
BUSINESS INTELLIGENCE & DATA ANALYSIS

PRACTICALS

TY BSc.IT SEMESTER 6

NAME: ALISHA SHAIKH

ROLL NO: B20222507

**COLLEGE: SASMIRA'S INSTITUTE OF
COMMERCE & SCIENCE**

YEAR: 2025

INDEX

<u>Sr. No</u>	<u>Name</u>	<u>Date</u>	<u>Signature</u>
1	Perform the analysis for the following:		
a	Import the data warehouse data in Microsoft Excel and create the Pivot table and Pivot Chart.		
b	Import the cube in Microsoft Excel and create the Pivot table and Pivot Chart to perform data analysis.		
2	Apply the what – if Analysis for data visualization. Design and generate necessary reports based on the data warehouse data. Use Excel.		
3	Perform the data classification using classification algorithm using R/Python.		
4	Perform the data clustering using clustering algorithm using R/Python.		
5	Perform the Linear regression on the given data warehouse data using R/Python.		
6	Perform the logistic regression on the given data warehouse data using R/Python.		
7	Write a Python program to read data from a CSV file, perform simple data analysis, and generate basic insights. (Use Pandas is a Python library).		
8	Create the Data staging area for the selected database using SQL.		

Practical No.1

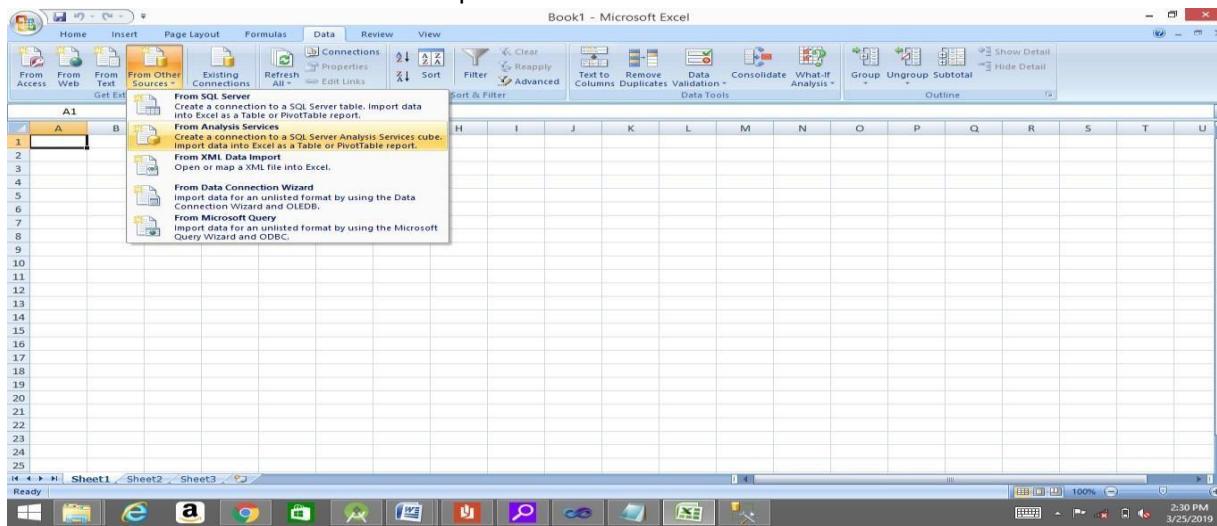
- a. Import the data warehouse data in Microsoft Excel and create the Pivot table and Pivot Chart.
 b. Import the cube in Microsoft Excel and create the Pivot table and Pivot Chart to Perform data analysis.

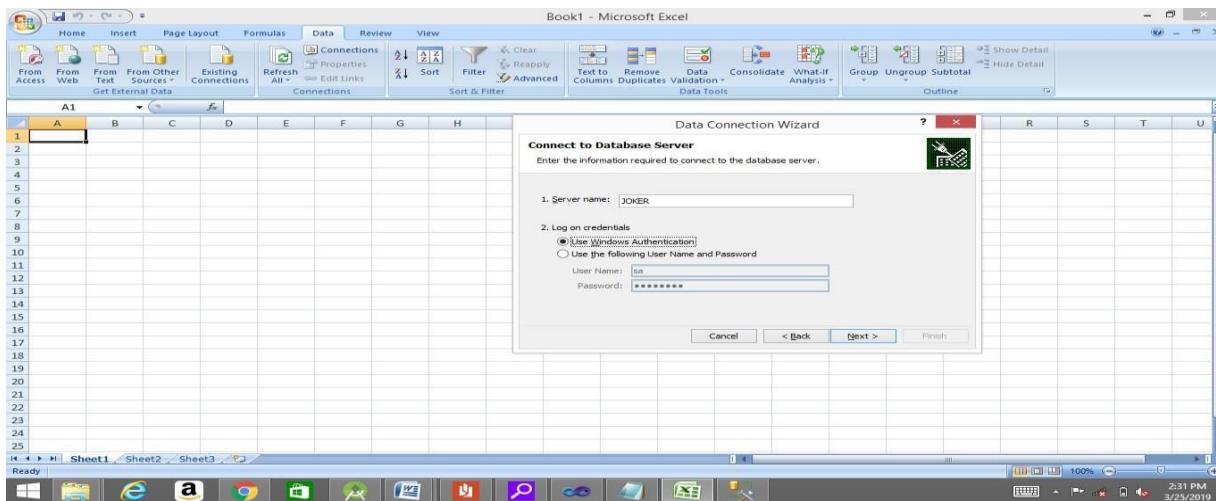
a. Import the data warehouse data in Microsoft Excel and create the Pivot table and Pivot Chart.

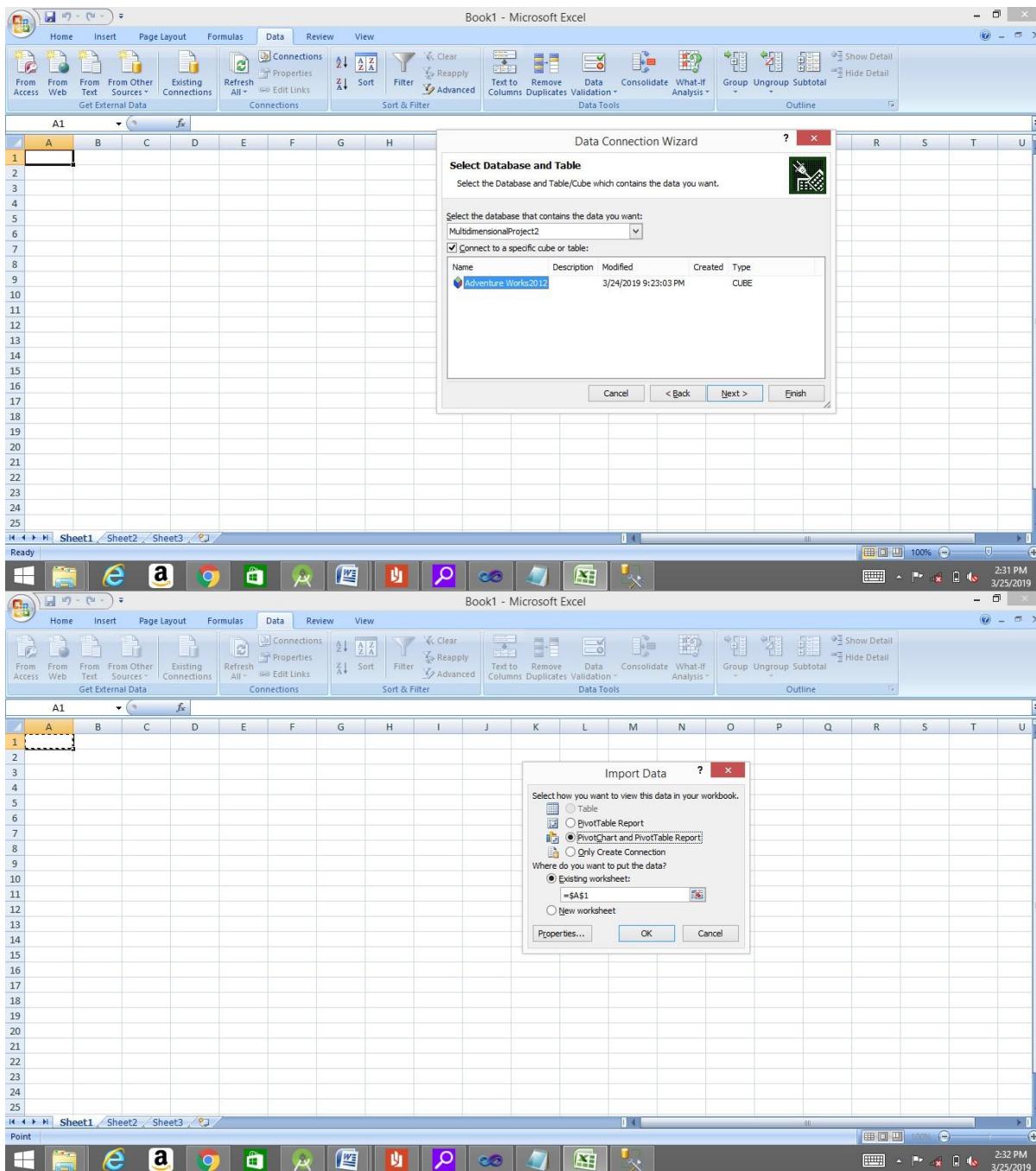
Pivot Tables

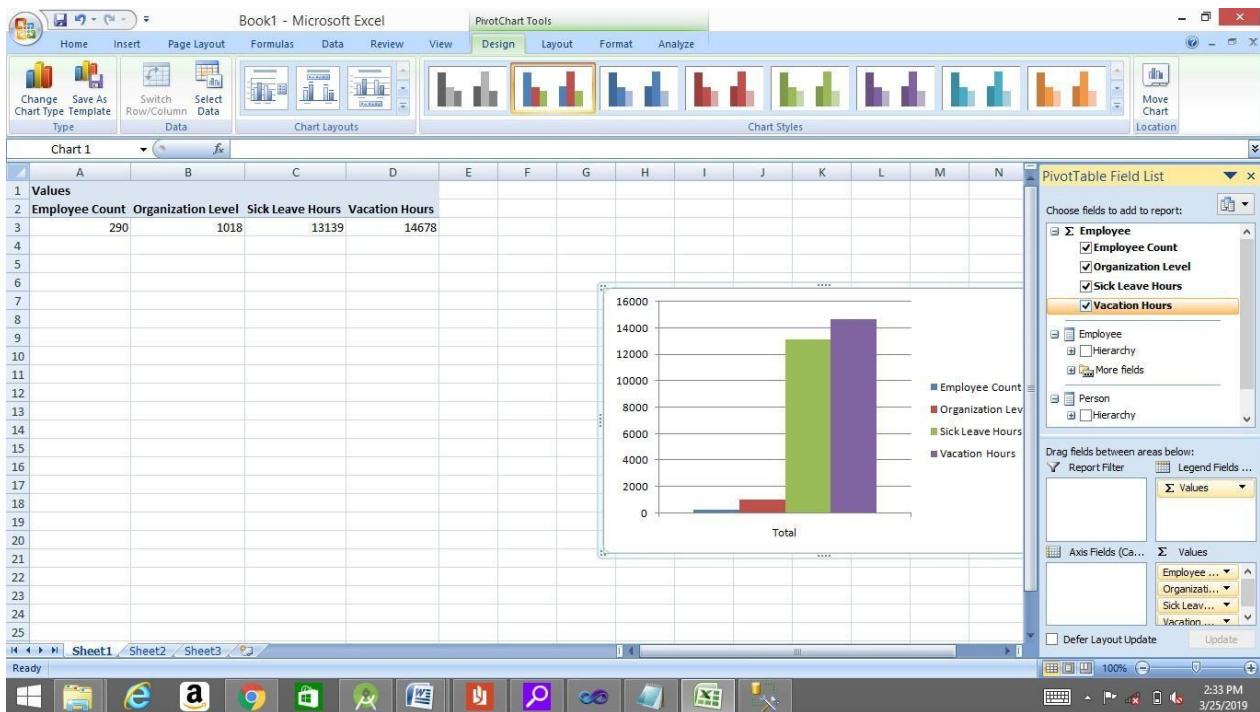
- **Pivot Tables** allow you to create a powerful view with data summarized **in a grid**, both **in** horizontal and vertical columns (also known as Matrix Views or Cross Tabs)
- A pivot chart is the visual representation of a pivot table in Excel. Pivot charts and pivot tables are connected with each other.

To create Pivot table and Pivot Chart Open Microsoft Excel









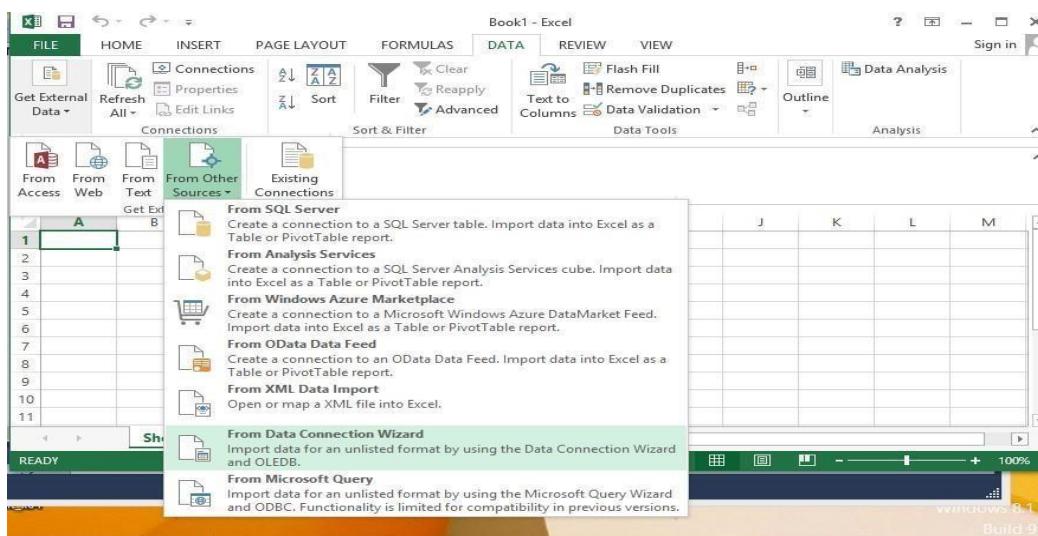
b. Import the cube in Microsoft Excel and create the Pivot table and Pivot Chart to Perform data analysis.

