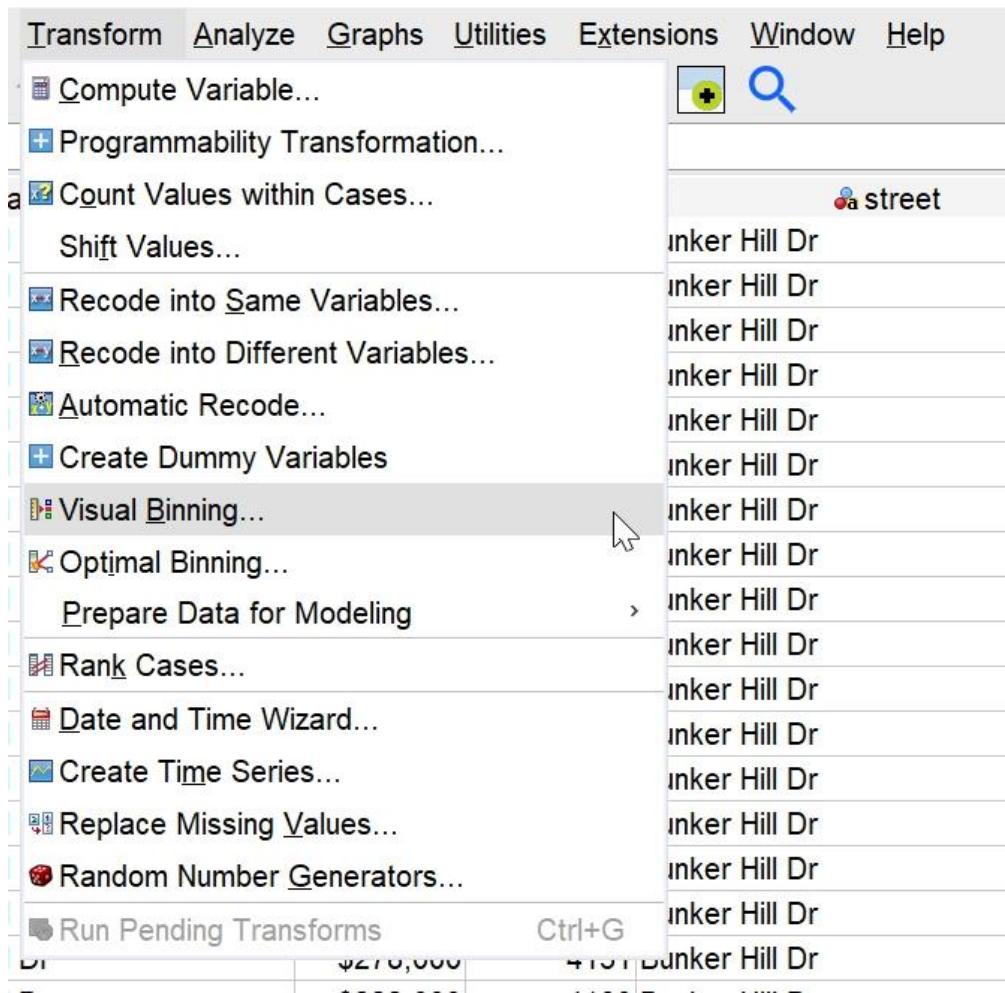
- In this practical we have successfully learnt how to make Bar and Pie Charts using the “Chart Builder” feature
- Also, we have learnt how to edit the ‘Properties’ of the graph i.e. changing border, colors, chart size, etc

Experiment - 4

Drawing of Histogram and Distribution Curve

STEPS

Command → Transform>Visual Binning



Command -> Click Continue

 Visual Binning X

i Select the variables whose values will be grouped into bins. Data will be scanned when you click Continue.

The Variables list below contains all numeric ordinal and scale variables.

Variables:  House Number...
 Closing Date [...]
 Days on Mark...
 Home Size in ...

Variables to Bin:  Purchase Pric...

Limit number of cases scanned to:

Continue **Cancel** **Help**



Command -> Click Make Cutpoints...

*fill the details as show below

 Make Cutpoints X

Equal Width Intervals
Intervals - fill in at least two fields

First Cutpoint Location:	\$100,000
Number of Cutpoints:	5
Width:	100000
Last Cutpoint Location: \$500,000	

Equal Percentiles Based on Scanned Cases
Intervals - fill in either field

Number of Cutpoints:	
Width(%):	

Cutpoints at Mean and Selected Standard Deviations Based on Scanned Cases

+/- 1 Std. Deviation
 +/- 2 Std. Deviation
 +/- 3 Std. Deviation

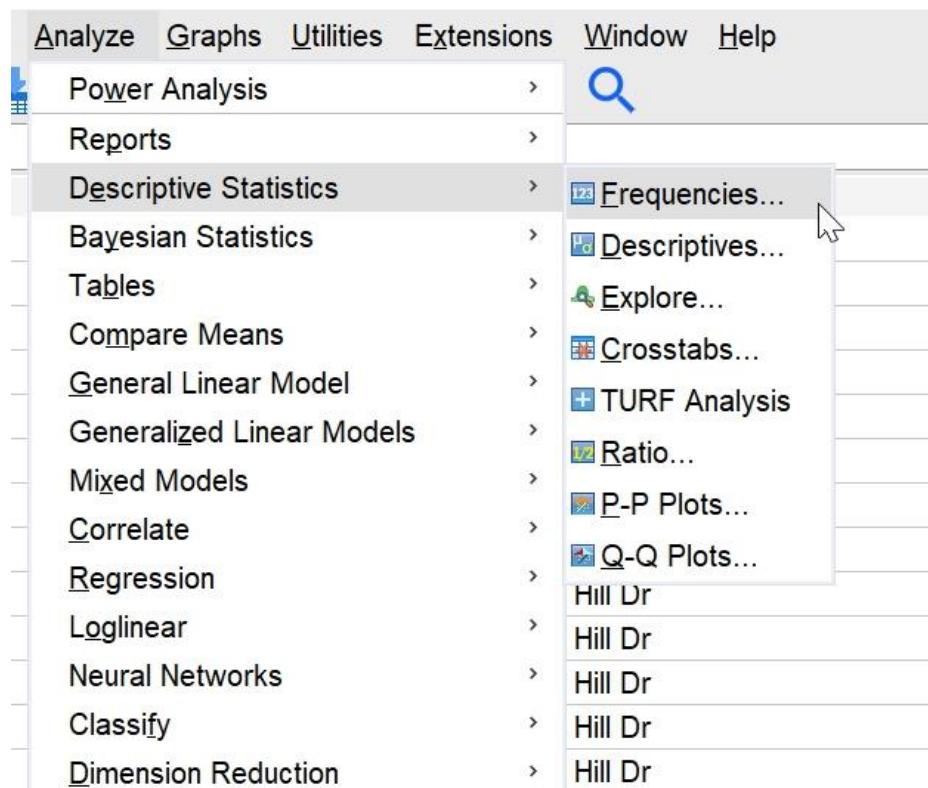
(i) Apply will replace the current cutpoint definitions with this specification.
A final interval will include all remaining values: N cutpoints produce N+1 intervals.

Apply Cancel Help

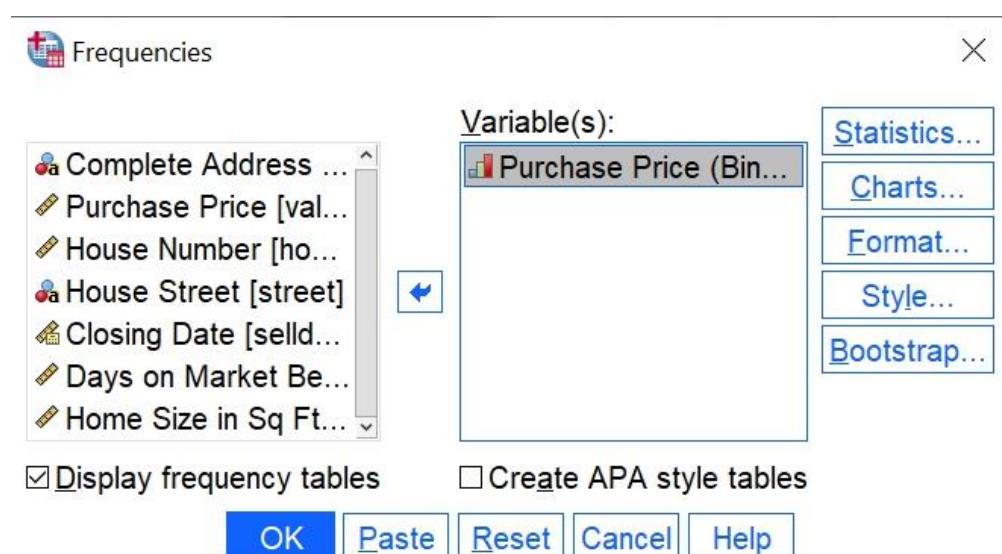
Data View (*new column has been created)

	address	value	houseenum	street	selldate	marktime	sqft	price	V
1	3520 Bunker Hill Dr	\$335,000	3520	Bunker Hill Dr	03/13/2000	91	3211	4	
2	3521 Bunker Hill Dr	\$321,000	3521	Bunker Hill Dr	06/21/2000	54	2810	4	
3	3530 Bunker Hill Dr	\$300,000	3530	Bunker Hill Dr	01/25/2000	13	3295	3	
4	3560 Bunker Hill Dr	\$325,000	3560	Bunker Hill Dr	01/01/2000	96	3034	4	
5	3620 Bunker Hill Dr	\$210,000	3620	Bunker Hill Dr	07/24/2000	18	2265	3	
6	3621 Bunker Hill Dr	\$416,000	3621	Bunker Hill Dr	03/28/2000	62	4420	5	
7	3630 Bunker Hill Dr	\$342,000	3630	Bunker Hill Dr	11/15/1999	18	4237	4	
8	3631 Bunker Hill Dr	\$347,000	3631	Bunker Hill Dr	06/13/2000	133	3933	4	
9	3651 Bunker Hill Dr	\$284,000	3651	Bunker Hill Dr	03/29/2000	103	2688	3	
10	3671 Bunker Hill Dr	\$290,000	3671	Bunker Hill Dr	01/01/2000	104	2790	3	
11	3700 Bunker Hill Dr	\$294,000	3700	Bunker Hill Dr	07/06/2000	46	3937	3	
12	3721 Bunker Hill Dr	\$235,000	3721	Bunker Hill Dr	09/08/1999	74	2504	3	
13	4020 Bunker Hill Dr	\$250,000	4020	Bunker Hill Dr	01/01/2000	10	1993	3	
14	4030 Bunker Hill Dr	\$290,000	4030	Bunker Hill Dr	07/31/2000	13	2894	3	
15	4141 Bunker Hill Dr	\$247,000	4141	Bunker Hill Dr	01/01/2000	34	3112	3	
16	4150 Bunker Hill Dr	\$232,000	4150	Bunker Hill Dr	02/03/2000	15	3127	3	
17	4151 Bunker Hill Dr	\$278,000	4151	Bunker Hill Dr	07/12/2000	66	2585	3	
18	4160 Bunker Hill Dr	\$222,000	4160	Bunker Hill Dr	01/01/2000	73	1740	3	
19	4161 Bunker Hill Dr	\$265,000	4161	Bunker Hill Dr	07/26/2000	55	2347	3	
20	4170 Bunker Hill Dr	\$300,000	4170	Bunker Hill Dr	08/25/2000	85	2939	3	
21	4171 Bunker Hill Dr	\$274,000	4171	Bunker Hill Dr	06/16/2000	85	2629	3	
22	4180 Bunker Hill Dr	\$265,000	4180	Bunker Hill Dr	08/07/2000	61	1622	3	
23	4201 Bunker Hill Dr	\$254,000	4201	Bunker Hill Dr	02/03/2000	1	2655	3	
24	4210 Bunker Hill Dr	\$262,000	4210	Bunker Hill Dr	02/04/2000	10	1765	3	

Command -> Analyze>Descriptive Statistics>Frequencies

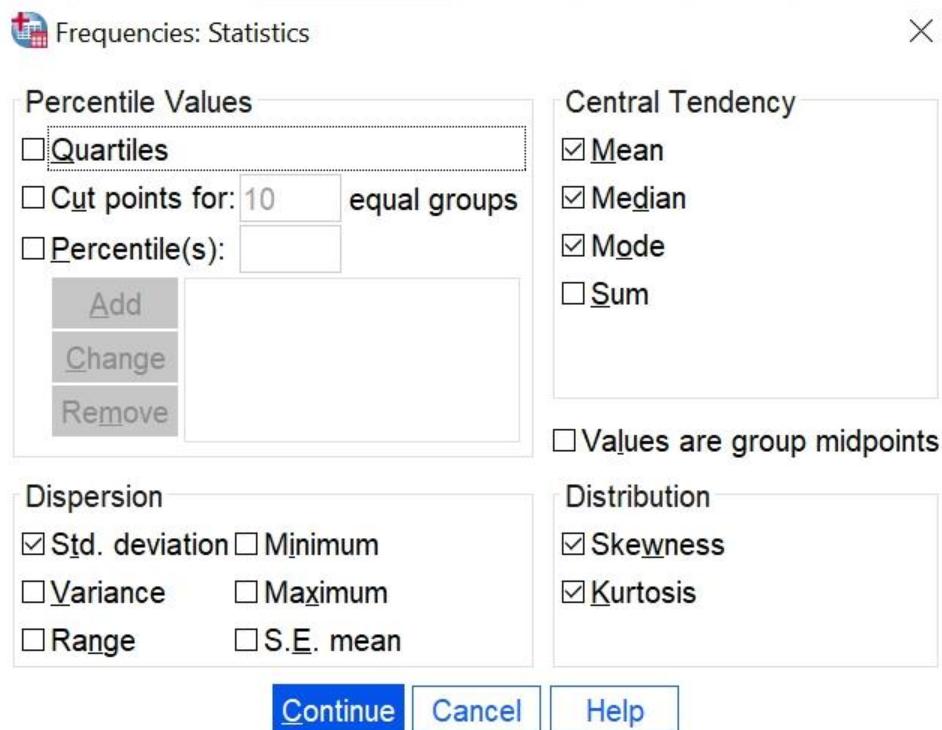


Command -> Click Purchase Price (Binned data)

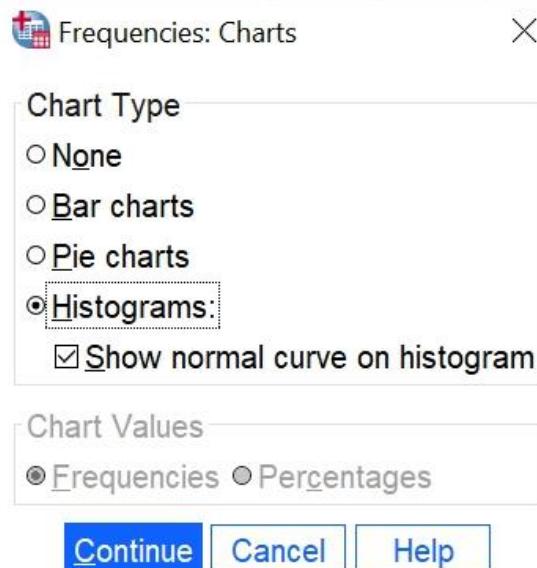


Command -> Click Statistics...

*fill the details as show below



Command -> Click Charts... *check “Histograms” along with “normal curve”



Command -> Click OK

Statistics

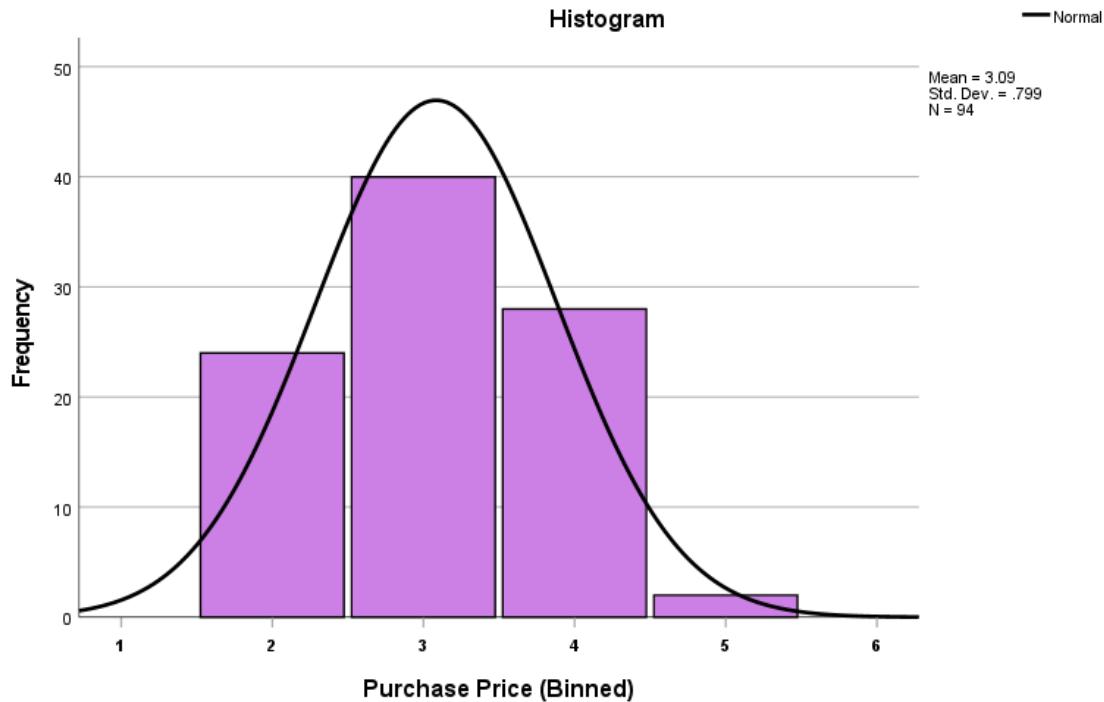
Purchase Price (Binned)

N	Valid	94
	Missing	0
Mean		3.09
Median		3.00
Mode		3
Std. Deviation		.799
Skewness		.103
Std. Error of Skewness		.249
Kurtosis		-.850
Std. Error of Kurtosis		.493

Purchase Price (Binned)

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	\$100,001 - \$200,000	24	25.5	25.5	25.5
	\$200,001 - \$300,000	40	42.6	42.6	68.1
	\$300,001 - \$400,000	28	29.8	29.8	97.9
	\$400,001 - \$500,000	2	2.1	2.1	100.0
	Total	94	100.0	100.0	

Histogram



Experiment - 5

Descriptive statistics like mean, mode, median, variance, skewness, kurtosis, etc.

THEORY

The various statistical constants are to be calculated and interpreted.

$$\text{sample mean } (\bar{x}) = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

Median

The point in the data set dividing it into two halves is called median.

Mode

It is the value which occurs with the greatest frequency (maybe more than one value).

Variance and standard deviation

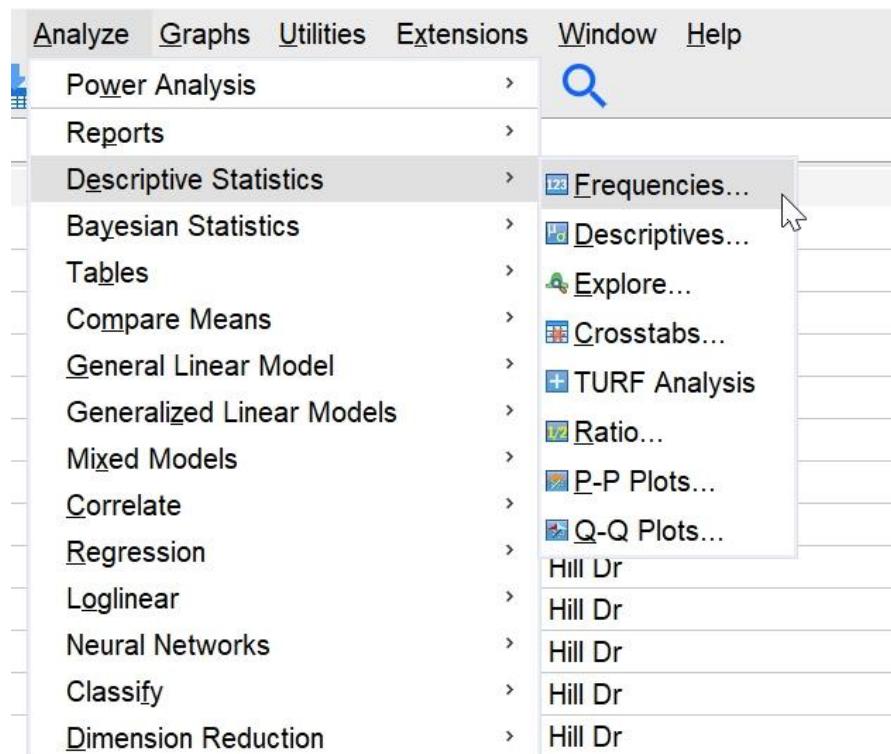
$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

STEPS

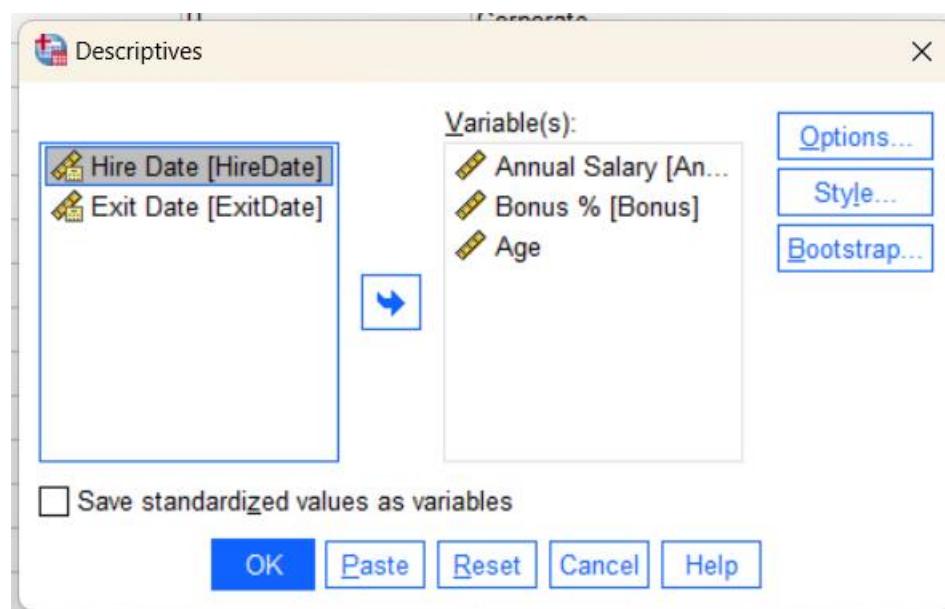
CASE – 1

- Statistical values when all types of cases together

Command -> Analyze>Descriptive Statistics>Descriptives...