Import the datawarehouse data in Microsoft Excel and create the Pivot table and Pivot Chart

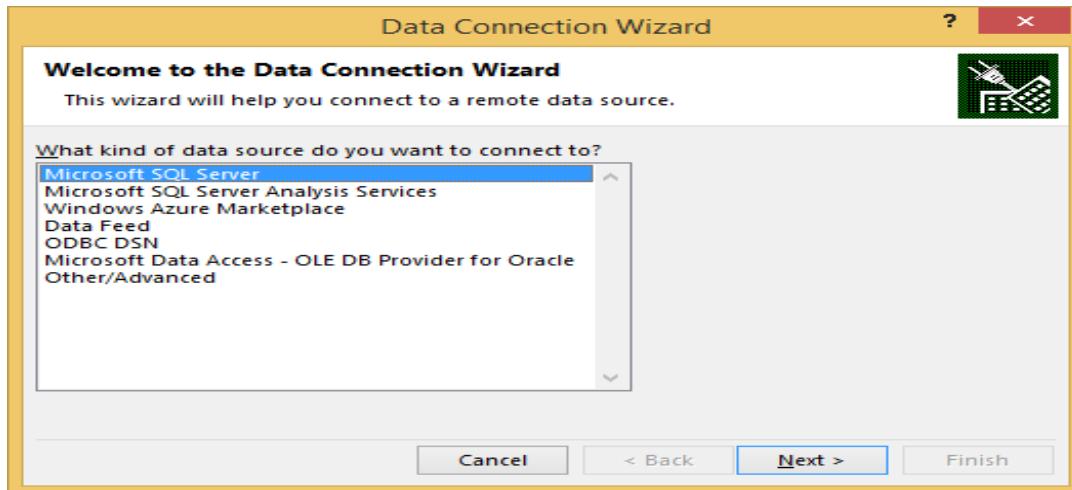
(Ms Office Professional is used to make sure Power View is enabled for visualization.)

Step 1: Open Excel 2013 (Professional)

Go to Data tab → Get External Data → From Other Sources → From Data Connection Wizard



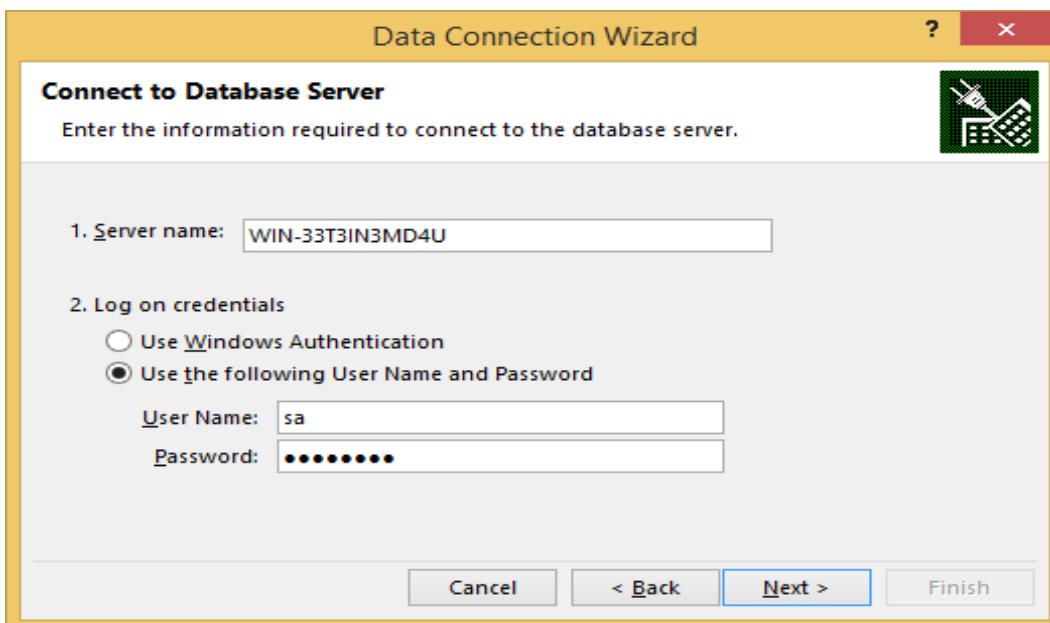
Step 2: In Data Connection Wizard → Select Microsoft SQL Server → Click on Next



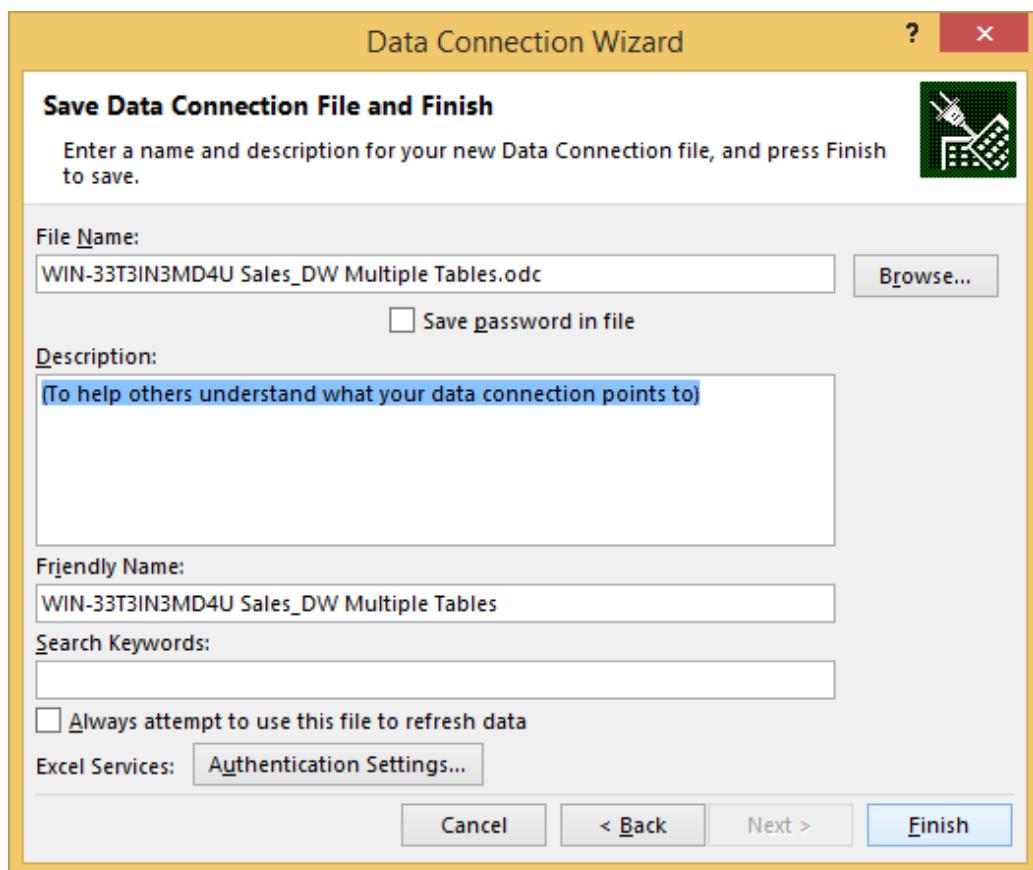
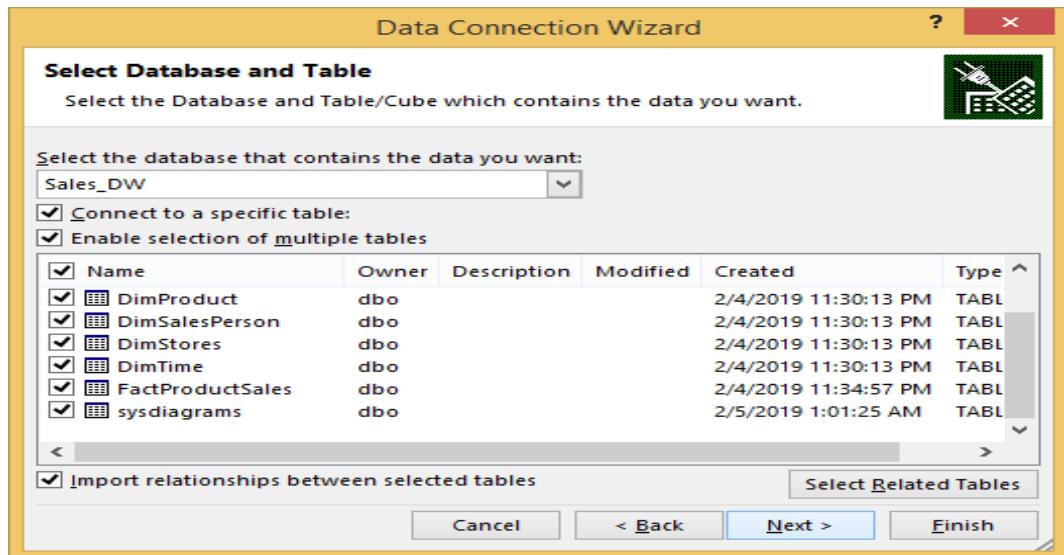
Step 3: In connect to Database Server provide Server name(Microsoft SQL Server Name)

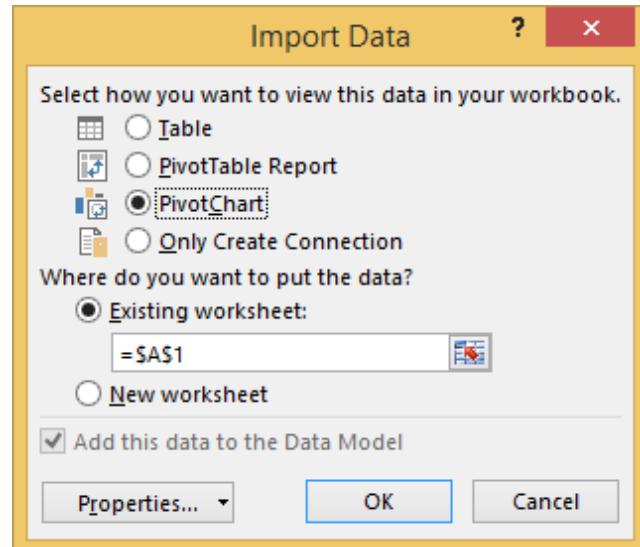
Provide password for sa account as given during installation of SQL Server 2012 full version)

Password: admin123 Click on Next

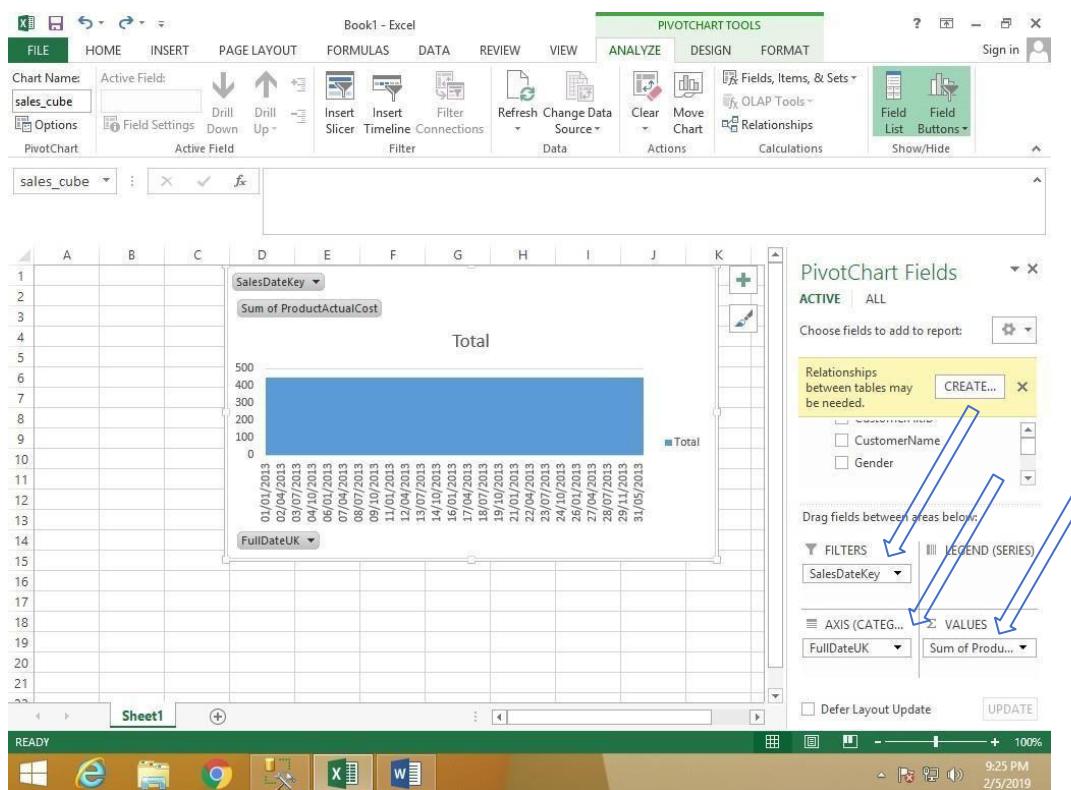


Step 4: In Select Database and Table → Select Sales_DW (already created in SQL) → check all dimensions and import relationships between selected tables

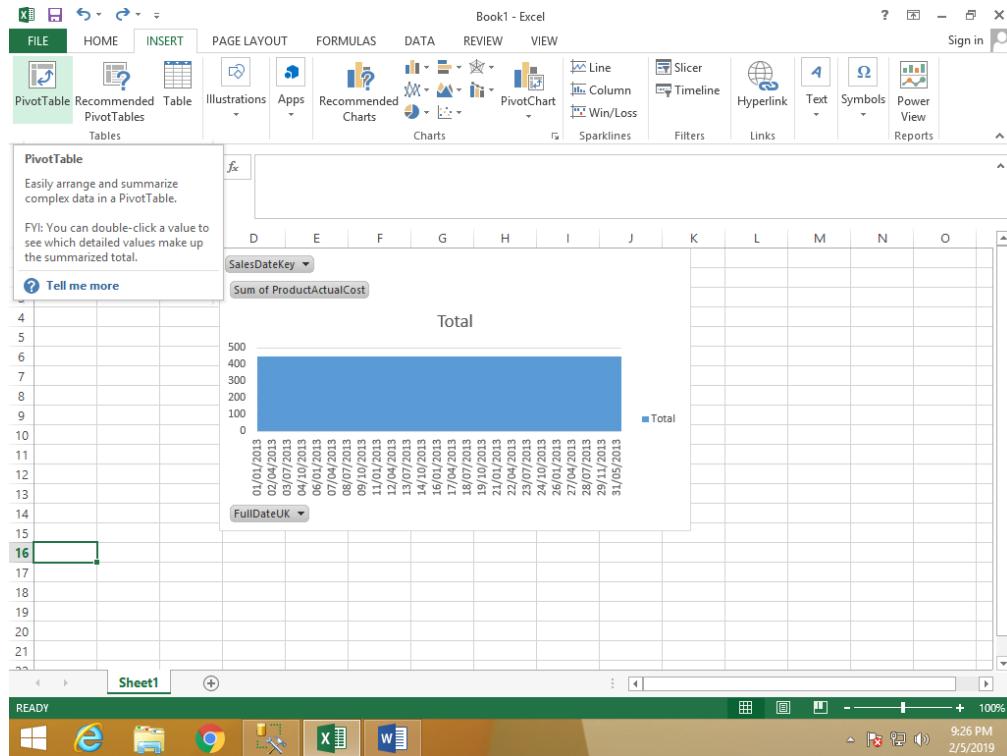




Step 7: In fields put SalesDateKey in filters, FullDateUK in axis and Sum of ProductActualCost in values



Step 8: In Insert Tab → go to Pivot Table



Practical No.2

Apply the what – if Analysis for data visualization. Design and generate necessary reports based on the data warehouse data. Use Excel.

What-If Analysis

- **What-If Analysis** is the process of changing the values in cells to see how those changes will affect the outcome of formulas on the worksheet.
- What-If Analysis in Excel allows you to try out different values (scenarios) for formulas.
- Three kinds of **What-If Analysis** tools come with **Excel**:
- Scenarios, Goal Seek, and Data Tables.
- Scenarios and Data tables take sets of input values and determine possible results.

Goal Seek

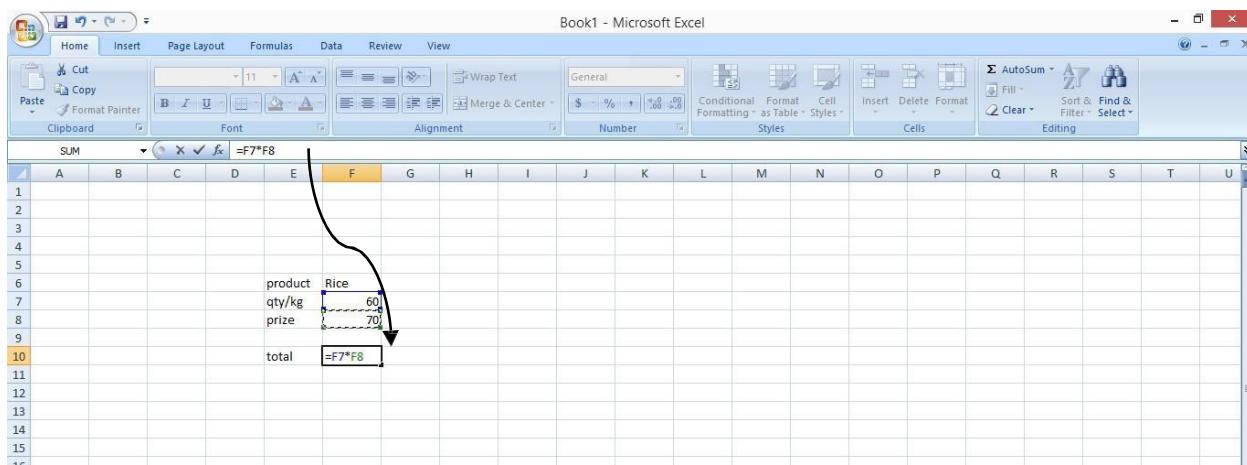
- If you know the result that you want from a formula, but you are not sure what input value the formula requires to get that result, you can use the [Goal Seek](#) feature.

Data Table

- If you have a formula that uses one or two variables, or multiple formulas that all use one common variable, you can use a [Data Table](#) to see all the outcomes in one place.

Scenario Manager

- A [Scenario Manager](#) is a set of values that Excel saves and can substitute automatically in cells on a worksheet. You can create and save different groups of values on a worksheet and then switch to any of these new scenarios to view different results.

a) Understanding Goal Seek Option

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The data is organized into columns A through M. Row 6 contains the header "product Rice". Rows 7 and 8 show data points: "qty/kg 60" and "prize 70" respectively. Row 10 is a summary row labeled "total" with the value "4200" in the adjacent cell. The formula bar at the top indicates the formula =F7*F8.

To Change qyt to achieve 6000/- . How much qty will be required.

Go to What-If analysis and select Goal Seek Option

The screenshot shows the same Microsoft Excel spreadsheet as before, but with the "Goal Seek" dialog box open. The dialog box has the following settings: "Set cell:" F10, "To value:" 6000, and "By changing cell:" \$F\$7. The background spreadsheet remains the same, with the total value still showing 4200.

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The formula bar indicates the formula $=F7*F8$ is entered into cell F10. The spreadsheet contains the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1																					
2																					
3																					
4																					
5																					
6						product	Rice														
7						qty/kg	85.71429														
8						prize	70														
9																					
10						total	6000														
11																					
12																					
13																					
14																					
15																					
16																					
17																					
18																					
19																					
20																					
21																					
22																					
23																					
24																					
25																					

The "Goal Seek Status" dialog box is open, showing the following information:

- Goal Seeking with Cell F10 found a solution.
- Target value: 6000
- Current value: 6000
- Buttons: Step, Pause, OK, Cancel

To Change prize to achieve 6000/- How much prize will be required.

Go to What-If analysis and select Goal Seek Option

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The formula bar indicates the formula $=F7*F8$ is entered into cell F8. The spreadsheet contains the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1																					
2																					
3																					
4																					
5						product	Rice														
6						qty/kg	60														
7						prize	70														
8																					
9						total	4200														
10																					
11																					
12																					
13																					
14																					
15																					
16																					
17																					
18																					
19																					
20																					
21																					
22																					
23																					
24																					
25																					

The "Goal Seek" dialog box is open, showing the following settings:

- Set cell: F10
- To value: 6000
- By changing cell: \$F\$8
- Buttons: OK, Cancel

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The formula bar displays the formula $=F7*F8$. The spreadsheet contains the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1																					
2																					
3																					
4																					
5																					
6						product	Rice														
7						qty/kg	60														
8						prize	100														
9																					
10						total		6000													
11																					
12																					
13																					
14																					
15																					
16																					
17																					
18																					
19																					
20																					
21																					
22																					
23																					
24																					
25																					

A "Goal Seek Status" dialog box is open, showing the following information:

- Goal Seeking with Cell F10 found a solution.
- Target value: 6000
- Current value: 6000
- Buttons: Step, Pause, OK, Cancel

b) Understanding Data Table Option

Generate Data table using qty.

For this, go to What-If analysis and select Data Table option.

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The formula bar displays the formula $=F7*F8$. The spreadsheet contains the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1																					
2																					
3																					
4																					
5																					
6						product	Rice														
7						qty/kg	60														
8						prize	70														
9																					
10						total	4200														
11																					
12																					
13																					
14																					
15																					
16																					
17																					
18																					
19																					
20																					
21																					
22																					
23																					
24																					
25																					

A data table has been generated in row K, showing the result of changing the quantity (qty) from 80 to 120. The table includes columns for product, qty, and total.

Book1 - Microsoft Excel

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The data is as follows:

	product	Rice	qty
7	qty/kg	60	4200
8	prize	70	
9			80
10			90
11			100
12			110
			120

A Data Table dialog box is open, showing the following input fields:

- Row input cell: \$F\$7
- Column input cell: \$F\$7

The status bar at the bottom indicates: Average: 783.3333333 Count: 6 Sum: 4700.

Book1 - Microsoft Excel

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The data is as follows:

	product	Rice	qty
7	qty/kg	60	4200
8	prize	70	
9			80
10			90
11			100
12			110
			120

A Data Table dialog box is open, showing the following input fields:

- Row input cell: \$F\$7
- Column input cell: \$F\$7

The status bar at the bottom indicates: Average: 783.3333333 Count: 6 Sum: 4700.

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The data is organized into columns A through U. Row 6 contains column headers: "product" (in A), "Rice" (in B), and "qty" (in C). Rows 7 through 11 show data points: (qty/kg, value), (prize, value), and (total, value). Row 12 is a summary row with values from rows 7 to 11. The formula bar at the top shows "L12". The status bar at the bottom indicates "Ready".

product	Rice	qty
qty/kg	60	4200
prize	70	5600
		6300
	4200	7000
		7700
		8400

Generate Data table using prize.

This screenshot shows the same Microsoft Excel spreadsheet as the first one, but it includes a new column "prize" (column H). The data now spans columns A through H. The "prize" column contains values 120, 130, 140, and 150. The formula bar at the top shows "H15". The status bar at the bottom shows "Average: 948", "Count: 5", and "Sum: 4740".

product	Rice	qty	prize	120	130	140	150
qty/kg	60	4200					
prize	70		120				
			130				
			140				
			150				
			4200				

Book1 - Microsoft Excel

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The Data tab is active, and the Data Tools group is open, displaying various data analysis tools like Connections, Sort & Filter, and Data Tools. The "Data Table..." option is highlighted. The spreadsheet contains a data table from row 6 to 16, spanning columns F to L. The table includes headers "product", "Rice", and "qty", and data rows for "qty/kg" (60, 80, 90, 100, 110, 120), "prize" (70, 5600, 6300, 7000, 7700, 8400), and a total row (4200). The "prize" column is selected. The status bar at the bottom shows "Average: 948", "Count: 5", "Sum: 4740", and "100%". The taskbar at the bottom shows various application icons.

Book1 - Microsoft Excel

This screenshot shows the same Microsoft Excel spreadsheet as the previous one, but the "Data Table" dialog box is now open. The dialog box has two input fields: "Row input cell" containing "\$F\$8" and "Column input cell" containing "\$F\$7". There are "OK" and "Cancel" buttons at the bottom. The rest of the spreadsheet and interface are identical to the first screenshot.

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The data is organized into several sections:

- Section 1 (Rows 6-12):** A header section with "product" and "Rice" in columns E and F respectively. Below it, "qty/kg" is in column E, and "prize" is in column F.
- Section 2 (Rows 7-11):** Data rows showing quantities and prices. Row 7: qty/kg 60, prize 4200. Row 8: prize 70, qty/kg 80, 5600. Row 9: prize 90, qty/kg 6300. Row 10: total 4200, prize 100, 7000. Row 11: total 4200, prize 110, 7700. Row 12: total 4200, prize 120, 8400.
- Section 3 (Row 15):** A summary row labeled "prize" with values 120, 130, 140, 150, and 4200.
- Section 4 (Row 16):** A summary row labeled "prize" with values 7200, 7800, 8400, 9000, and 4200.

The status bar at the bottom indicates: Average: 4126.666667 Count: 9 Sum: 37140 100% 2:51 PM 3/25/2019

Generate Data table using qty and prize both.

This screenshot shows the same data as the first one, but the structure is slightly different. The "prize" column is moved to the left of the "qty/kg" column in the header section. The "total" rows are also placed before the "prize" summary row.

The status bar at the bottom indicates: Average: 568.8888889 Count: 9 Sum: 5120 100% 2:52 PM 3/25/2019

Book1 - Microsoft Excel

Average: 568.8888889 Count: 9 Sum: 5120

2:53 PM 3/25/2019

	product	Rice	qty	
7	qty/kg	60		4200
8	prize	70		80 5600
9				90 6300
10		total	4200	100 7000
11				110 7700
12				120 8400
14			qty	
15		prize	4200 120 130 140 150	
16			80	
17			90	
18			100	
19			110	

Book1 - Microsoft Excel

Average: 568.8888889 Count: 9 Sum: 5120

2:53 PM 3/25/2019

	product	Rice	qty	
7	qty/kg	60		4200
8	prize	70		80 5600
9				90 6300
10		total	4200	100 7000
11				110 7700
12				120 8400
14			qty	
15		prize	4200 120 130 140 150	
16			80	
17			90	
18			100	
19			110	

The screenshot shows a Microsoft Excel spreadsheet titled "Book1 - Microsoft Excel". The data is organized into two main sections:

- Data Section:** Rows 6 to 12 show a table with columns for "product" (Rice), "qty/kg" (60, 70), and "total" (4200). Row 13 is blank.
- Pivot Table Section:** Rows 14 to 19 show a pivot table with a single row for "prize" and five columns for "qty" (4200, 120, 130, 140, 150). The data is summarized by "prize" levels: 80, 90, 100, 110, and 120, with corresponding quantities: 9600, 10800, 12000, 13200, and 14300 respectively.

The Excel ribbon is visible at the top, and the status bar at the bottom shows "Average: 8412.8 Count: 25 Sum: 210320". The system tray at the bottom right indicates the date and time as 3/25/2019 2:53 PM.

c) Understanding Scenario Manager Option

This screenshot shows the same Microsoft Excel spreadsheet as above, but with a focus on the "Data" tab's ribbon group. The "Scenario Manager..." button is highlighted with a yellow box. The data in the table remains the same, with the value in cell F7 changed to 60.

The status bar at the bottom shows "Average: 65 Count: 2 Sum: 130". The system tray at the bottom right shows the date and time as 3/25/2019 2:55 PM.

product Rice

qty/kg 60

prize 70

total 4200

product Rice

qty/kg 60

prize 70

total 4200

Book1 - Microsoft Excel

Scenario Manager

Scenarios:

- monday

Changing cells: \$F\$7:\$F\$8

Comment: Created by Windows User on 3/25/2019

	product	Rice
7	qty/kg	60
8	prize	70
10	total	4200

Book1 - Microsoft Excel

	product	Rice
7	qty/kg	4
8	prize	5
10	total	280

Book1 - Microsoft Excel

Home Insert Page Layout Formulas Data Review View

Connections Refresh Properties All Sort Filter Advanced

Text to Columns Remove Duplicates Validation Data Tools

Group Ungroup Subtotal Outline

F7 f_x 4

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1																					
2																					
3																					
4																					
5																					
6																					
7																					
8																					
9																					
10																					
11																					
12																					
13																					
14																					
15																					
16																					
17																					
18																					
19																					
20																					
21																					
22																					
23																					
24																					
25																					

Scenario Manager

Scenarios: monday

Add... Delete Edit... Merge... Summary...

Changing cells: \$F\$7:\$F\$8

Comment: Created by Windows User on 3/25/2019

Show Close

Sheet1 Sheet2 Sheet3 Sheet4

Average: 4.5 Count: 2 Sum: 9 100% 3:02 PM 3/25/2019

Book1 - Microsoft Excel

Home Insert Page Layout Formulas Data Review View

Connections Refresh Properties All Sort Filter Advanced

Text to Columns Remove Duplicates Validation Data Tools

Group Ungroup Subtotal Outline

F7 f_x 60

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1																					
2																					
3																					
4																					
5																					
6																					
7																					
8																					
9																					
10																					
11																					
12																					
13																					
14																					
15																					
16																					
17																					
18																					
19																					
20																					
21																					
22																					
23																					
24																					
25																					

Sheet1 Sheet2 Sheet3 Sheet4

Average: 65 Count: 2 Sum: 130 100% 3:02 PM 3/25/2019

Practical No.3

Perform the data classification using classification algorithm using R/Python.