Choose the required variable > option and check the desired boxes



Document output window

```
DESCRIPTIVES VARIABLES=AnnualSalary Bonus Age  
/STATISTICS=MEAN STDDEV MIN MAX.
```

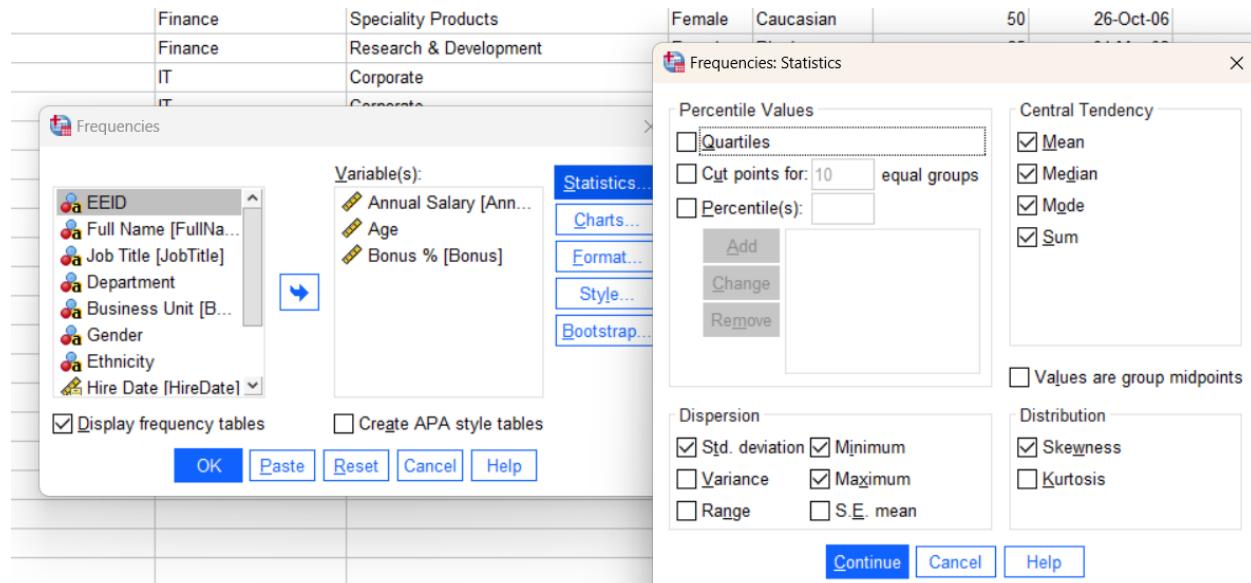
Descriptives

Descriptive Statistics

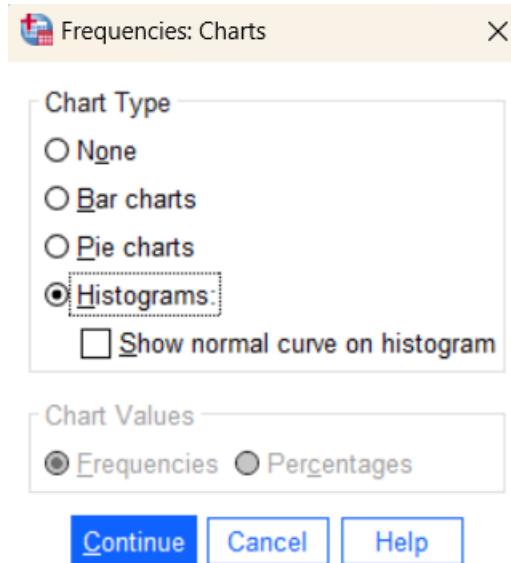
	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	19	41336.0	249270.0	124754.421	50210.0387
Bonus %	19	0.0%	30.0%	9.421%	9.3530%
Age	19	25	65	45.00	14.903
Valid N (listwise)	19				

```
FREQUENCIES VARIABLES=AnnualSalary Age Bonus  
/STATISTICS=STDDEV MINIMUM MAXIMUM MEAN MEDIAN MODE SUM SKEWNESS SESKEW  
/ORDER=ANALYSIS.
```

Command -> Analyze>Descriptive Statistics>Frequencies...



Command -> Choose the required variable > Histogram > Continue



Document output window

Frequencies

Statistics

		Annual Salary	Age	Bonus %
N	Valid	19	19	19
	Missing	0	0	0
Mean		124754.421	45.00	9.421%
Median		113527.000	50.00	9.000%
Mode		41336.0 ^a	27 ^a	0.0%
Std. Deviation		50210.0387	14.903	9.3530%
Skewness		.592	-.079	.658
Std. Error of Skewness		.524	.524	.524
Minimum		41336.0	25	0.0%
Maximum		249270.0	65	30.0%
Sum		2370334.0	855	179.0%

a. Multiple modes exist. The smallest value is shown

Frequency Table

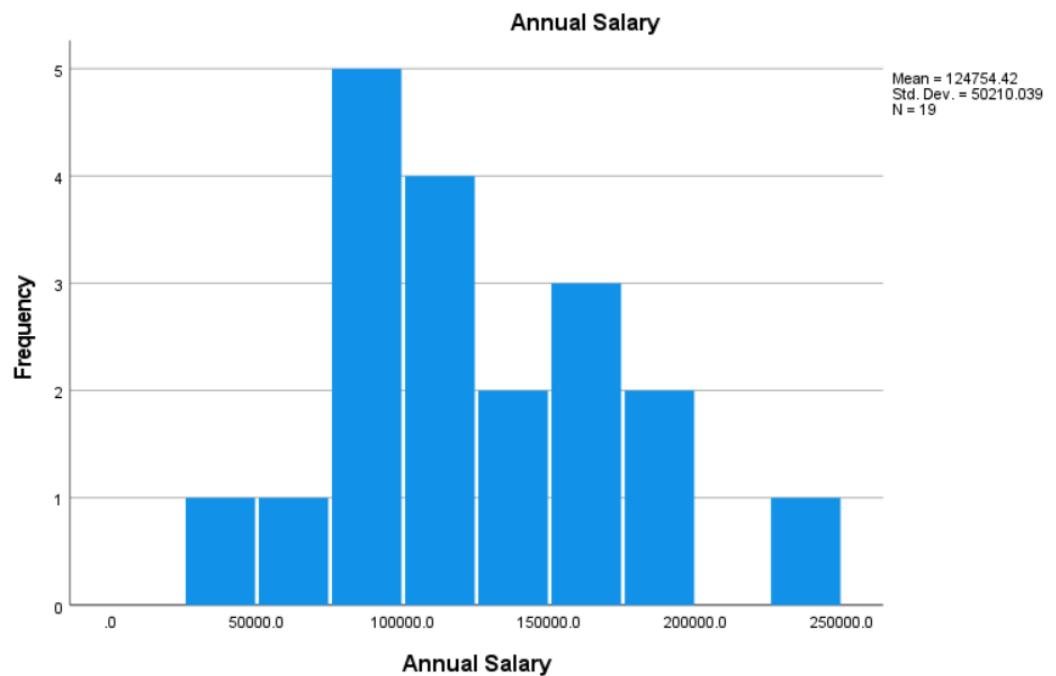
Annual Salary

	Frequency	Percent	Valid Percent	Cumulative Percent
Valid	41336.0	1	5.3	5.3
	50994.0	1	5.3	5.3
	77203.0	1	5.3	5.3
	84913.0	1	5.3	5.3
	95409.0	1	5.3	5.3
	97078.0	1	5.3	5.3
	99975.0	1	5.3	5.3
	105086.0	1	5.3	5.3
	109851.0	1	5.3	5.3
	113527.0	1	5.3	5.3
	119746.0	1	5.3	5.3
	141604.0	1	5.3	5.3
	146742.0	1	5.3	5.3
	154828.0	1	5.3	5.3
	157333.0	1	5.3	5.3
	163099.0	1	5.3	5.3
	175837.0	1	5.3	5.3
	186503.0	1	5.3	5.3
	249270.0	1	5.3	5.3
Total	19	100.0	100.0	

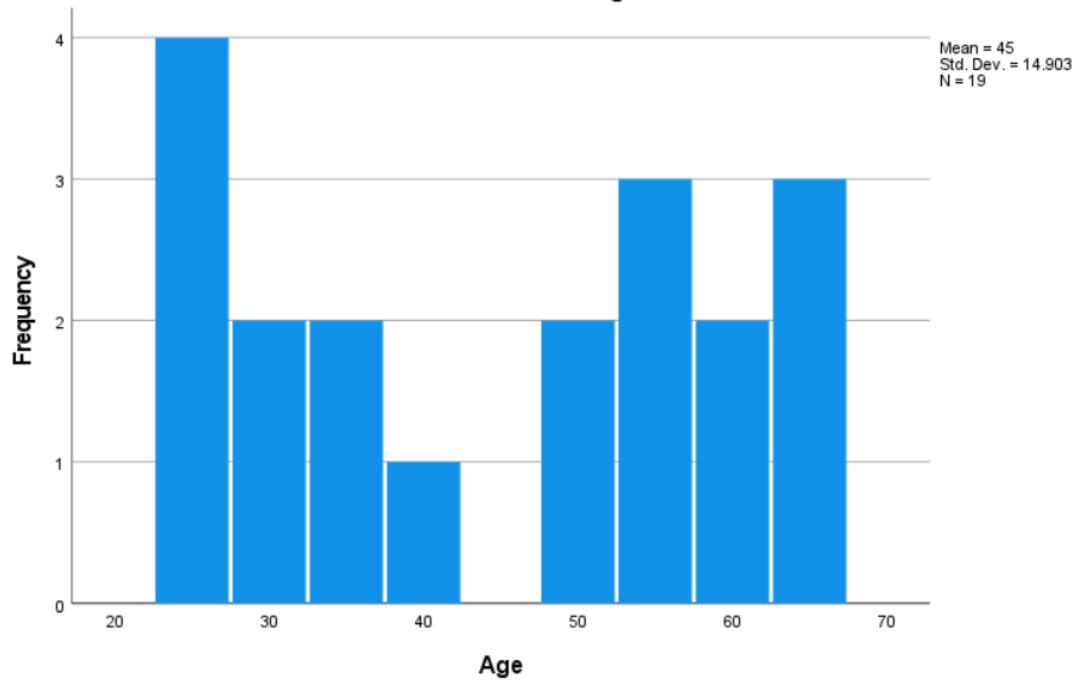
Age

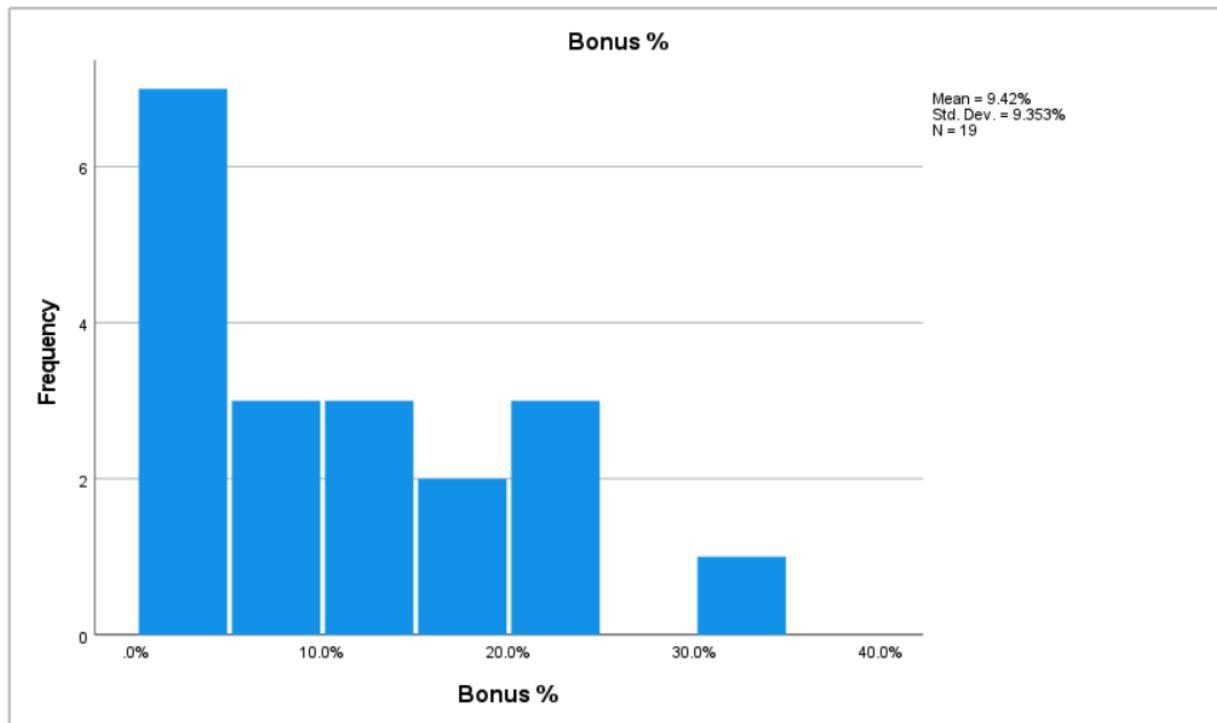
	Frequency	Percent	Valid Percent	Cumulative Percent
Valid	25	1	5.3	5.3
	26	1	5.3	10.5
	27	2	10.5	21.1
	29	1	5.3	26.3
	31	1	5.3	31.6
	34	1	5.3	36.8
	36	1	5.3	42.1
	41	1	5.3	47.4
	50	1	5.3	52.6
	51	1	5.3	57.9
	55	2	10.5	68.4
	57	1	5.3	73.7
	59	2	10.5	84.2
	64	2	10.5	94.7
	65	1	5.3	100.0
Total	19	100.0	100.0	