Classification

- **Classification** is a **data mining** function that assigns items in a collection to target categories or classes.
- The goal of **classification** is to accurately predict the target class for each case in the **data**.
- For example, a **classification** model could be used to identify loan applicants as risky and safe.
- Classifier
- Prediction

Consider the annual rainfall details at a place starting from January 2012. We create an R time series object for a period of 12 months and plot it.

```
import pandas as pd
```

```
df= pd.read_csv("/content/shows.csv")
```

```
print(df)
```

	Age	Experience	Rank	Nationality	Go
0	36	10	9	UK	NO
1	42	12	4	USA	NO
2	23	4	6	N	NO
3	52	4	4	USA	NO
4	43	21	8	USA	YES
5	44	14	5	UK	NO
6	66	3	7	N	YES
7	35	14	9	UK	YES
8	52	13	7	N	YES
9	35	5	9	N	YES
10	24	3	5	USA	NO
11	18	3	7	UK	YES
12	45	9	9	UK	YES

```
from sklearn import tree
```

```
import pydotplus
```

```
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
import matplotlib.image as pltimage
```

```
#change string values into numerical values:
d= {'UK':0, 'USA':1, 'N':2}
```

```
print (d)
```

```
→ { 'UK': 0, 'USA': 1, 'N': 2}
```

```
df['Nationality']= df['Nationality'].map(d)
```

```
df
```

	Age	Experience	Rank	Nationality	Go	
0	36	10	9	0	NO	
1	42	12	4	1	NO	
2	23	4	6	2	NO	
3	52	4	4	1	NO	
4	43	21	8	1	YES	
5	44	14	5	0	NO	
6	66	3	7	2	YES	
7	35	14	9	0	YES	
8	52	13	7	2	YES	
9	35	5	9	2	YES	
10	24	3	5	1	NO	
11	18	3	7	0	YES	
12	45	9	9	0	YES	

```
d={'YES':1, 'NO': 0}
```

```
df['Go']=df['Go'].map(d)
```

```
print(df)
```

	Age	Experience	Rank	Nationality	Go	
0	36	10	9	0	0	
1	42	12	4	1	0	
2	23	4	6	2	0	
3	52	4	4	1	0	
4	43	21	8	1	1	
5	44	14	5	0	0	
6	66	3	7	2	1	
7	35	14	9	0	1	
8	52	13	7	2	1	
9	35	5	9	2	1	
10	24	3	5	1	0	
11	18	3	7	0	1	
12	45	9	9	0	1	

```
features=['Age','Experience','Rank','Nationality']
```

```
x=df[features]
```

```
y=df['Go']
```

```
print(x)
```

	Age	Experience	Rank	Nationality
0	36	10	9	0
1	42	12	4	1
2	23	4	6	2
3	52	4	4	1
4	43	21	8	1
5	44	14	5	0
6	66	3	7	2
7	35	14	9	0
8	52	13	7	2
9	35	5	9	2
10	24	3	5	1
11	18	3	7	0
12	45	9	9	0

```
print(y)
```

0	0
1	0
2	0
3	0
4	1
5	0
6	1
7	1
8	1
9	1
10	0
11	1
12	1

Name: Go, dtype: int64

```
df.shape
```

```
(13, 5)
```

```
dtree=DecisionTreeClassifier()
```

```
dtree=dtree.fit(x,y)
```

Decision Tree

```

import pandas as pd
-----
from sklearn import tree
import pydotplus
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
import matplotlib.image as pltimage
-----
df= pd.read_csv("/content/shows.csv")
print(df)

```

	Age	Experience	Rank	Nationality	Go
0	36	10	9	UK	NO
1	42	12	4	USA	NO
2	23	4	6	N	NO
3	52	4	4	USA	NO
4	43	21	8	USA	YES
5	44	14	5	UK	NO
6	66	3	7	N	YES
7	35	14	9	UK	YES
8	52	13	7	N	YES
9	35	5	9	N	YES
10	24	3	5	USA	NO
11	18	3	7	UK	YES
12	45	9	9	UK	YES

```

# Change string values into numerical values:
d = {'UK': 0, 'USA': 1, 'N': 2}
df['Nationality'] = df['Nationality'].map(d)
d = {'YES': 1, 'NO': 0}
df['Go'] = df['Go'].map(d)
-----
print(df)

```

	Age	Experience	Rank	Nationality	Go
0	36	10	9	0	0
1	42	12	4	1	0
2	23	4	6	2	0
3	52	4	4	1	0
4	43	21	8	1	1
5	44	14	5	0	0
6	66	3	7	2	1
7	35	14	9	0	1
8	52	13	7	2	1
9	35	5	9	2	1
10	24	3	5	1	0
11	18	3	7	0	1
12	45	9	9	0	1

```
print(x)
```

```
print(y)
```

```
df.shape
```

```
   ➔    Age  Experience  Rank  Nationality
  0     36          10     9           0
  1     42          12     4           1
  2     23           4     6           2
  3     52           4     4           1
  4     43          21     8           1
  5     44          14     5           0
  6     66           3     7           2
  7     35          14     9           0
  8     52          13     7           2
  9     35           5     9           2
 10    24           3     5           1
 11    18           3     7           0
 12    45           9     9           0
 0     0
 1     0
 2     0
 3     0
 4     1
 5     0
 6     1
 7     1
 8     1
 9     1
10    0
11    1
12    1
Name: Go, dtype: int64
(13, 5)
```

```
features = ['Age','Experience','Rank','Nationality']
```

```
x=df[features]
```

```
y=df['Go']
```

```
print(x)
```

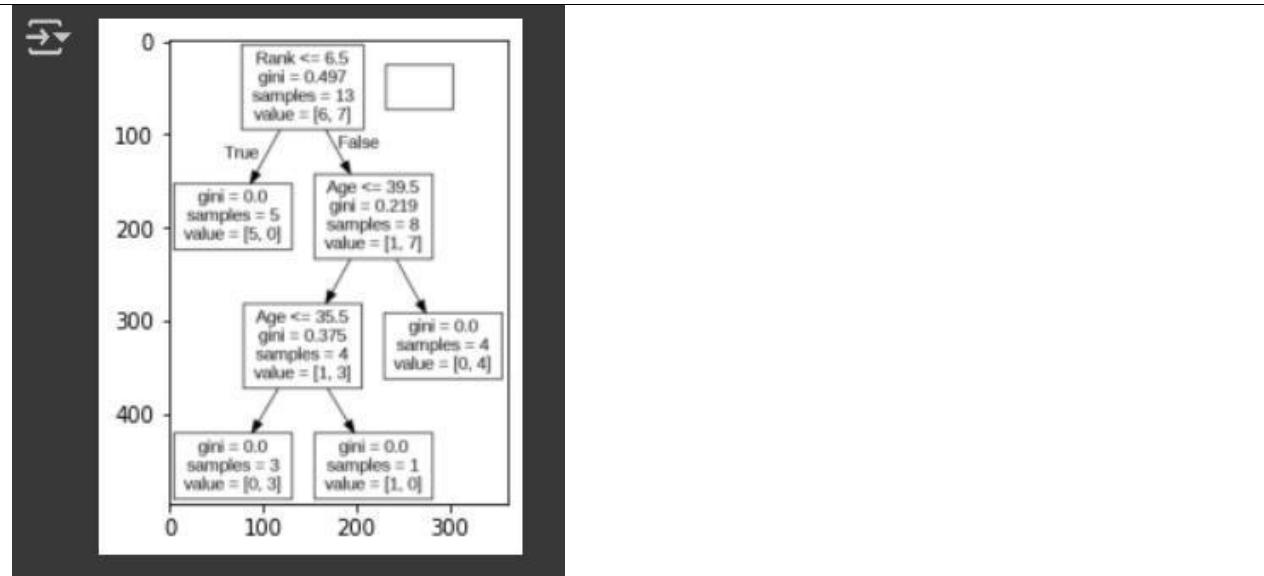
```
print(y)
```

```
→ Age  Experience  Rank  Nationality
  0   36           10    9      0
  1   42           12    4      1
  2   23           4     6      2
  3   52           4     4      1
  4   43           21    8      1
  5   44           14    5      0
  6   66           3     7      2
  7   35           14    9      0
  8   52           13    7      2
  9   35           5     9      2
  10  24           3     5      1
  11  18           3     7      0
  12  45           9     9      0
  0   0
  1   0
  2   0
  3   0
  4   1
  5   0
  6   1
  7   1
  8   1
  9   1
  10  0
  11  1
  12  1
Name: Go, dtype: int64
```

```
df.shape
```

```
→ (13, 5)
```

```
dtree=DecisionTreeClassifier()
dtree=dtree.fit(x.values,y.values)
data=tree.export_graphviz(dtree,out_file=None,feature_names=features)
graph=pydotplus.graph_from_dot_data(data)
graph.write_png('mydecisiontree.png')
img=pltimage.imread('mydecisiontree.png')
imgplot=plt.imshow(img)
plt.show()
```



```
print(dtree.predict([[40,10,7,1]]))
```

→ [1]

```
print(dtree.predict([[40,10,6,1]]))
```

→ [0]

```
features = ['Age','Experience','Rank','Nationality']
```

```
x=df[features]
```

```
print(x)
```

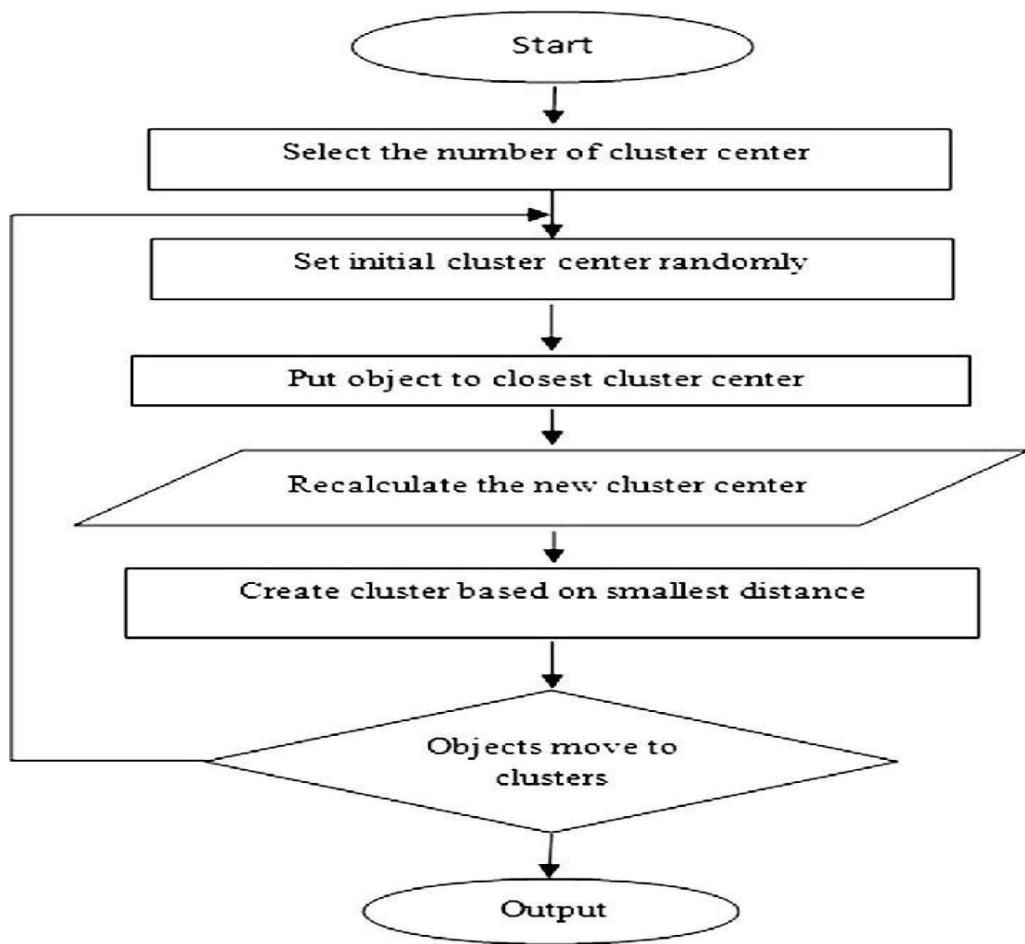
	Age	Experience	Rank	Nationality
0	36	10	9	0
1	42	12	4	1
2	23	4	6	2
3	52	4	4	1
4	43	21	8	1
5	44	14	5	0
6	66	3	7	2
7	35	14	9	0
8	52	13	7	2
9	35	5	9	2
10	24	3	5	1
11	18	3	7	0
12	45	9	9	0

Practical No.4

Perform the data clustering using clustering algorithm using R/Python.

Clustering

- Clustering is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group than those in other groups.
- Let's understand this with an example.
- Suppose, you are the head of a rental store and wish to understand preferences of your Customers to scale up your business.
- Is it possible for you to look at details of each costumer and devise a unique business strategy for each one of them?
- Definitely not. But, what you can do is to cluster all of your Customers into say 10 groups based on their purchasing habits and use a separate strategy for costumers in each of these 10 groups. And this is what we call clustering.

K-Means Clustering

```
import pandas as pd
import numpy as np
import random as rd
import matplotlib.pyplot as plt

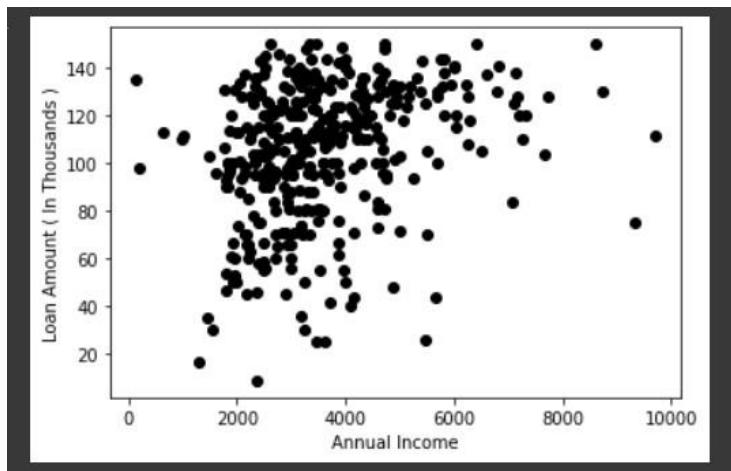
data = pd.read_csv("/content/clustering.csv")
```

data

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CosapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status
0	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0	1.0	Rural	N
1	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	360.0	1.0	Urban	Y
2	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	360.0	1.0	Urban	Y
3	LP001008	Male	No	0	Graduate	No	6000	0.0	141.0	360.0	1.0	Urban	Y
4	LP001013	Male	Yes	0	Not Graduate	No	2333	1516.0	95.0	360.0	1.0	Urban	Y
...
376	LP002953	Male	Yes	3+	Graduate	No	5703	0.0	128.0	360.0	1.0	Urban	Y
377	LP002974	Male	Yes	0	Graduate	No	3232	1950.0	108.0	360.0	1.0	Rural	Y
378	LP002978	Female	No	0	Graduate	No	2900	0.0	71.0	360.0	1.0	Rural	Y
379	LP002979	Male	Yes	3+	Graduate	No	4106	0.0	40.0	180.0	1.0	Rural	Y
380	LP002990	Female	No	0	Graduate	Yes	4583	0.0	133.0	360.0	0.0	Semirurban	N

381 rows × 13 columns

```
X=data[["LoanAmount","ApplicantIncome"]]
#Visualize the data points (ap,la)
plt.scatter(X["ApplicantIncome"],X["LoanAmount"],c="black")
plt.xlabel("Annual Income")
plt.ylabel("Loan Amount ( In Thousands )")
plt.show()
```

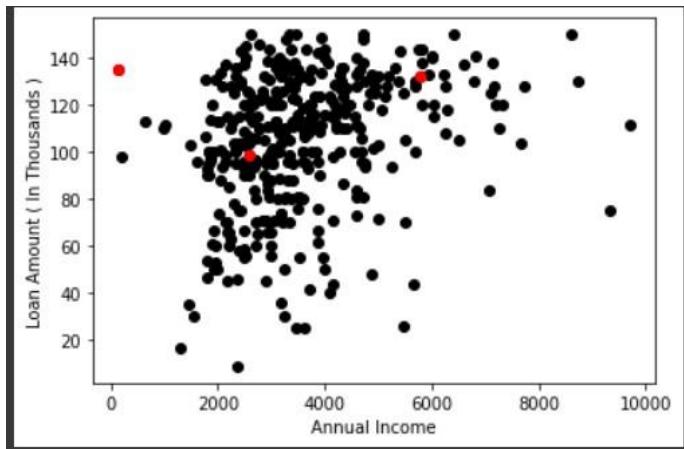


```
# Step 1 and Step 2 Choose the number of clusters(K) and select random
centriod for each cluster
# number of clusters
K=3
# Select random observations as centriods
Centroids = (X.sample(n=K))
```

```

plt.scatter(X["ApplicantIncome"],X["LoanAmount"],c="black")
plt.scatter(Centroids["ApplicantIncome"],Centroids["LoanAmount"],c="red")
plt.xlabel("Annual Income")
plt.ylabel("Loan Amount ( In Thousands )")
plt.show()

```



```

diff = 1
j=0

while(diff!=0):
    XD=X
    i=1
    for index1,row_c in Centroids.iterrows():
        ED=[]
        for index2,row_d in XD.iterrows():
            d1=(row_c["ApplicantIncome"]-row_d["ApplicantIncome"])**2
            d2=(row_c["LoanAmount"]-row_d["LoanAmount"])**2
            d=np.sqrt(d1+d2)
            ED.append(d)
        X[i]=ED
        i=i+1

    C=[]
    for index,row in X.iterrows():
        min_dist=row[1]
        pos=1
        for i in range(K):
            if row[i+1] < min_dist:
                min_dist = row[i+1]
                pos=i+1
        C.append(pos)
    X["Cluster"]=C

```

```

Centroids_new = X.groupby(["Cluster"]).mean()[["LoanAmount", "ApplicantIncome"]]
if j == 0:
    diff=1
    j=j+1
else:
    diff = (Centroids_new['LoanAmount'] -
Centroids['LoanAmount']).sum() + (Centroids_new['ApplicantIncome'] -
Centroids['ApplicantIncome']).sum()
    print(diff.sum())
Centroids =
X.groupby(["Cluster"]).mean()[["LoanAmount", "ApplicantIncome"]]

```

```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:14: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:26: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

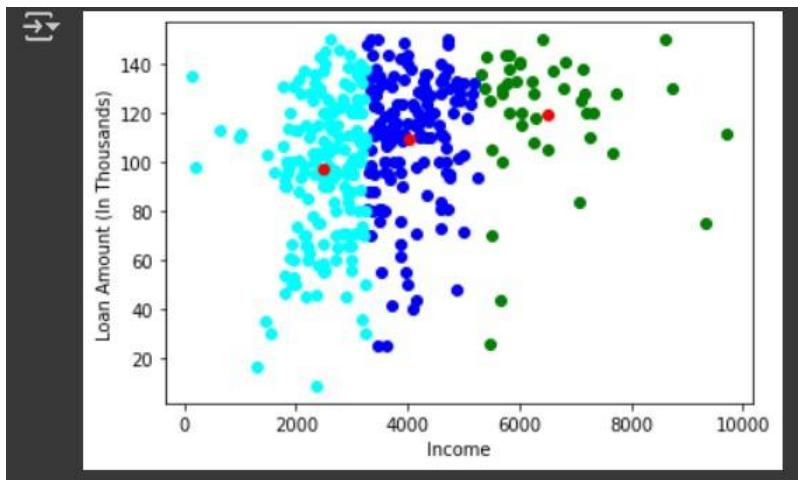
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
760.9691465468836
570.3181122036852
400.02659764803184
495.26981705062997
186.27599218629547
227.7568592993333
255.37125148586168
244.66095351174067
229.06905235705375
218.24897861156342
107.07928213052429
52.84741626127729
98.54724443834282
90.64953219227577
18.274686272279013
9.21023994083339
18.345487493007468
46.27013250786139
0.0

```

```

color=['blue','green','cyan']
for k in range(K):
    data=X[X["Cluster"]==k+1]
    plt.scatter(data["ApplicantIncome"],data["LoanAmount"],c=color[k])
plt.scatter(Centroids["ApplicantIncome"],Centroids["LoanAmount"],c='red')
plt.xlabel('Income')
plt.ylabel('Loan Amount (In Thousands)')
plt.show()

```



Practical No.5

Perform the Linear regression on the given data warehouse data using R/Python.

Regression

- In statistical modeling, regression analysis is a set of statistical processes for estimating the relationships between a dependent variable and one or more independent variables

Linear Regression

- In Linear Regression these two variables are related through an equation, where exponent (power) of both these variables is 1.
- $y = ax + b$ is an equation for linear regression.
- Where, y is the response variable, x is the predictor variable and a and b are constants which are called the coefficients.

lm() Function

- In R, the lm(), or “linear model,” function can be used to create a simple regression model. The lm() function accepts a number of arguments (“Fitting Linear Models,” n.d.).

```

x<- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)
.
y<- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)

# Apply the lm() function.
relation <- lm(y~x)

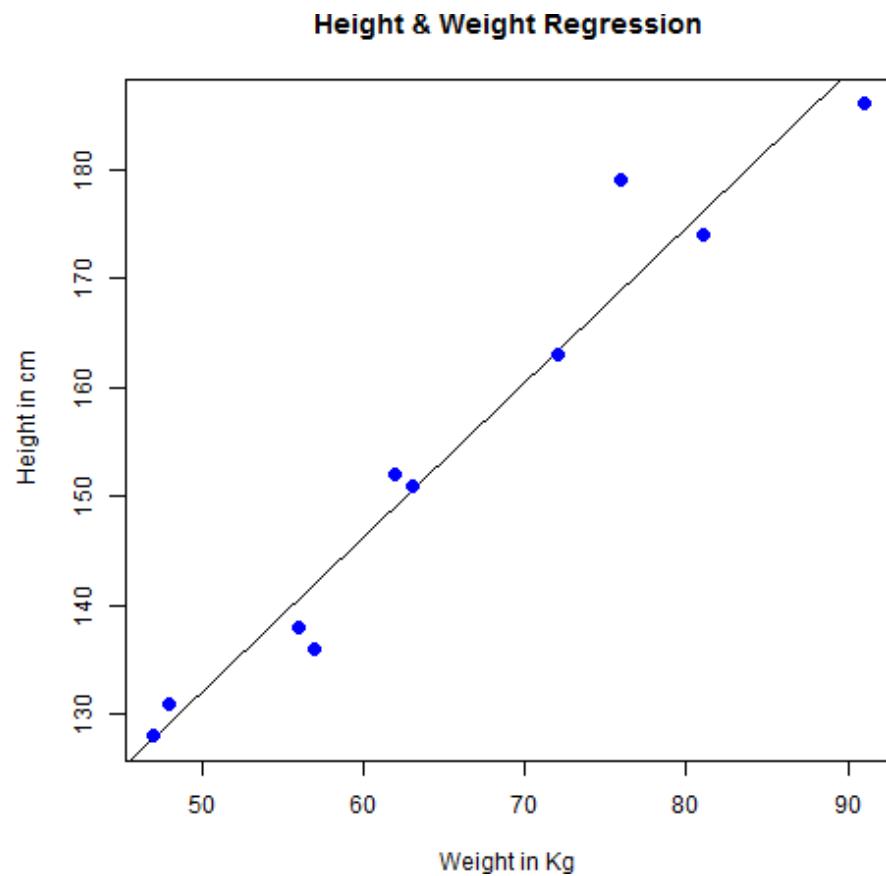
# Find weight of a person with height 170.
a <- data.frame(x = 170)
result <- predict(relation,a)
print(result)

# Give the chart file a name.
png(file = "linearregression.png")

# Plot the chart.
plot(y,x,col = "blue",main = "Height & Weight Regression",
abline(lm(x~y)),cex = 1.3,pch = 16,xlab = "Weight in Kg",ylab = "Height in cm")

# Save the file.
dev.off()

```



Practical No.6

Perform the logistic regression on the given data warehouse data using R/Python.

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

```
x = dataset.iloc[:, [2, 3]].values  
y = dataset.iloc[:, 4].values
```

X

```
array([[ 19, 190000],  
       [ 35, 200000],  
       [ 26, 430000],  
       [ 27, 570000],  
       [ 19, 760000],  
       [ 27, 580000],  
       [ 27, 840000],  
       [ 32, 1500000],  
       [ 25, 330000],  
       [ 35, 650000],  
       [ 26, 800000],  
       [ 26, 520000],  
       [ 20, 860000],  
       [ 32, 180000],  
       [ 18, 820000],  
       [ 29, 800000],  
       [ 47, 250000],  
       [ 45, 260000],  
       [ 46, 280000],  
       [ 48, 290000],  
       [ 45, 220000],  
       [ 47, 490000],  
       [ 48, 410000],
```

y

```

from sklearn.model_selection import train_test_split
xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size = 0.25, random_state = 0)

-----
from sklearn.preprocessing import StandardScaler
sc_x = StandardScaler()
xtrain = sc_x.fit_transform(xtrain)
xtest = sc_x.transform(xtest)
print (xtrain)
[[ 0.58164944 -0.88670699]
 [-0.60673761  1.46173768]
 [-0.01254409 -0.5677824 ]
 [-0.60673761  1.89663484]
 [ 1.37390747 -1.40858358]
 [ 1.47293972  0.99784738]
 [ 0.08648817 -0.79972756]
 [-0.01254409 -0.24885782]
 [-0.21060859 -0.5677824 ]
 [-0.21060859 -0.19087153]
 [-0.30964085 -1.29261101]
 [-0.30964085 -0.5677824 ]
 [ 0.38358493  0.09905991]
 [ 0.8787462 -0.59677555]
 [ 2.06713324 -1.17663843]
 [ 1.07681071 -0.13288524]
 [ 0.68068169  1.78066227]
 [-0.70576986  0.56295021]
 [ 0.77971394  0.35999821]
 [ 0.8787462 -0.53878926]
 [-1.20093113 -1.58254245]

```

```

from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(xtrain, ytrain)

```

```

→ LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                      intercept_scaling=1, l1_ratio=None, max_iter=100,
                      multi_class='auto', n_jobs=None, penalty='l2',
                      random_state=0, solver='lbfgs', tol=0.0001, verbose=0,
                      warm_start=False)

```

```
y_pred = classifier.predict(xtest)
```

```

from sklearn.metrics import confusion_matrix
cm = confusion_matrix(ytest, y_pred)
print ("Confusion Matrix : \n", cm)

```

```

→ Confusion Matrix :
 [[65  3]
 [ 8 24]]

```

```
from sklearn.metrics import accuracy_score
print ("Accuracy : ", accuracy_score(ytest, y_pred))
```

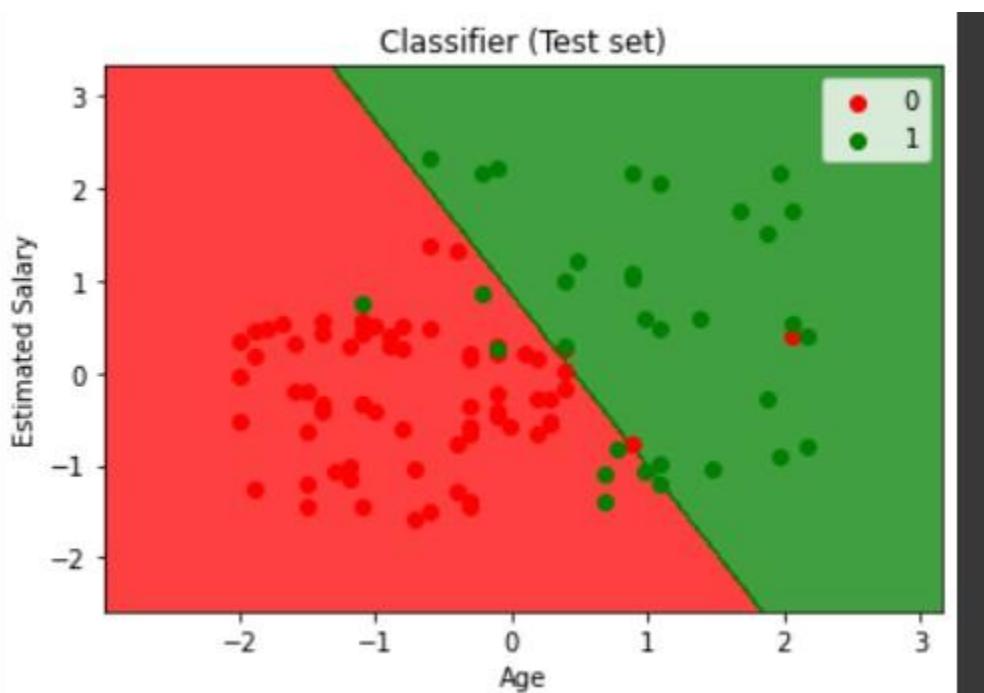
 Accuracy : 0.89

```
from matplotlib.colors import ListedColormap
X_set, y_set = xtest, ytest
X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1,
                                 stop = X_set[:, 0].max() + 1, step = 0.01),
                     np.arange(start = X_set[:, 1].min() - 1,
                               stop = X_set[:, 1].max() + 1, step = 0.01))

plt.contourf(X1, X2, classifier.predict(
    np.array([X1.ravel(), X2.ravel()]).T).reshape(
        X1.shape), alpha = 0.75, cmap = ListedColormap(('red', 'green')))

plt.xlim(X1.min(), X1.max())
plt.ylim(X2.min(), X2.max())

for i, j in enumerate(np.unique(y_set)):
    plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
                c = ListedColormap(('red', 'green'))(i), label = j)
plt.title('Classifier (Test set)')
plt.xlabel('Age')
plt.ylabel('Estimated Salary')
plt.legend()
plt.show()
```