Histogram



Age



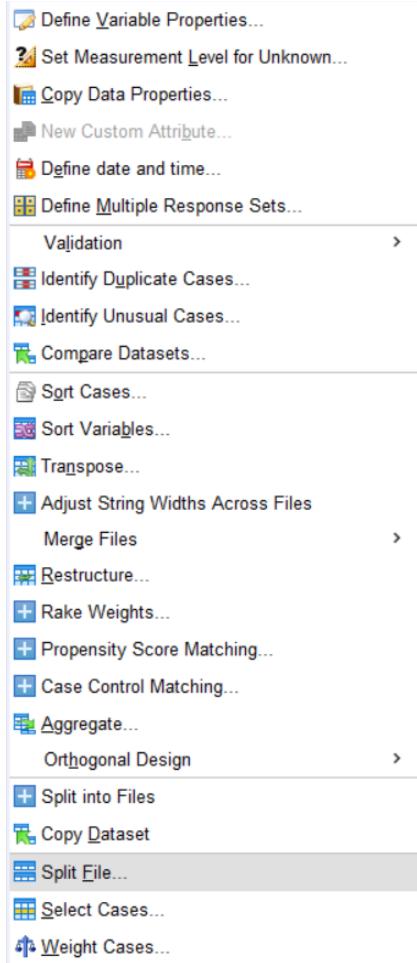


CASE – 2

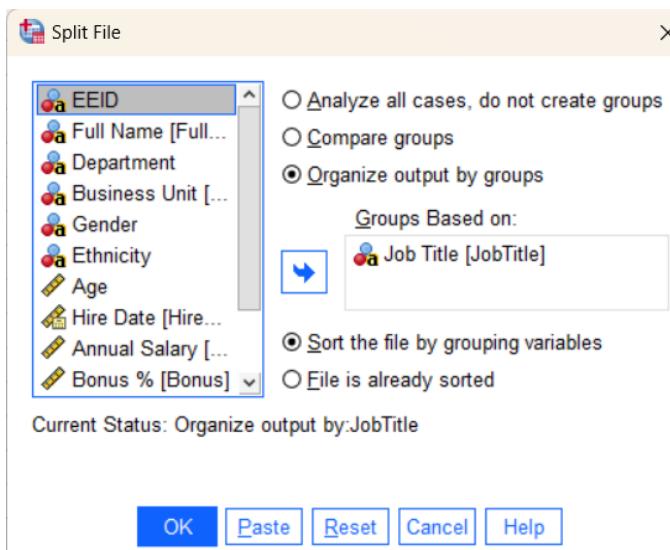
- Statistical values for all type of cases separately.

Suppose we are interested to have the statistical details with respect to all type of “Job Title” separately then follow

Command -> [Data>Split file](#)



Select the desired variable (e.g. "Job Title") for which the splitter information is to be obtained, also select organize output by groups.



Document output window

```
SORT CASES BY JobTitle.  
SPLIT FILE SEPARATE BY JobTitle.  
DESCRIPTIVES VARIABLES=AnnualSalary Bonus Age  
/STATISTICS=MEAN STDDEV MIN MAX.
```

► Descriptives

Job Title = Account Representative

Descriptive Statistics^a

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	1	50994.0	50994.0	50994.000	.
Bonus %	1	0.0%	0.0%	0.000%	.
Age	1	57	57	57.00	.
Valid N (listwise)	1				

a. Job Title = Account Representative

Job Title = Analyst

Descriptive Statistics^a

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	1	41336.0	41336.0	41336.000	.
Bonus %	1	0.0%	0.0%	0.000%	.
Age	1	25	25	25.00	.
Valid N (listwise)	1				

a. Job Title = Analyst

Job Title = Computer Systems Manager

Descriptive Statistics^a

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	1	84913.0	84913.0	84913.000	.
Bonus %	1	7.0%	7.0%	7.000%	.
Age	1	26	26	26.00	.
Valid N (listwise)	1				

a. Job Title = Computer Systems Manager

Job Title = Controls Engineer

Descriptive Statistics^a

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	1	109851.0	109851.0	109851.000	.
Bonus %	1	0.0%	0.0%	0.000%	.
Age	1	27	27	27.00	.
Valid N (listwise)	1				

a. Job Title = Controls Engineer

Double-click to activate

Job Title = Director

Descriptive Statistics^a

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	3	163099.0	186503.0	175146.333	11717.2765
Bonus %	3	20.0%	24.0%	21.333%	2.3094%
Age	3	50	65	59.67	8.386
Valid N (listwise)	3				

a. Job Title = Director

Job Title = Manager**Descriptive Statistics^a**

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	3	105086.0	119746.0	112786.333	7358.0120
Bonus %	3	6.0%	10.0%	8.333%	2.0817%
Age	3	27	59	38.33	17.926
Valid N (listwise)	3				

a. Job Title = Manager

Job Title = Sr. Analyst**Descriptive Statistics^a**

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	3	77203.0	97078.0	89896.667	11024.6664
Bonus %	3	0.0%	0.0%	0.000%	0.0000%
Age	3	31	55	40.00	13.077
Valid N (listwise)	3				

a. Job Title = Sr. Analyst

Job Title = Sr. Manger**Descriptive Statistics^a**

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	4	141604.0	157333.0	150126.750	7260.0641
Bonus %	4	10.0%	15.0%	13.250%	2.3629%
Age	4	36	64	51.50	11.676
Valid N (listwise)	4				

a. Job Title = Sr. Manger

Job Title = Technical Architect**Descriptive Statistics^a**

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	1	99975.0	99975.0	99975.000	.
Bonus %	1	0.0%	0.0%	0.000%	.
Age	1	59	59	59.00	.
Valid N (listwise)	1				

a. Job Title = Technical Architect

Job Title = Vice President**Descriptive Statistics^a**

	N	Minimum	Maximum	Mean	Std. Deviation
Annual Salary	1	249270.0	249270.0	249270.000	.
Bonus %	1	30.0%	30.0%	30.000%	.
Age	1	41	41	41.00	.
Valid N (listwise)	1				

a. Job Title = Vice President

CONCLUSION

- From this practical we have learnt to how to analyze data using “Descriptive Statistics like mean, mode, median, variance, skewness, kurtosis, etc.
- Also we have tackled two cases ([cases together](#) and [cases separately](#)) by using the “Data>Split File” features to analyze data together as well separated on some particular basis (here “Job Title”)

Experiment - 6

Correlation between two Random Variables

Theory and Methodology

As a measure of intensity or degree of linear relationship between two variables Karl Pearson coefficient of correlation denoted as r_{xy} or $r(X, Y)$ between two random variables X & Y is defined as

If $(x_i, y_i), i = 1, 2, \dots, n$ is the bivariate distribution, then

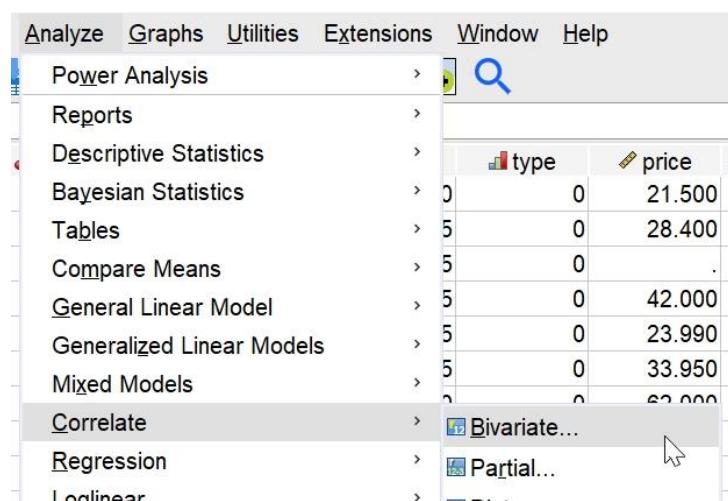
$$r_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i(x_i - \bar{x})^2} \sqrt{\sum_i(y_i - \bar{y})^2}} = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sqrt{\sum x^2 - \frac{(\sum x)^2}{n}} \sqrt{\sum y^2 - \frac{(\sum y)^2}{n}}}$$

Open the desired file in SPSS editor [File->Open->file path](#)

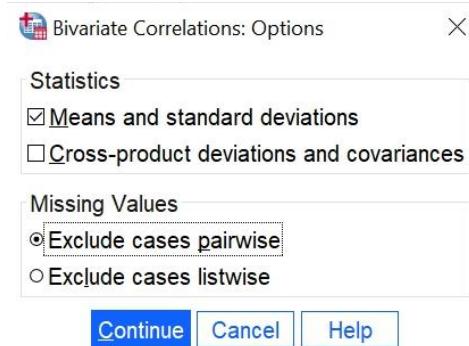
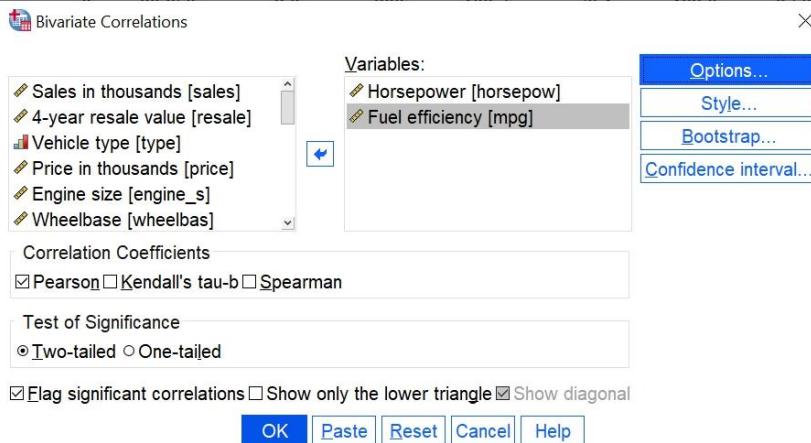
(car_sales.sav from repository is used and the variables Horsepower of the engine & Fuel efficiency is considered)

Case1. Missing values dealt as Exclude cases pairwise

Now follow the path [Analyze->Correlate->Bivariate](#)



Click option and check the desired boxes.



Document Output Window

Descriptive Statistics			
	Mean	Std. Deviation	N
Horsepower	185.95	56.700	156
Fuel efficiency	23.84	4.283	154

Correlations			
		Horsepower	Fuel efficiency
Horsepower	Pearson Correlation	1	-.611**
	Sig. (2-tailed)		<.001
	N	156	154
Fuel efficiency	Pearson Correlation	-.611**	1
	Sig. (2-tailed)	<.001	
	N	154	154

**. Correlation is significant at the 0.01 level (2-tailed).

Case2. Missing values dealt as Exclude cases list wise

Now follow the path [Analyze->Correlate->Bivariate](#)



Document Output Window

Descriptive Statistics

	Mean	Std. Deviation	N
Horsepower	185.66	57.006	154
Fuel efficiency	23.84	4.283	154

Correlations^b

		Horsepow er	Fuel efficiency
Horsepower	Pearson Correlation	1	-.611**
	Sig. (2-tailed)		<.001
Fuel efficiency	Pearson Correlation	-.611**	1
	Sig. (2-tailed)		<.001

**. Correlation is significant at the 0.01 level (2-tailed).

b. Listwise N=154

Case2. Missing values before calculation of correlation*

(*Note : if the missing values are dealt before finding the correlation by series mean, then the total number of cases has become 157)

Document Output Window

Descriptive Statistics

	Mean	Std. Deviation	N
Horsepower	186.74	57.379	157
Fuel efficiency	23.86	4.269	157

Correlations

		Horsepowe r	Fuel efficiency
Horsepower	Pearson Correlation	1	-.584**
	Sig. (2-tailed)		<.001
Fuel efficiency	N	157	157
	Pearson Correlation	-.584**	1
	Sig. (2-tailed)	<.001	
	N	157	157

**. Correlation is significant at the 0.01 level (2-tailed).

Experiment - 7

Regression Analysis

Theory and Methodology

Let (x_i, y_i) , $i = 1, 2, \dots, n$ be a bivariate sample, assuming one of them Independent say(x) and other dependent on first say(y), we predict the value of y by fitting a curve to data. If the curve is line then it is called line of regression and there is said to be a linear regression.

The curve fit for Y on x is called regression curve of y on x and similarly curve fit for X on y is called regression curve of x on y

Let the line of regression of Y on x be

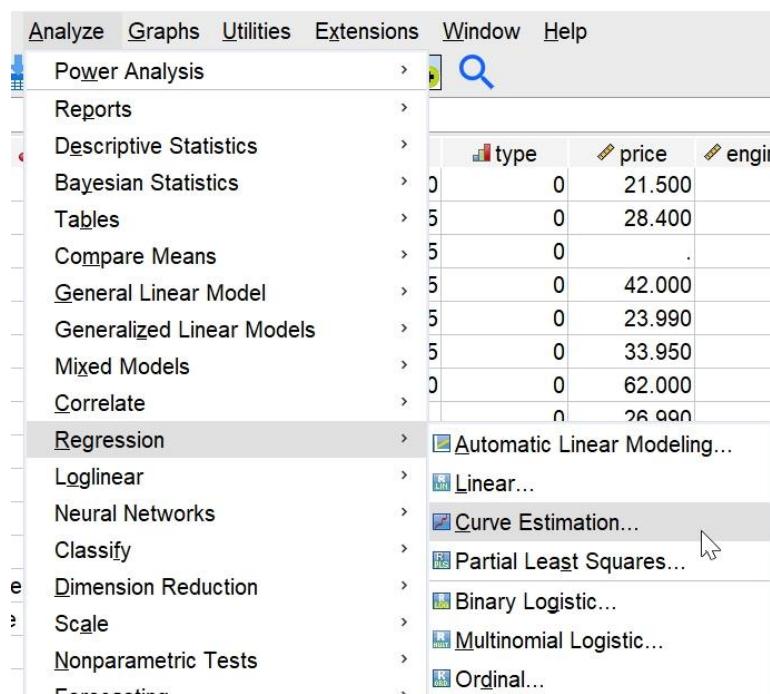
$$Y = a + bx$$

Open the desired file in SPSS editor [File->Open->file path](#)

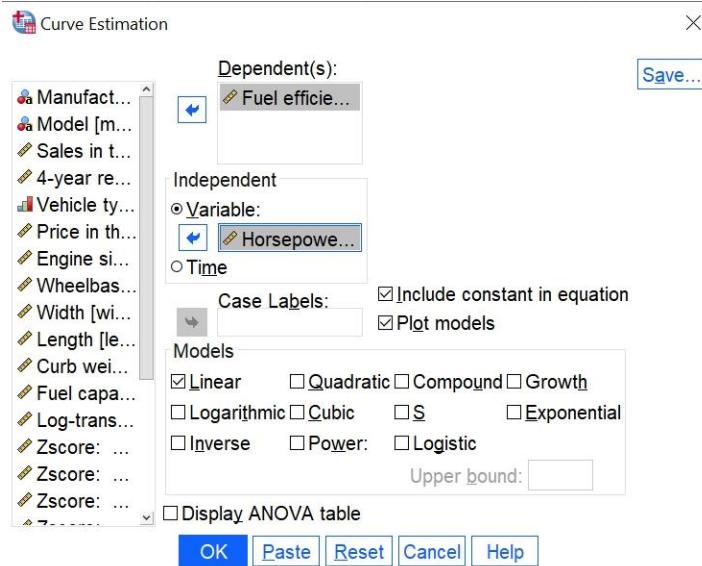
(car_sales.sav from repository is used and the variables Horsepower of the engine & Fuel efficiency is considered)

Case1. Using curve fitting path

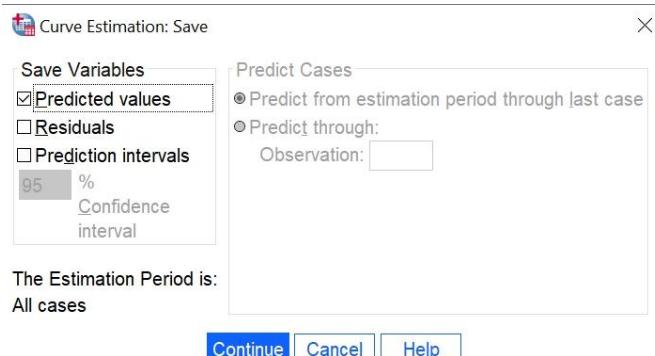
Now follow the path [Analyze->Regression->Curve Estimation](#)



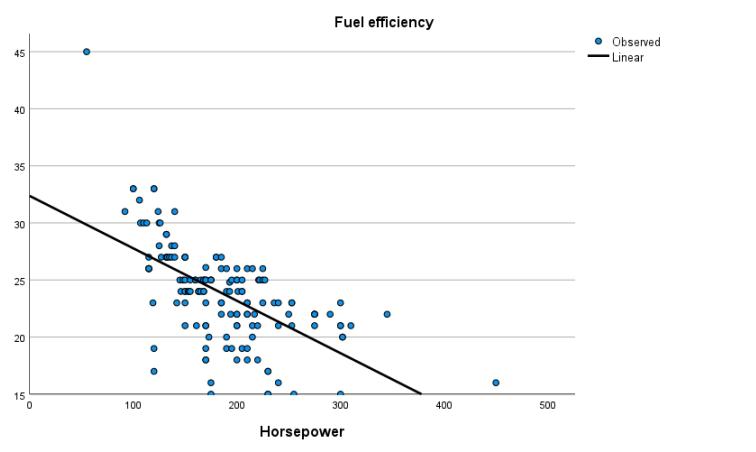
Choose dependent and independent variables and check the other desired boxes (Linear Model). Click save



Check the predicted values



Document Output Window



Model Description

Model Name	MOD_2
Dependent Variable	1
Equation	Linear
Independent Variable	Horsepower
Constant	Included
Variable Whose Values Label Observations in Plots	Unspecified

Case Processing Summary

	N
Total Cases	157
Excluded Cases ^a	3
Forecasted Cases	0
Newly Created Cases	0

a. Cases with a missing value in any variable are excluded from the analysis.

Variable Processing Summary

	Variables	
	Dependent Fuel efficiency	Independent Horsepower
Number of Positive Values	154	156
Number of Zeros	0	0
Number of Negative Values	0	0
Number of Missing Values	User-Missing System-Missing	0 3
		1

Model Summary and Parameter Estimates

Dependent Variable: Fuel efficiency

Equation	R Square	Model Summary			Parameter Estimates	
		F	df1	df2	Sig.	Constant
Linear	.374	90.743	1	152	<.001	32.371
						-.046

The independent variable is Horsepower.

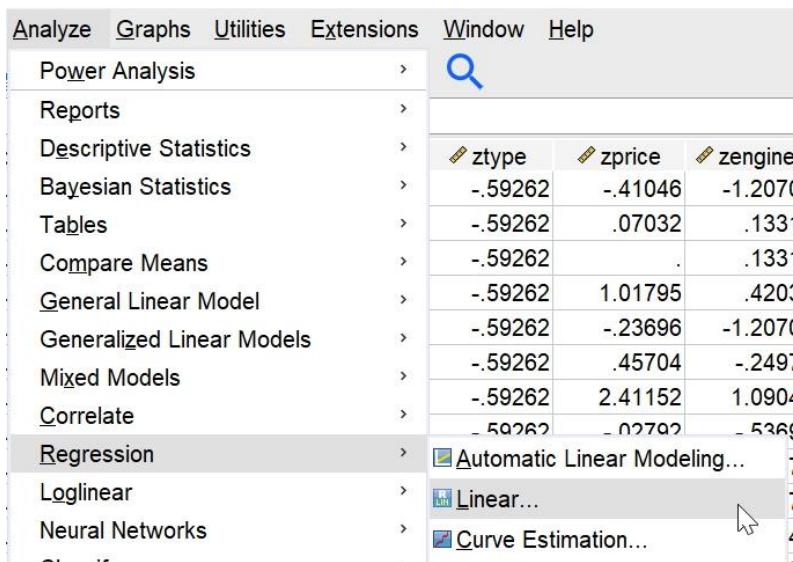
Also in data sheet a new variable for predicted values has been created.

FIT_1

25.94066
22.03627
22.03627
22.72528
25.48132
23.18462
18.13189
24.56264
23.50616
23.50616
24.33297
21.34727
22.95495
22.95495

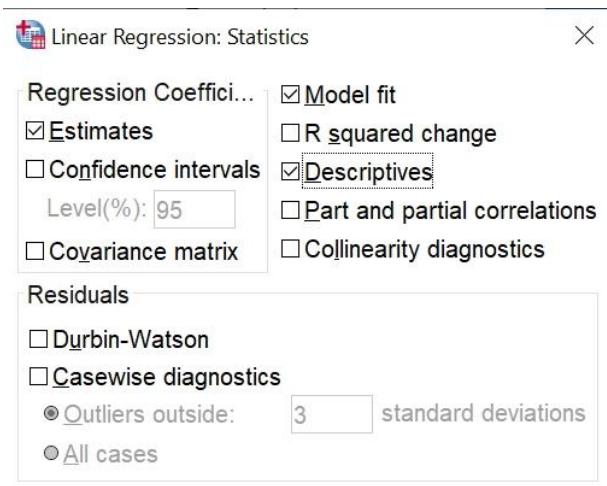
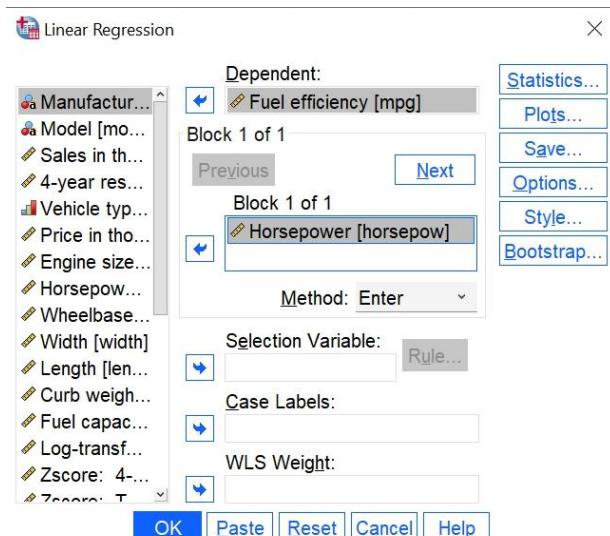
Case2. Using Linear Model Approach