Practical No.7

Write a Python program to read data from a CSV file, perform simple data analysis, and generate basic insights. (Use Pandas is a Python library).

```
import pandas as pd

# Step 1: Read the CSV file into a DataFrame
df = pd.read_csv("/content/data.csv")

# Show the first 5 rows of the DataFrame
print("First 5 rows of the data:")
print(df.head())

# Step 3: Data summary and descriptive statistics
print("\nData Summary (info):")
print(df.info())

print("\nDescriptive Statistics:")
print(df.describe())

# Step 4: Count of unique values for categorical columns (Gender)
print("\nCount of unique values in 'Gender':")
print(df['Gender'].value_counts())

# Step 5: Average Salary by Gender
print("\nAverage Salary by Gender:")
print(df.groupby('Gender')['Salary'].mean())

# Step 6: Check for missing values
print("\nMissing values in each column:")
print(df.isnull().sum())

# Step 7: Find the highest and lowest salary
print("\nHighest Salary:")
print(df[df['Salary'] == df['Salary'].max()])

print("\nLowest Salary:")
print(df[df['Salary'] == df['Salary'].min()])

# Step 9: Filter data (e.g., Employees with Salary > 50000)
print("\nEmployees with Salary greater than 50000:")
print(df[df['Salary'] > 50000])
```

```

First 5 rows of the data:
   Name  Age  Gender  Salary
0  John   28    Male  50000
1  Jane   34  Female  60000
2  Doe    23    Male  45000
3  Anna   41  Female  70000
4  Tom    29    Male  52000

Data Summary (info):
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5 entries, 0 to 4
Data columns (total 4 columns):
 #   Column  Non-Null Count  Dtype  
--- 
 0   Name      5 non-null     object  
 1   Age       5 non-null     int64   
 2   Gender    5 non-null     object  
 3   Salary    5 non-null     int64   
dtypes: int64(2), object(2)
memory usage: 292.0+ bytes
None

Descriptive Statistics:
              Age        Salary
count  5.000000  5.000000
mean   31.000000 55400.000000
std    6.819091  9787.747443
min    23.000000 45000.000000
25%   28.000000 50000.000000
50%   29.000000 52000.000000
75%   34.000000 60000.000000
max   41.000000 70000.000000

Count of unique values in 'Gender':
Gender
Male      3
Female    2
Name: count, dtype: int64

```

```

Name: count, dtype: int64

Average Salary by Gender:
Gender
Female  65000.0
Male    49000.0
Name: Salary, dtype: float64

Missing values in each column:
Name      0
Age       0
Gender    0
Salary    0
dtype: int64

Highest Salary:
   Name  Age  Gender  Salary
3  Anna   41  Female  70000

Lowest Salary:
   Name  Age  Gender  Salary
2  Doe   23    Male  45000

Employees with Salary greater than 50000:
   Name  Age  Gender  Salary
1  Jane   34  Female  60000
3  Anna   41  Female  70000
4  Tom    29    Male  52000

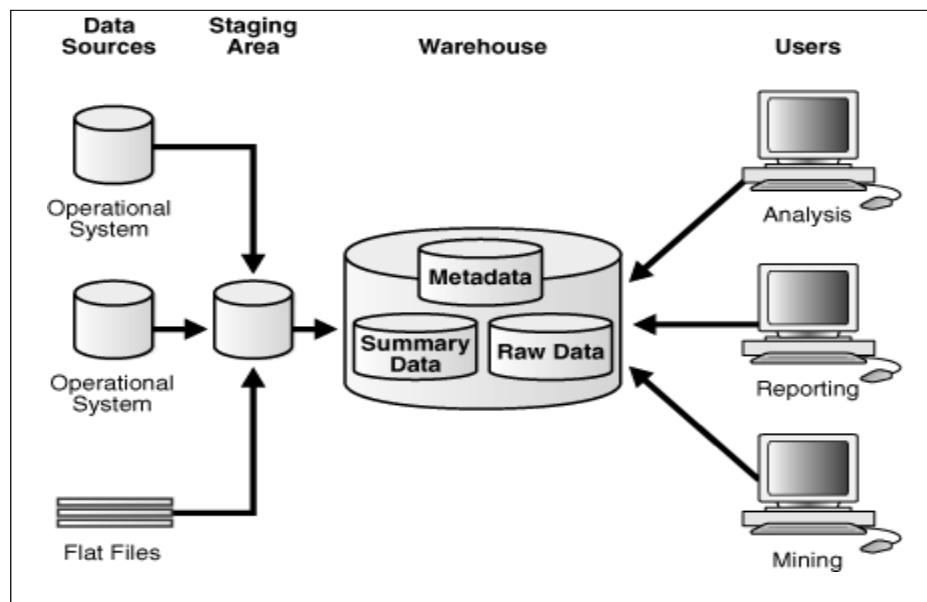
```

Practical No.9

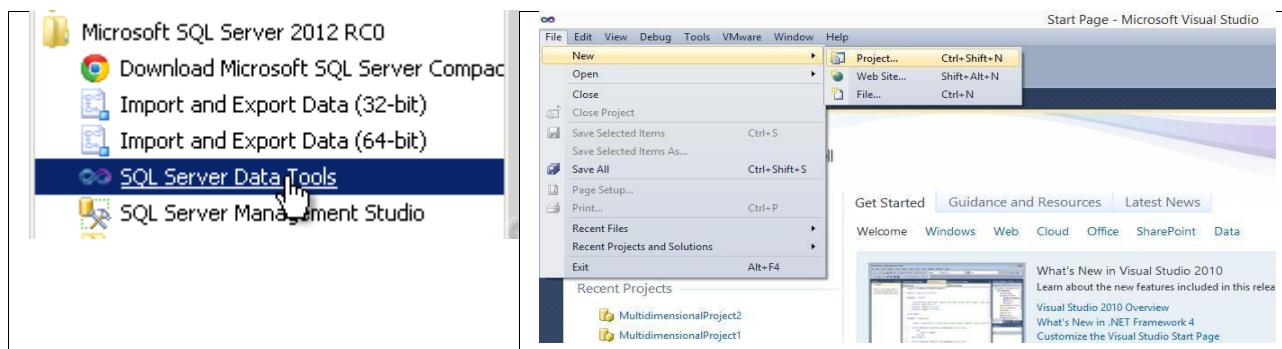
- Create the Data staging area for the selected database.
- Create the cube with suitable dimension and fact tables based on ROLAP, MOLAP and HOLAP model.

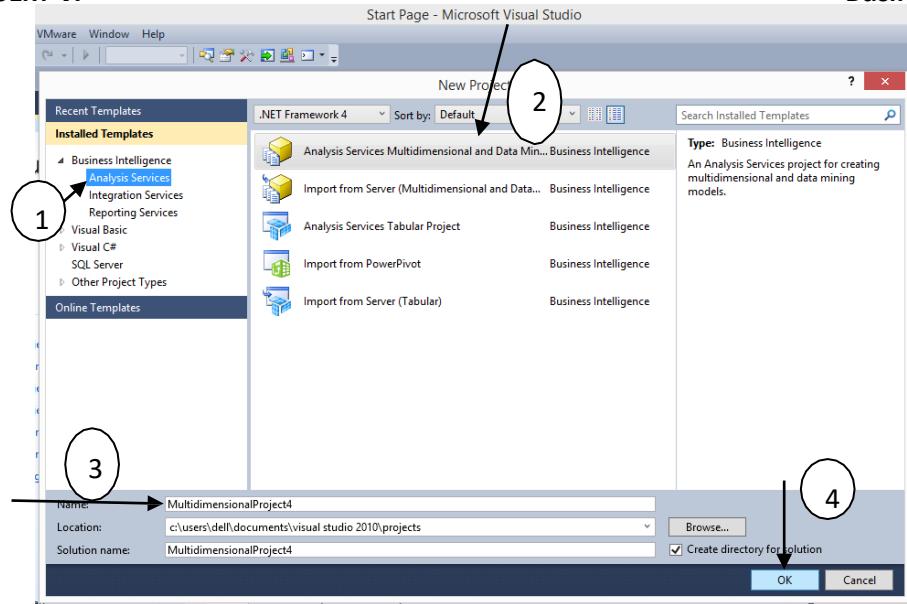
Practical No.3(a). Create the Data staging area for the selected database.**Staging Area:**

A staging area, or landing zone, is an intermediate storage area used for data processing during the extract, transform and load (ETL) process. The data staging area sits between the data source(s) and the data target(s), which are often data warehouses, data marts, or other data repositories



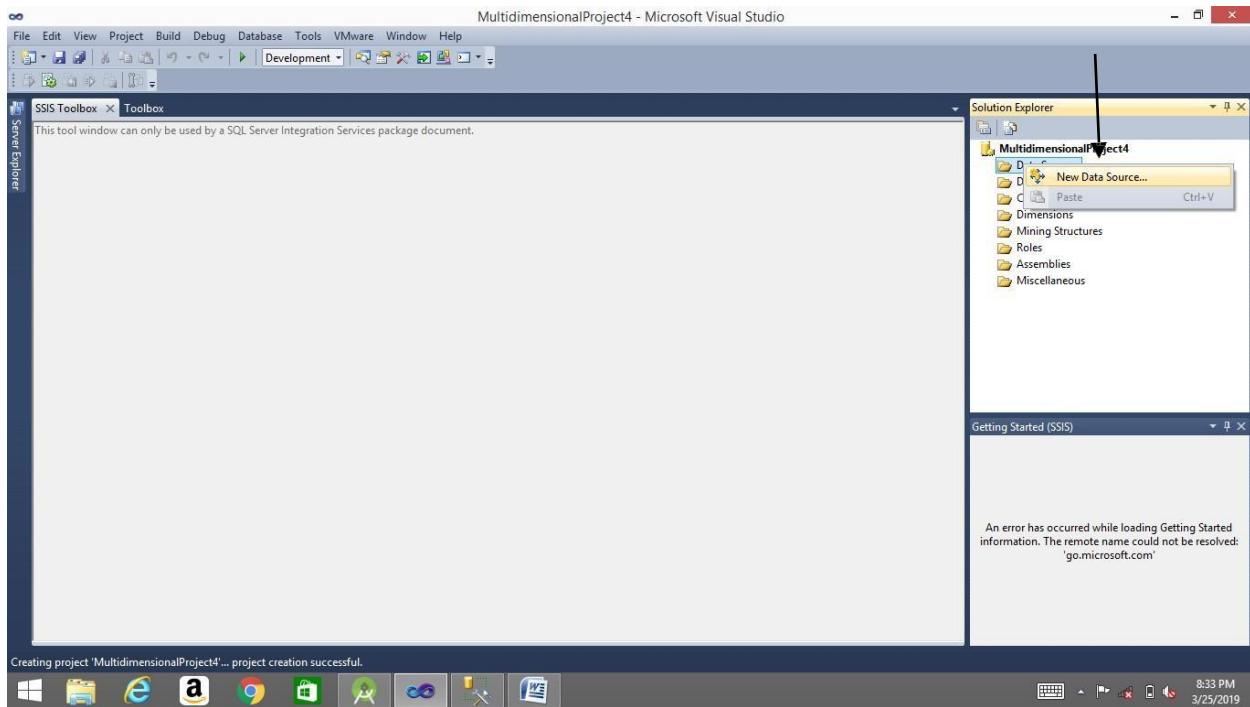
To create a staging area, we have to use SSDT tool and create a new project as

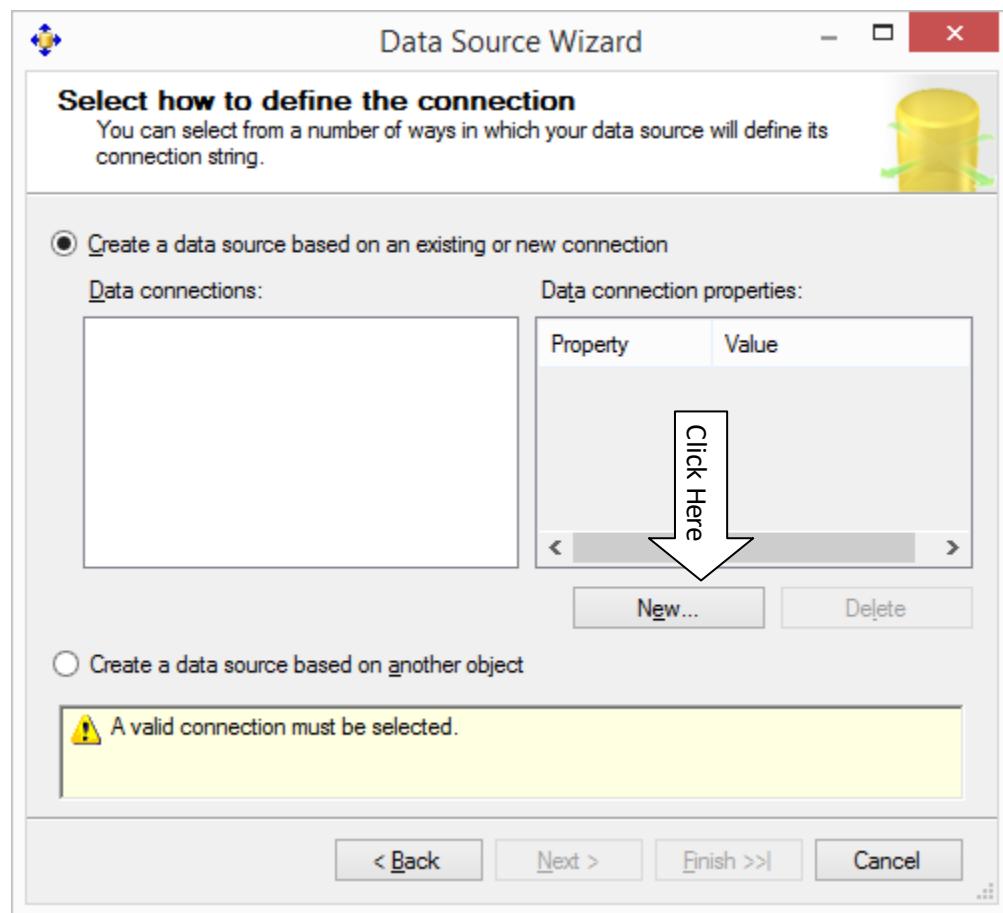
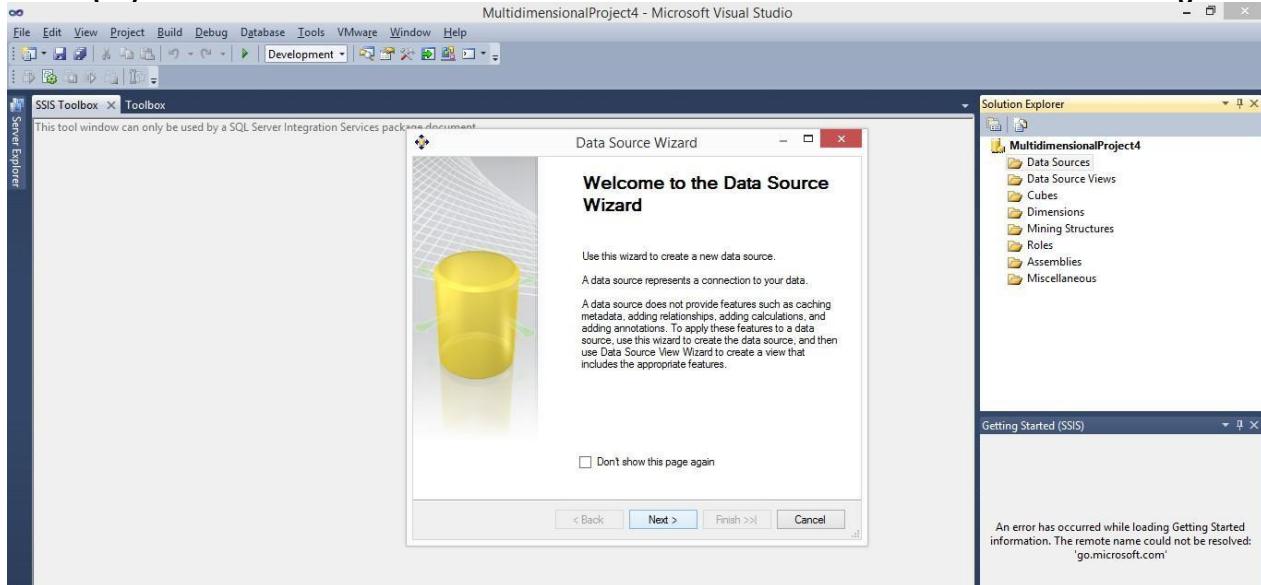




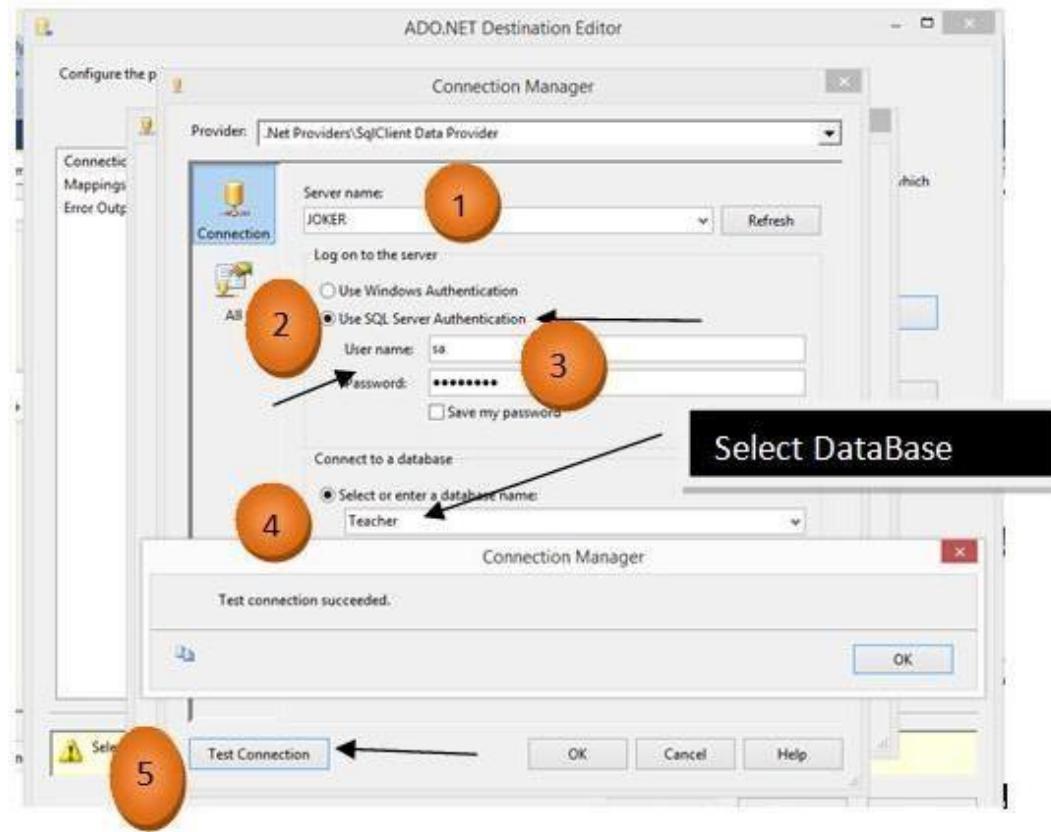
Step1 : Create a Data Source For Staging Area.

Right Click on Data Source Option and Select New Data Source from Solution Explorer as

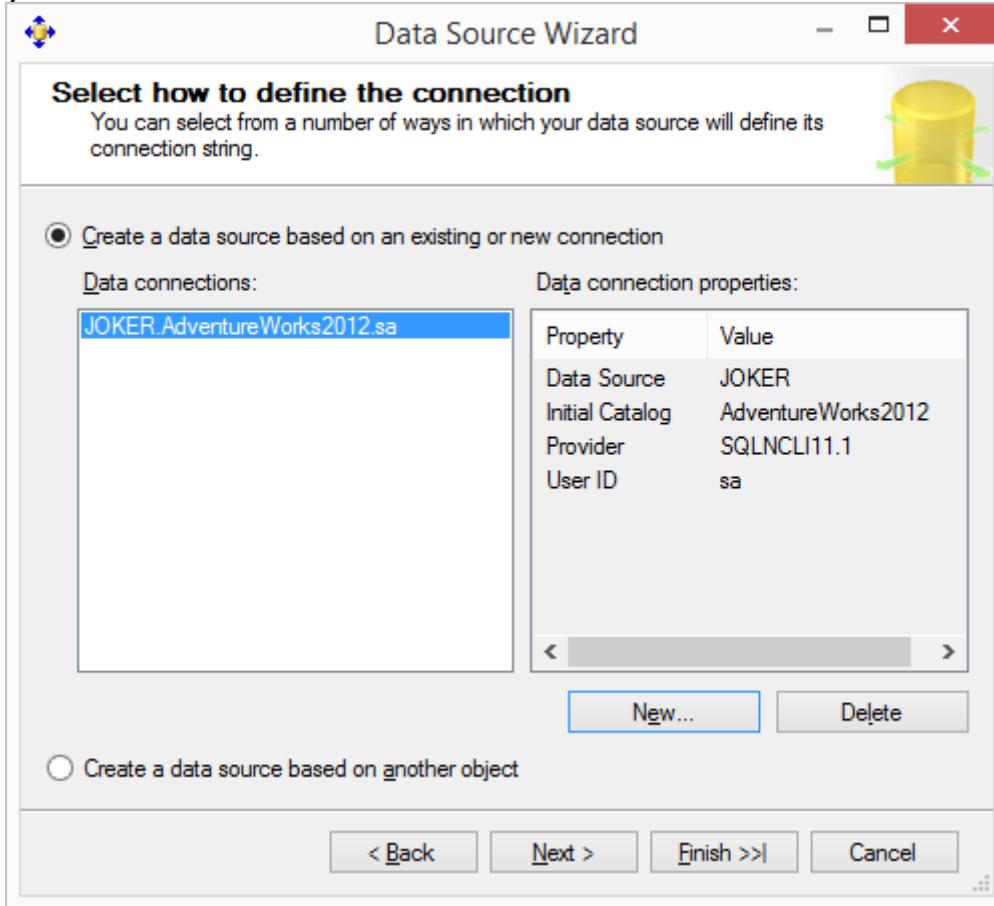




Enter The SQL Server instance name = JOKER in this case

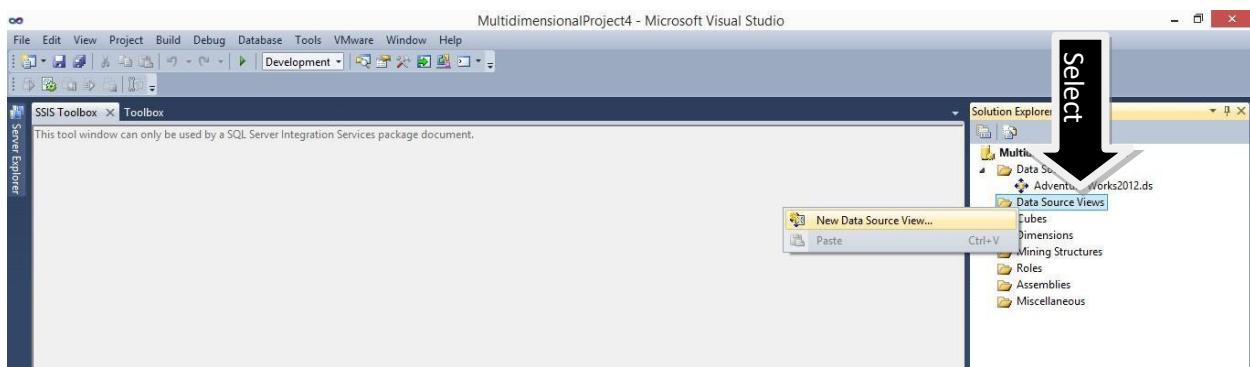


Click Ok

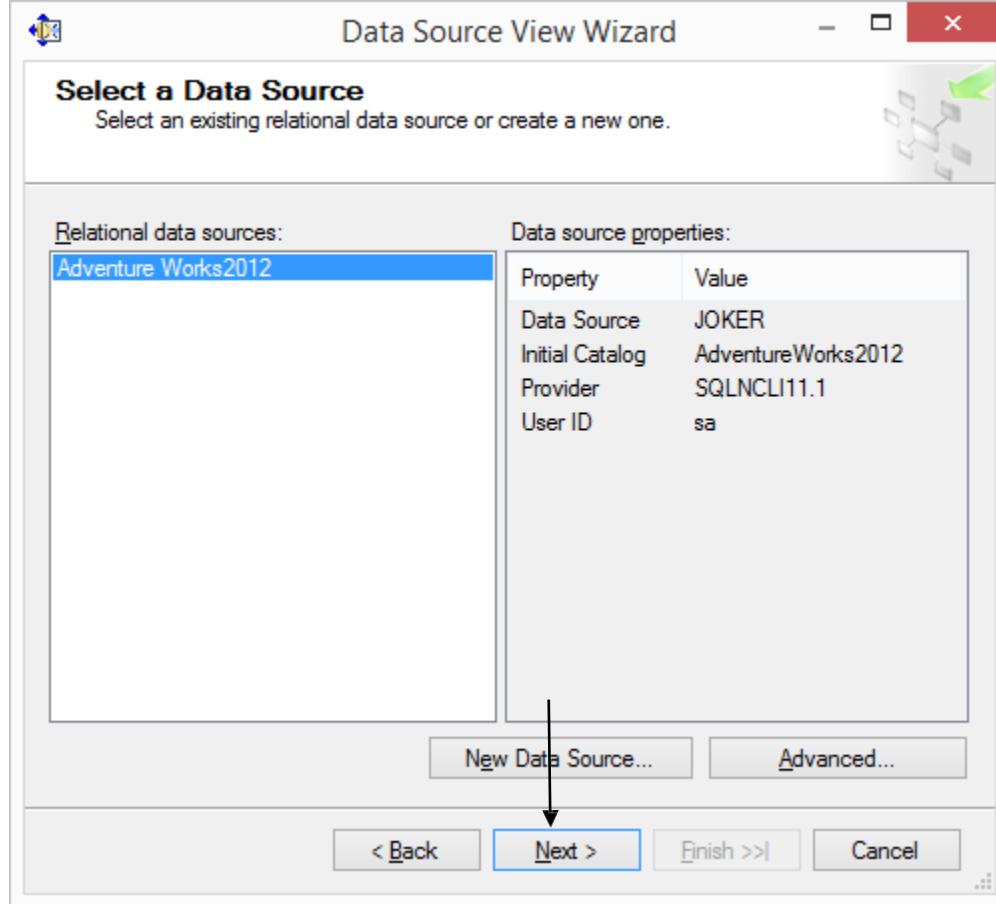


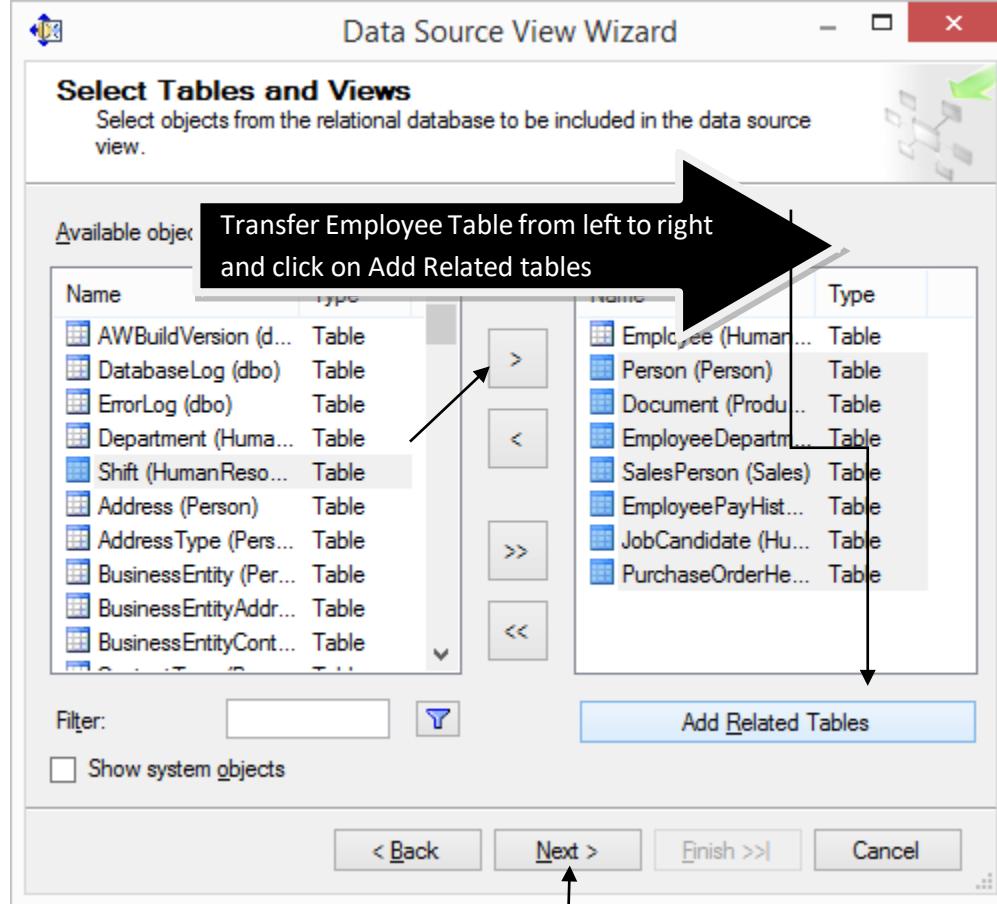
Step2: Now we will create data source View to visualize the Data Source.