Now follow the path **Analyze->Regression->Linear**



Choose dependent and independent variables and click statistics,

Check the desired boxes and continue



Click Save and check the desired boxes

Linear Regression: Save

Predicted Values	Residuals
<input checked="" type="checkbox"/> Unstandardized	<input checked="" type="checkbox"/> Unstandardized
<input type="checkbox"/> Standardized	<input type="checkbox"/> Standardized
<input type="checkbox"/> Adjusted	<input type="checkbox"/> Studentized
<input type="checkbox"/> S.E. of mean predictions	<input type="checkbox"/> Deleted
	<input type="checkbox"/> Studentized deleted
Distances	Influence Statistics
<input type="checkbox"/> Mahalanobis	<input type="checkbox"/> DfBetas
<input type="checkbox"/> Cook's	<input type="checkbox"/> Standardized DfBetas
<input type="checkbox"/> Leverage values	<input type="checkbox"/> DfFits
Prediction Intervals	<input type="checkbox"/> Standardized DfFits
<input type="checkbox"/> Mean	<input type="checkbox"/> Covariance ratios
<input type="checkbox"/> Individual	
Confidence Interval: 95 %	
Coefficient statistics	
<input type="checkbox"/> Create coefficient statistics	
<input checked="" type="radio"/> Create a new dataset	
Dataset name:	<input type="text"/>
<input type="radio"/> Write a new data file	
<input type="button" value="File..."/>	
Export model information to XML file	
<input type="text"/>	<input type="button" value="Browse..."/>
<input checked="" type="checkbox"/> Include the covariance matrix	

Click Option

Linear Regression: Options

Stepping Method Criteria
<input checked="" type="radio"/> Use probability of F
Entry: .05 Removal: .10
<input type="radio"/> Use F value
Entry: 3.84 Removal: 2.71
<input checked="" type="checkbox"/> Include constant in equation
Missing Values
<input type="radio"/> Exclude cases listwise
<input checked="" type="radio"/> Exclude cases pairwise
<input type="radio"/> Replace with mean

Document Output Window

Descriptive Statistics					
	Mean	Std. Deviation	N		
Fuel efficiency	23.84	4.283	154		
Horsepower	185.95	56.700	156		

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1049.053	1	1049.053	90.743 <.001 ^b
	Residual	1757.226	152	11.561	
	Total	2806.279	153		

a. Dependent Variable: Fuel efficiency
b. Predictors: (Constant), Horsepower

Coefficients ^a					
Model	Unstandardized Coefficients		Standardized Coefficients		Sig.
	B	Std. Error	Beta	t	
1	(Constant)	32.431	.942		34.421 <.001
	Horsepower	-.046	.005	-.611	-9.526 <.001

a. Dependent Variable: Fuel efficiency

Residuals Statistics ^a					
	Minimum	Maximum	Mean	Std. Deviation	N
Predicted Value	11.65	29.89	23.84	2.619	156
Residual	-9.889	15.109	-.014	3.389	154
Std. Predicted Value	-4.657	2.309	.000	1.000	156
Std. Residual	-2.908	4.444	-.004	.997	154

a. Dependent Variable: Fuel efficiency

Correlations			
	Fuel efficiency	Horsepower	
Pearson Correlation	Fuel efficiency	1.000	-.611
	Horsepower	-.611	1.000
Sig. (1-tailed)	Fuel efficiency	.	<.001
	Horsepower	.	.
N	Fuel efficiency	154	154
	Horsepower	154	156

Variables Entered/Removed^a

Model	Variables Entered	Variables Removed	Method
1	Horsepower ^b	.	Enter

a. Dependent Variable: Fuel efficiency

b. All requested variables entered.

Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.611 ^a	.374	.370	3.400

a. Predictors: (Constant), Horsepower

b. Dependent Variable: Fuel efficiency

Two new variables for predicted values (PRE_1) and difference between predicted and actual (RES_1) have been created.

PRE_1	RES_1
25.96548	2.03452
22.04006	2.95994
22.04006	3.95994
22.73278	-.73278
25.50367	1.49633
23.19460	-1.19460
18.11464	2.88536
24.58004	1.51996
23.51787	.48213
23.51787	1.28213
24.34913	.65087

Conclusion:

The Data set of Car sales from repository is used, and a regression line fitted between the variables (Horse power-independent variable and Fuel efficiency-dependent variable). the value of constant is.... and coefficient is Explain further the regression line properties exhibited in the case.

Precautions:

1. Variable Properties are to be defined carefully.
2. Missing values if any should be taken care of.

Experiment - 8

Hypothesis Testing ‘T’- test

Theory and Methodology

Null Hypothesis : Any hypothesis we wish to test and is denoted by H_0 (hypothesis of no difference)

Alternative Hypothesis : Any hypothesis which is complementary to null hypothesis is called and alternative hypothesis, denoted by H_1

‘t’ - test for single mean

If a random sample $x_i, i = 1, 2, \dots, n$ of size ‘n’ and is drawn from a normal population with mean μ then we have the statistic

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} , \quad \text{or} \quad t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n-1}}}$$

‘t’ – test for difference of means

$$t = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Paired ‘t’- test for difference of means

Here we consider the increments $d_i = x_i - y_i, i = 1, 2, \dots, n$ under the null hypothesis H_0 that increments are due to fluctuations of sampling, that is, drug is not responsible for these instruments;

The test statistic $t = \frac{\bar{d}}{s/\sqrt{n}}$ where

$$\bar{d} = \frac{\sum d_i}{n} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (d_i - \bar{d})^2 = \frac{1}{n-1} \left[\sum d_i^2 - \frac{(\sum d_i)^2}{n} \right] \quad \text{with } n-1 \text{ df}$$

Case1. Test for single mean

Open the desired file in SPSS editor [File > Open > file path](#)

(car_sales.sav from repository is used and the variables Fuel efficiency is considered)

H_0 : Average is 23 km/lt or ($= 23$ km per lt) H_1 : $\mu \neq 23$ km per lt)

The screenshot shows the SPSS menu bar with 'Analyze', 'Graphs', 'Utilities', 'Extensions', 'Window', and 'Help'. Under 'Analyze', 'Compare Means' is selected, which has 'One-Sample T Test...' highlighted. Below the menu is a table with columns: type, price, engine_s, and hp.

One-Sample T Test Dialog:

- Test Variable(s):** Fuel efficiency [mpg]
- Test Value:** 23
- Options...** button is visible.
- Buttons at the bottom:** OK, Paste, Reset, Cancel, Help.

One-Sample T Test: Options Dialog:

- Confidence Interval Percentage:** 95%
- Missing Values:**
 - Exclude cases analysis by analysis
 - Exclude cases listwise
- Buttons at the bottom:** Continue, Cancel, Help.

Document Output View

One-Sample Statistics

	N	Mean	Std. Deviation	Std. Error Mean
Fuel efficiency	154	23.84	4.283	.345

One-Sample Test

Test Value = 23					
t	df	Sig. (2-tailed)	Mean Difference	95% Confidence Interval of the Difference	
				Lower	Upper
Fuel efficiency	2.444	153	.016	.844	.16 1.53

One-Sample Effect Sizes

	Standardizer ^a	Point Estimate	95% Confidence Interval	
			Lower	Upper
Fuel efficiency	Cohen's d	4.283	.197	.037 .356
	Hedges' correction	4.304	.196	.037 .354

a. The denominator used in estimating the effect sizes.

Cohen's d uses the sample standard deviation.

Hedges' correction uses the sample standard deviation, plus a correction factor.

*Note the p value is 0.016 which is less than the assigned value 0.05 hence null hypothesis is rejected, i.e., claim $\mu = 23 \text{ km per lt}$ cannot be accepted.

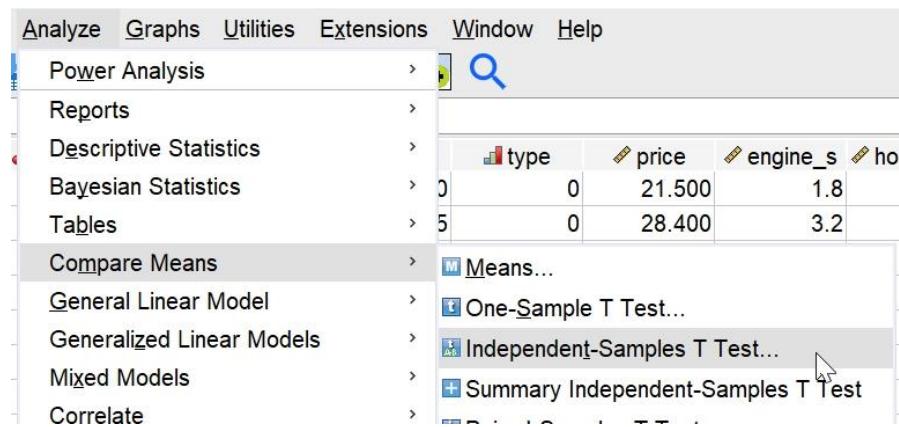
Case2. Test for difference of means

Open the desired file in SPSS editor **File > Open > file path**

(car_sales.sav from repository is used and the variables Fuel efficiency is considered)

$H_0 : \text{Fuel efficiency of Automobiles} = \text{Fuel efficiency of Trucks} (A = \mu_T)$

Now follow the path **Analyze>Compare Means>Independent Sample Test**



Transfer the the desired variable to test variable then transfer the variable w.r.t. which the means are to be tested (in present case Automobiles & truck are considered)

The screenshot shows the 'Independent-Samples T Test' dialog box and its 'Define Groups' sub-dialog.

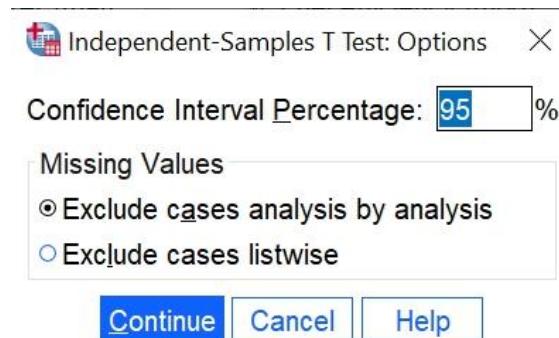
Main Dialog: Independent-Samples T Test

- Test Variable(s):** Fuel efficiency [mpg]
- Grouping Variable:** type(0 1)
- Options...** button
- Bootstrap...** button
- OK**, **Paste**, **Reset**, **Cancel**, **Help** buttons

Sub-DIALOG: Define Groups

- Use specified values** radio button selected
- Group 1:** 0
- Group 2:** 1
- Cut point:** (empty input field)
- Continue**, **Cancel**, **Help** buttons

Enter the desired level of significance



Document Output View

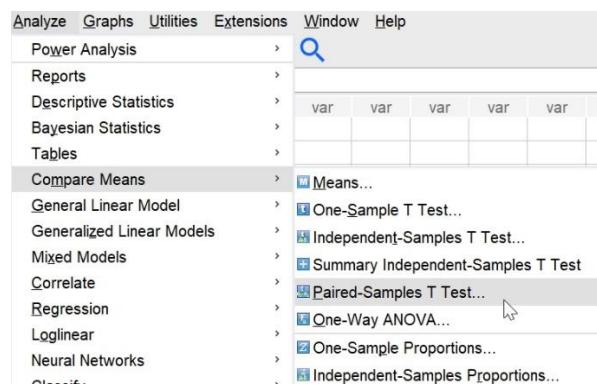
Group Statistics								
	Vehicle type	N	Mean	Std. Deviation	Std. Error Mean			
Fuel efficiency	Automobile	114	25.30	3.646	.341			
	Truck	40	19.70	3.107	.491			

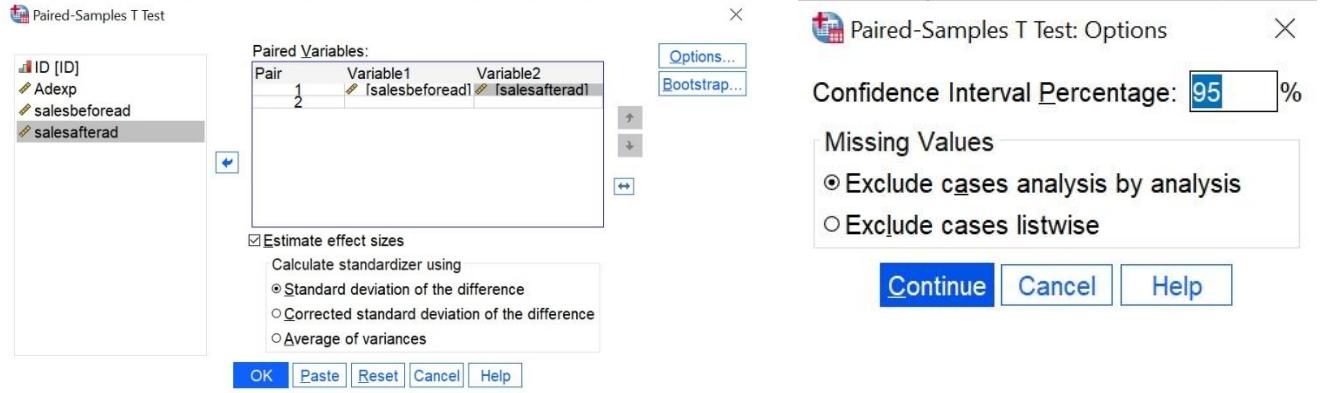
Independent Samples Test								
	Levene's Test for Equality of Variances			df	Sig. (2-tailed)	t-test for Equality of Means		
	F	Sig.	t			Mean Difference	Std. Error Difference	95% Confidence Interval of the Difference
Fuel efficiency	Equal variances assumed	.004	.948	152	<.001	5.597	.646	4.321 6.874
	Equal variances not assumed			9.356	79.405	<.001	.598	4.407 6.788

*(Levene's test is to decide the equivalence of the variances) in the present case in first row reading of p-value is $0.948 > 0.05$ hence equivalence of variances can be accepted (Now we will read all information from first row). Also for the t-test p-value is almost zero <0.05 hence null hypothesis is rejected.

Case3. Paired 't'-test for difference of means

Open the desired file in SPSS editor [File > Open > file path](#)





$H_0: \text{avg sales before ad} = \text{avg sales after ad}$ ($b = \mu a$) or $\{\mu b - \mu a = 0\}$

Document Output View

Paired Samples Statistics

		Mean	N	Std. Deviation	Std. Error Mean
Pair 1	salesbeforead	72.9974	20	2.75234	.61544
	salesafterad	74.1420	20	2.74647	.61413

Paired Samples Correlations

		N	Correlation	Sig.
Pair 1	salesbeforead & salesafterad	20	.990	<.001

Paired Samples Test

	Mean	Std. Deviation	Std. Error Mean	95% Confidence Interval of the Difference				df	Sig. (2-tailed)
				Lower	Upper	t			
Pair 1	salesbeforead - salesafterad	-1.14459	.39674	.08871	-1.33027	-.95891	-12.902	19	<.001

Precautions:

- Variable Properties are to be defined carefully.
- Missing values if any should be taken care of.
- Null hypothesis is defined carefully and level of significance chosen appropriately.

Experiment – 9

Chi Square Test

Theory and Methodology

Chi Square test in SPSS is test for association (or independence) of two nominal variables.

Case1. Chi-squared test of independence in contingency tables

A cell containing information of two characteristics resulting into a bivariate data. The table displaying data is called contingency table. Horizontal and vertical display of data is referred as $r \times c$ table.

Our interest is to see whether two characteristics are independent, thus

H_0 : Two characteristics are independent.

H_1 : they are dependent.

O_{ij} : observed frequency of i^{th} row and j^{th} column respectively.

e_{ij} : expected frequency of i^{th} row and j^{th} column respectively.

Then

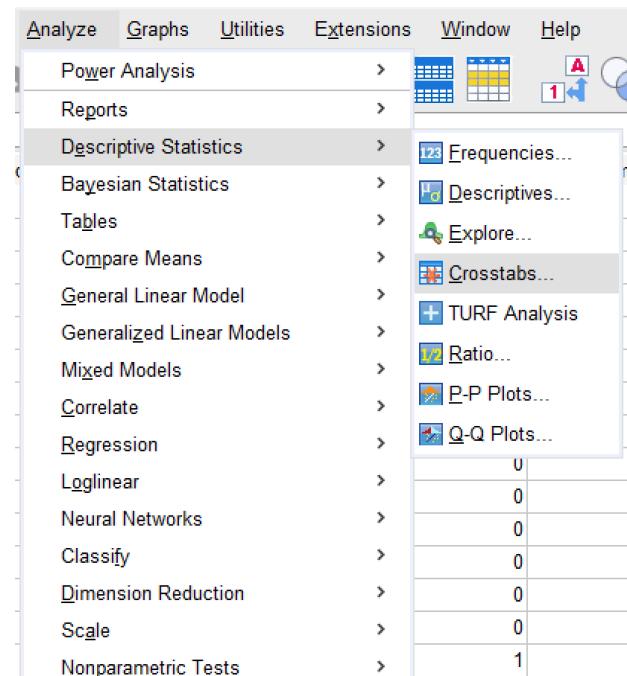
$$\text{So, } \chi^2 = \sum_i \sum_j \frac{(O_{ij} - e_{ij})^2}{e_{ij}} \quad \text{with } (r - 1) \times (c - 1) \text{ degree of freedom (df).}$$

Open the desired file in SPSS editor [File > Open > file path](#)

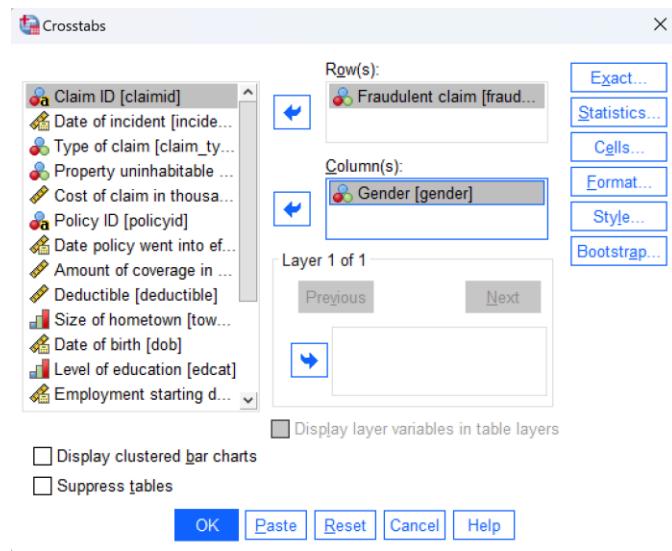
(insurance_claims.sav from repository is used and Fraudulent claim & Gender variables are considered for their association)

H_0 : Fraudulent claims are independent of Gender

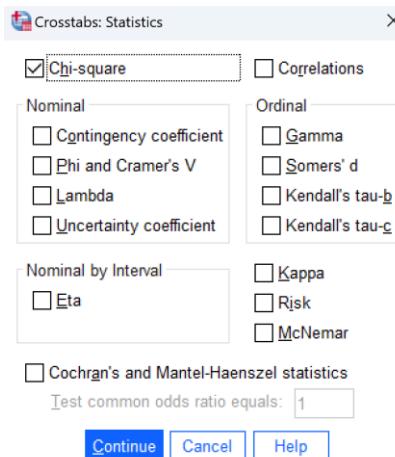
H_1 : Fraudulent claim is Gender dependent



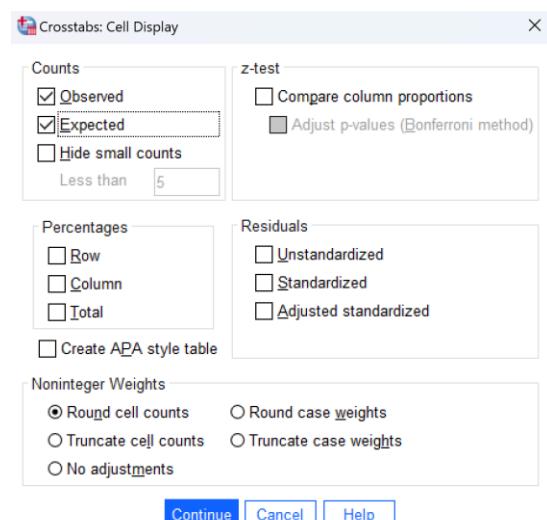
Transfer the desired variables in rows & column box (their order does not matter).



Then click Statistics tab and check Chi-square box.



click continue. Then click Cells button



Document Output View

→ Crosstabs

Case Processing Summary

	Valid		Cases Missing		Total	
	N	Percent	N	Percent	N	Percent
Fraudulent claim *	4415	100.0%	0	0.0%	4415	100.0%
Gender						

Fraudulent claim * Gender Crosstabulation

		Gender		Total
		Male	Female	
Fraudulent claim	No	Count	1964	1988
		Expected Count	1946.9	2005.1
	Yes	Count	211	252
		Expected Count	228.1	234.9
Total		Count	2175	2240
		Expected Count	2175.0	2240.0

Chi-Square Tests

	Value	df	Asymptotic Significance (2-sided)	Exact Sig. (2-sided)	Exact Sig. (1-sided)
Pearson Chi-Square	2.820 ^a	1	.093		
Continuity Correction ^b	2.657	1	.103		
Likelihood Ratio	2.824	1	.093		
Fisher's Exact Test				.095	.051
Linear-by-Linear Association	2.819	1	.093		
N of Valid Cases	4415				

a. 0 cells (0.0%) have expected count less than 5. The minimum expected count is 228.09.

b. Computed only for a 2x2 table

The Pearson Chi-Square p-value 0.093 is greater than 0.05, therefore Null Hypothesis is not rejected. It is concluded that fraudulent claim is not gender biased.

Case2. Chi-square test for goodness of fit

If a population has a specified theoretical distribution. The test is based on how good a fit we have between the frequency of occurrence of observations in an observed sample and the expected frequencies obtained from the hypothesized distribution.

A goodness of fit test between observed and expected frequencies is based on the quantity

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - e_i)^2}{e_i}$$

O_i : Observed Frequency of the cell

e_i : Expected Frequency of the cell

Where χ^2 is a value of a random variable whose sampling distribution is approximated very closely by the Chi square distribution with $v = k - 1$ degree of freedom (df).

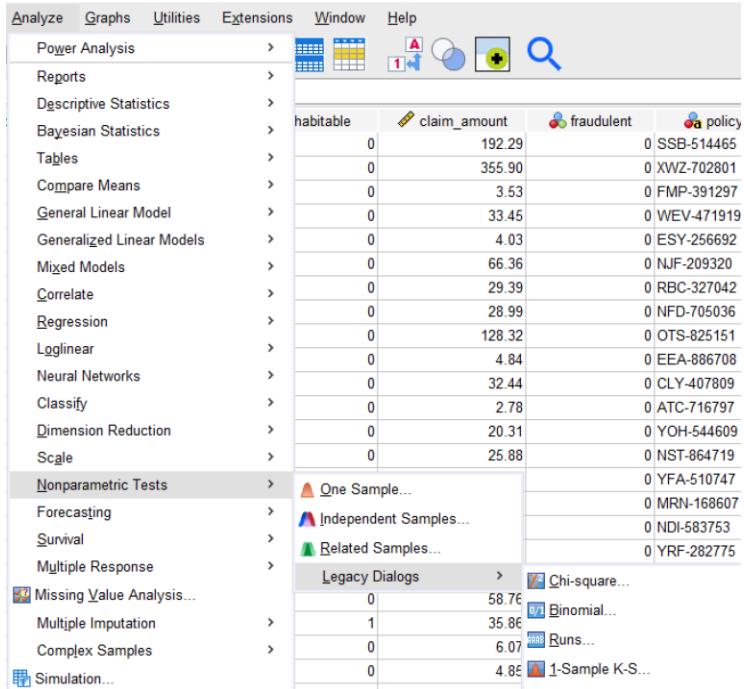
Critical values

If the observed frequencies differ considerably from the expected frequencies, the χ^2 value will be large and fit is poor. A good fit leads to the acceptance of H_0 , whereas a poor fit leads to rejection.

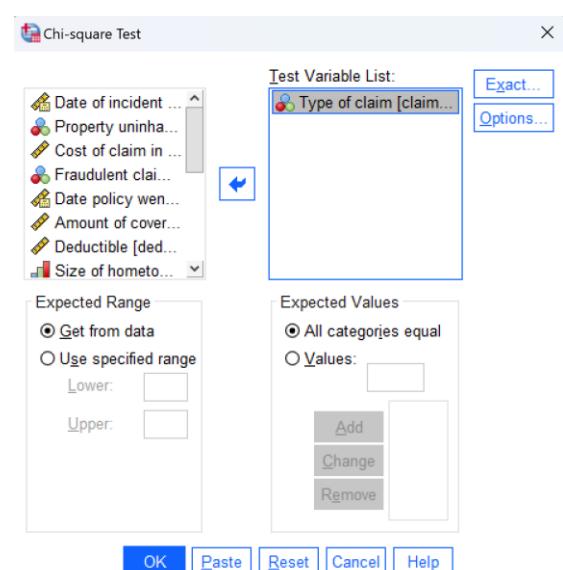
H_0 : Claims are equally distributed

H_1 : It is not equally distributed

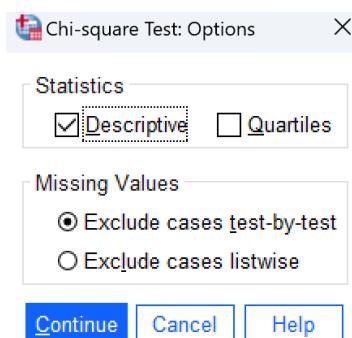
Now follow the path **Analyze > Nonparametric Tests > Legacy dialogs > Chi Square**



Transfer the desired variable to Test variable list box



Click **option** and check the desired box (Descriptive under statistics)



Document Output View

→ NPar Tests

Descriptive Statistics

	N	Mean	Std. Deviation	Minimum	Maximum
Type of claim	4415	3.06	1.535	1	5

Chi-Square Test

Frequencies

Type of claim

	Observed N	Expected N	Residual
Wind/Hail	1054	883.0	171.0
Water damage	627	883.0	-256.0
Fire/Smoke	1039	883.0	156.0
Contamination	404	883.0	-479.0
Theft/Vandalism	1291	883.0	408.0
Total	4415		

Test Statistics

Type of claim

Chi-Square	583.259 ^a
df	4
Asymp. Sig.	<.001

a. 0 cells (0.0%) have expected frequencies less than 5. The minimum expected cell frequency is 883.0.

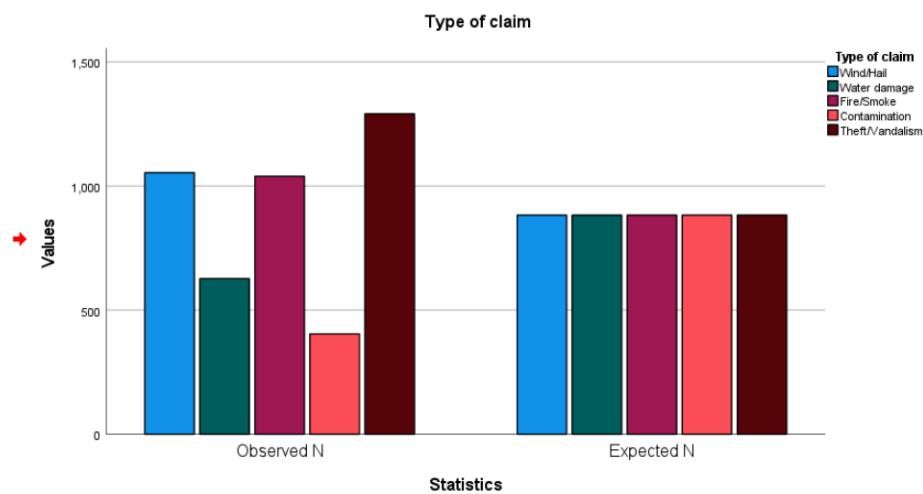
The Test Statistics shows the p-value almost zero which is less than 0.05, therefore the null hypothesis is rejected in favor of alternate, i.e., type of claim are not equally distributed. This can be visualized through graph also.

Double click the frequency table to bring into editable mode

Select both columns of observed and expected frequency and right-click
 >Create Graph>Bar

Type of claim			
	Observed N	Expected N	Residual
Wind/Hail	1054	883.0	171.0
Water damage	627	883.0	-256.0
Fire/Smoke	1039	883.0	156.0
Contamination	404	883.0	-479.0
Theft/Vandalism	1291	883.0	408.0
Total	4415		

- [Cut](#) Ctrl+X
- [Copy](#) Ctrl+C
- [Paste](#) Ctrl+V
- [Delete](#) Delete
- [Select Table](#)
- [Select cells with similar significance](#)
- [Sort Rows](#) >
- [Create Graph](#) > Bar
- [Table Properties...](#)



Conclusion:

Give information of data set used for testing. Also mentioning the p-value write the decision about rejection or non-rejection of null hypothesis.

Precautions:

1. Variable Properties are to be defined carefully.
2. Missing values if any should be taken care of.
3. Null hypothesis is defined carefully and level of significance chosen appropriately.

Experiment - 10

ANOVA One-Way

Theory and Methodology

Let k random samples of size $n_1, n_2 \dots n_k$ are selected from normal populations with mean $\mu_1, \mu_2 \dots \mu_k$, and common variance σ^2 . And let the means are depending on one factor then One Way ANOVA is applied

$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ and H_1 : at least two means are different.

Procedure: the total variation can be attributed to variation within the class and variation between the class.

Total sum of squares = Sum of square within the class + sum of square between the class

$$TSS = SSB + SSW$$

The degree of freedom (df) for TSS is $(N - 1)$, for SSB is $(k - 1)$, & for SSW is $(N - k)$.

Now mean square (MS) for each variation is obtained by dividing each sum of square by its respective degree of freedom (df), i.e.,

$$MSS = \frac{TSS}{N - 1}, \quad MSB = \frac{SSB}{k - 1}, \quad MSW = \frac{SSW}{N - k}$$

The value of F - Statistic will be

$$F = \frac{\sum_i \sum_j n_i (x_i - \bar{x})^2}{\sum_i \sum_j (x_{ij} - \bar{x})^2} \times \frac{N - k}{k - 1}, \quad v_1 = k - 1, \quad v_2 = N - k$$

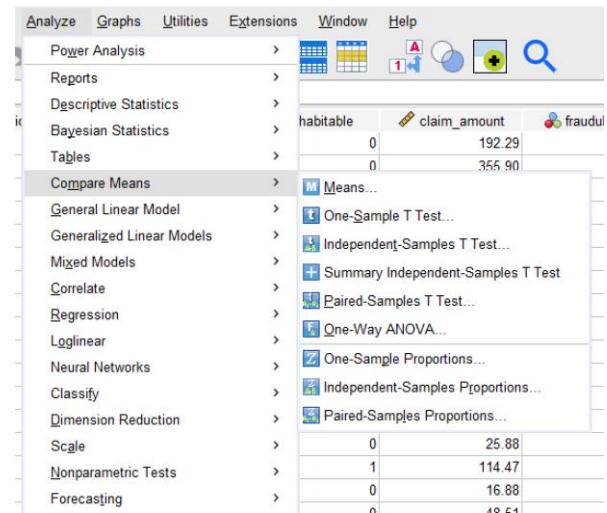
Open the desired file in SPSS editor [File > Open > file path](#)

(insurance_claims.sav from repository is used and claim amount & Type of claim variables are considered) (see the data 5-types of claims are there)

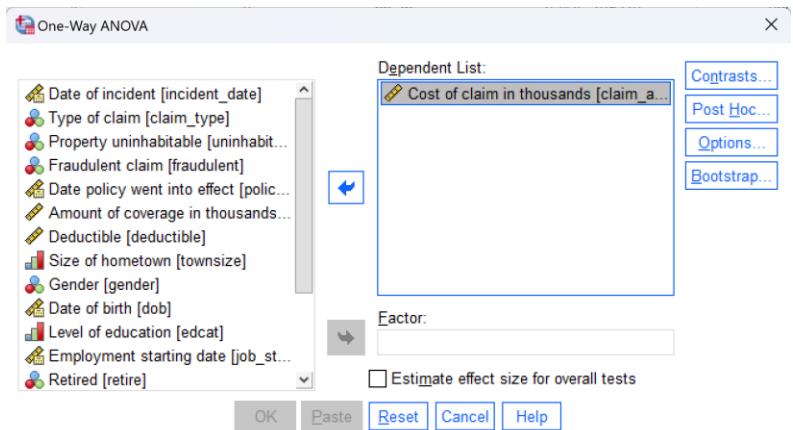
H_0 : average claim amount in each category is same. i.e., $\mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$

H_1 : At least one of these differs

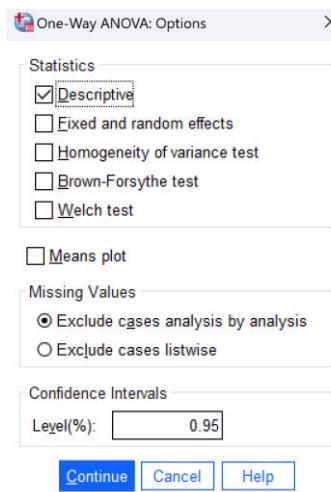
Now follow the path [Analyze > Compare Means](#)



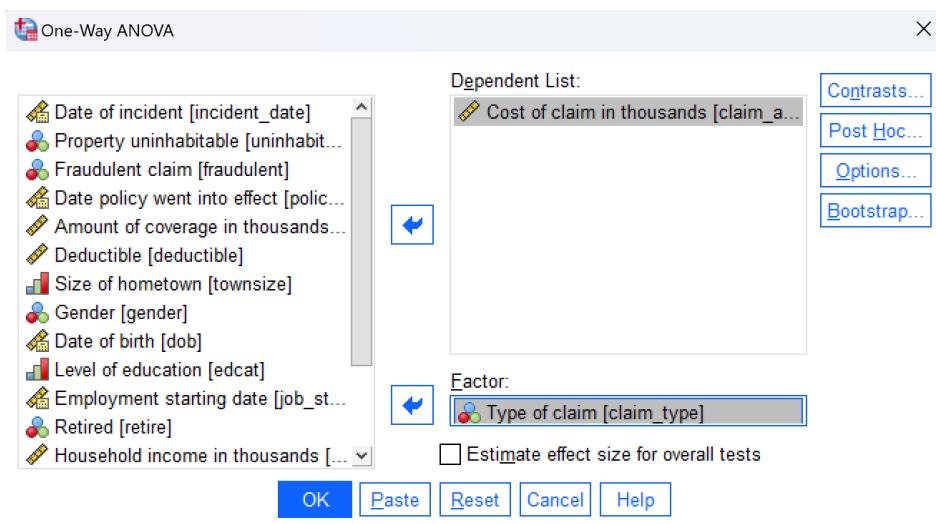
Transfer Cost of claim to dependent list and Type of claim to factor list



Click **Option** and select the Descriptive to display statistics



Continue & click OK.



Document Output View

→ Oneway

Descriptives

Cost of claim in thousands

	N	Mean	Std. Deviation	Std. Error	95% Confidence Interval for Mean			Minimum	Maximum
					Lower Bound	Upper Bound			
Wind/Hail	1054	16.7819	20.51640	.63195	15.5419	18.0219		1.46	184.38
Water damage	627	35.2896	47.55906	1.89933	31.5598	39.0195		2.31	504.94
Fire/Smoke	1039	171.5793	204.53371	6.34537	159.1281	184.0305		11.59	1635.00
Contamination	404	202.2067	227.66948	11.32698	179.9393	224.4740		18.69	1662.00
Theft/Vandalism	1291	17.4795	25.27650	.70348	16.0994	18.8596		1.57	576.91
Total	4415	73.0109	144.40137	2.17323	68.7503	77.2715		1.46	1662.00

ANOVA

Cost of claim in thousands

	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	25043718.32	4	6260929.579	412.125	<.001
Within Groups	66995929.16	4410	15191.821		
Total	92039647.48	4414			

Since the p-value of F-statistics is almost zero which is less than 0.05, hence the Null Hypothesis is rejected, i.e., at least one of the mean is different.

Conclusion:

Give information of data set used for testing. Also mentioning the p-value write the decision about rejection or non-rejection of null hypothesis.

Precautions:

1. Variable Properties are to be defined carefully.
2. Missing values if any should be taken care of.
3. Null hypothesis is defined carefully and level of significance chosen appropriately.