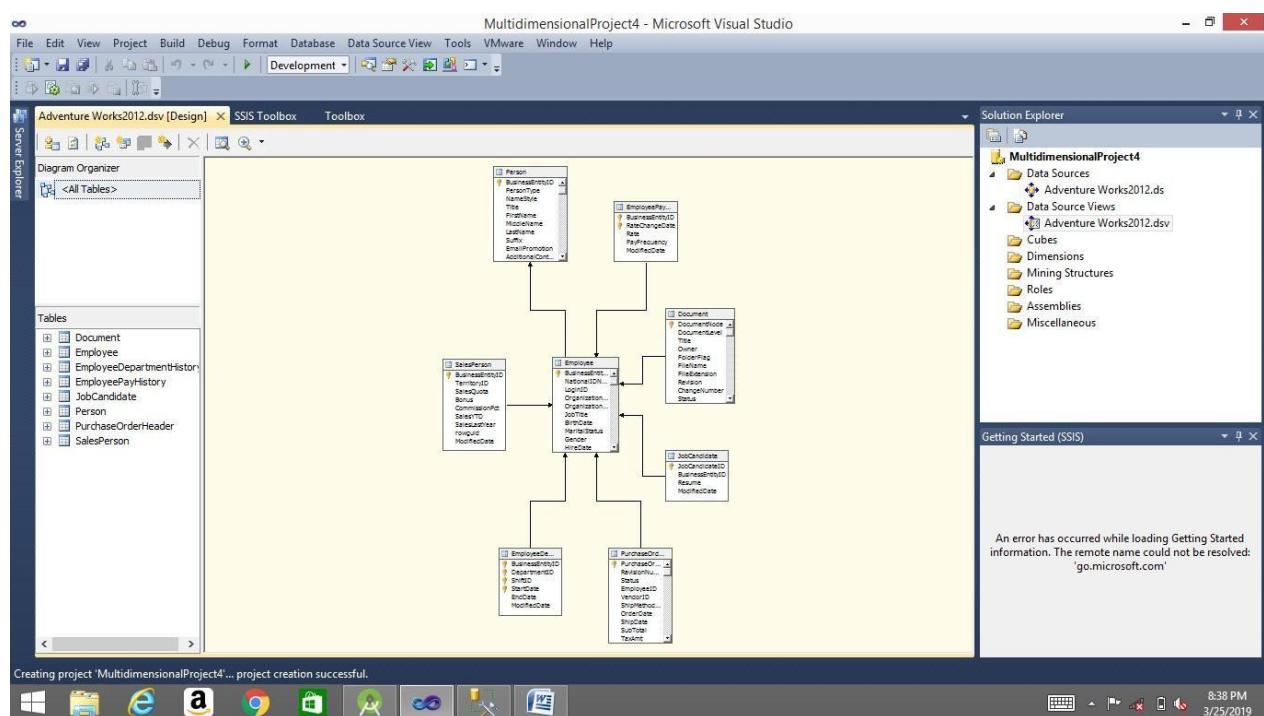
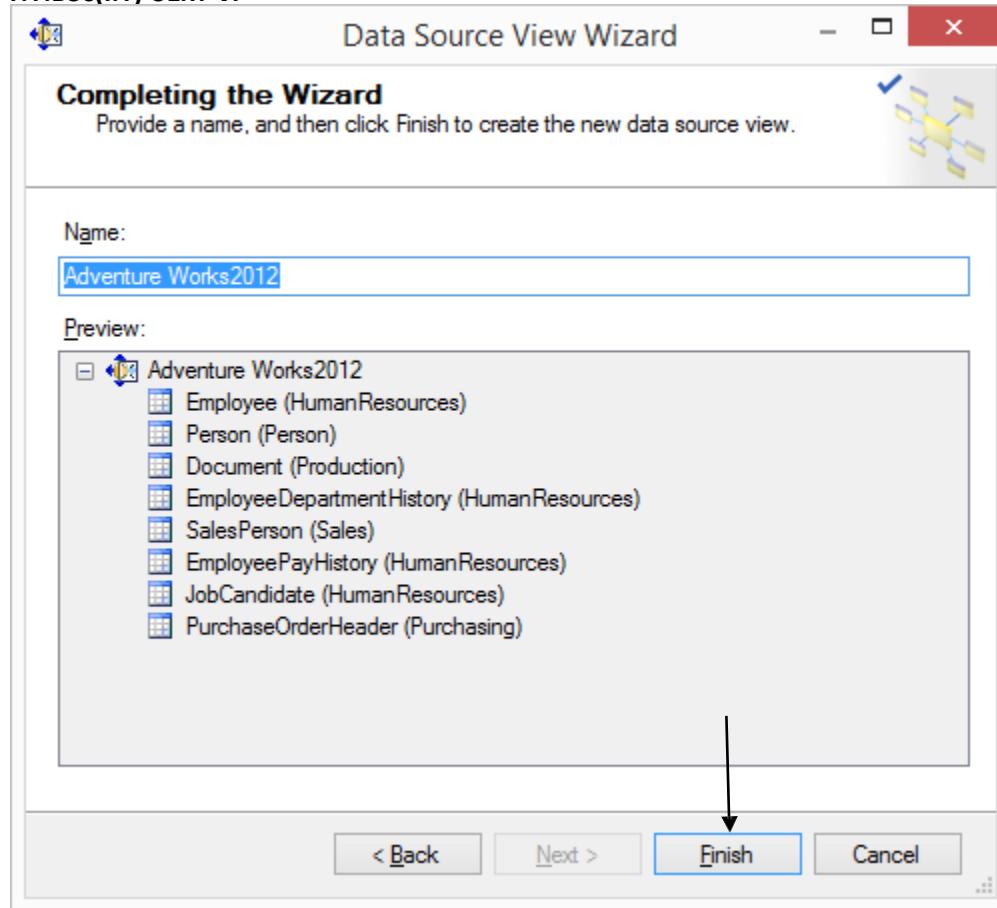
For this right click on Data Source View Option and select New Data Source View











This is a Staging area showing in START SCHEMA.

Practical No.3(a) done.

Practical No.3(b). Create the cube with suitable dimension and fact tables based on ROLAP, MOLAP and HOLAP model.

Dimension Table

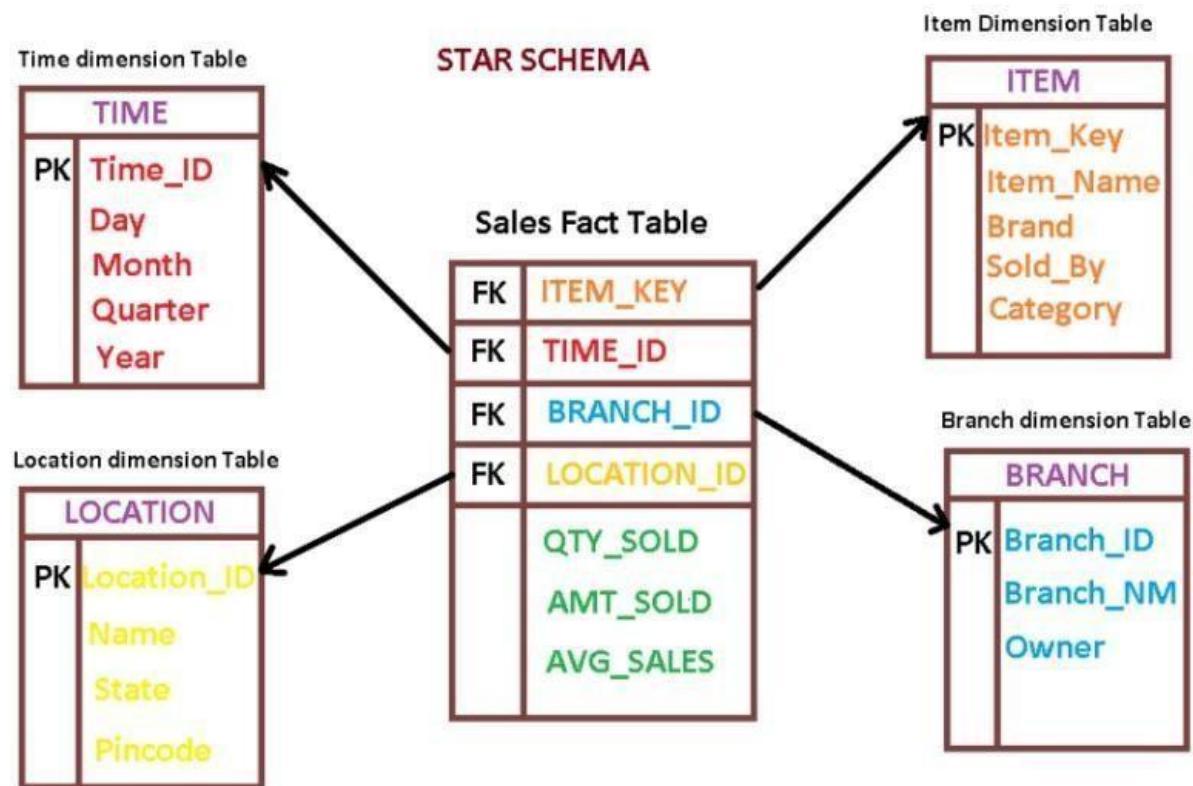
A Dimension Table is a table in a star schema of a data warehouse. Data warehouses are built using dimensional data models which consist of fact and dimension tables. Dimension tables are used to describe dimensions; they contain dimension keys, values and attributes.

Fact Table

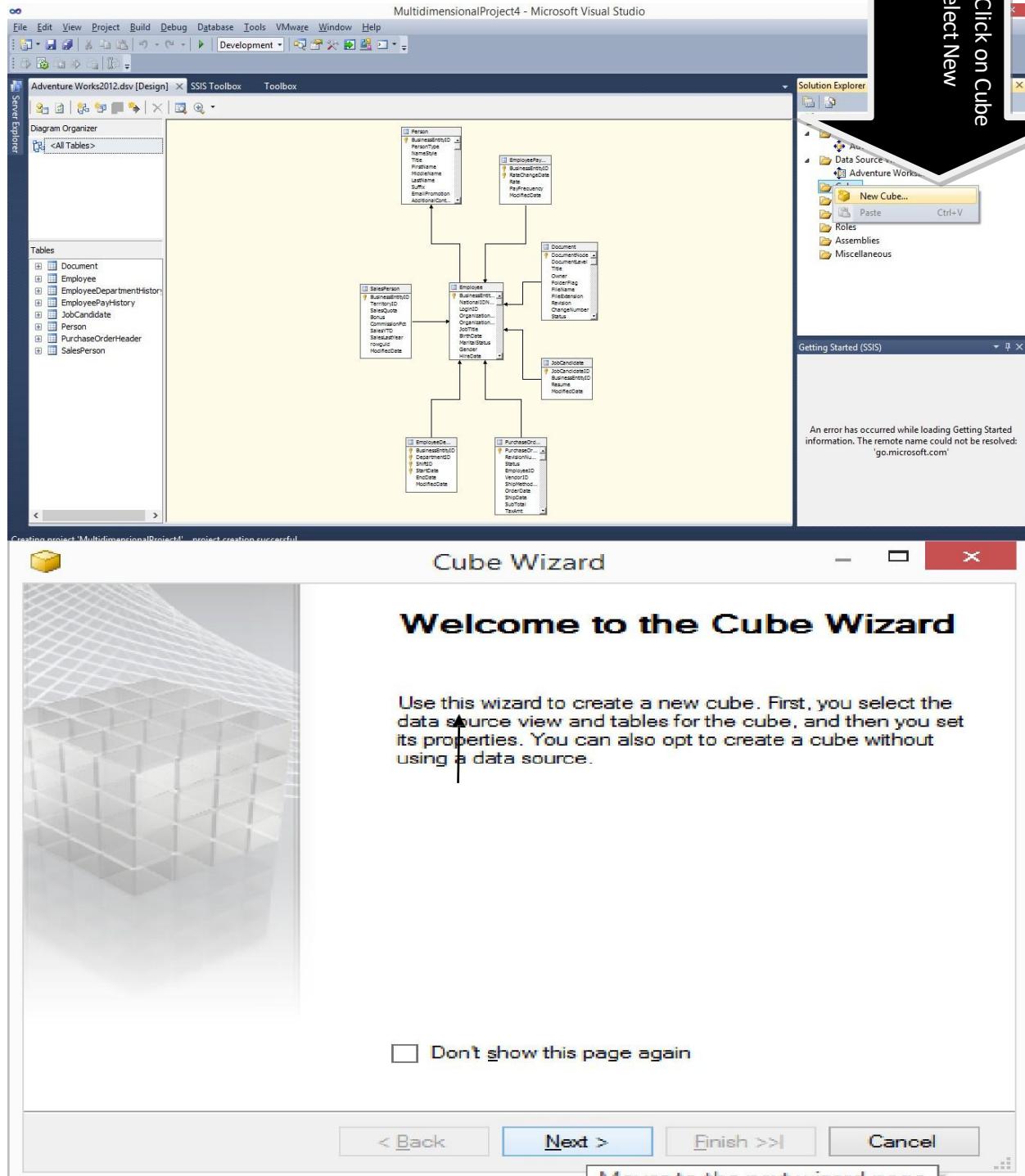
A fact table is found at the center of a star schema or snowflake schema surrounded by dimension tables. A fact table consists of facts of a particular business process e.g., sales revenue by month by product. Facts are also known as measurements or metrics. A fact table record captures a measurement or a metric.

Cube

Cubes are data processing units composed of fact tables and dimensions from the data warehouse. They provide multidimensional views of data, querying and analytical capabilities to clients.



To create a cube right click on Cube and select New Cube as



Cube Wizard

Select Creation Method
Cubes can be created by using existing tables, creating an empty cube, or generating tables in the data source.

How would you like to create the cube?

Use existing tables
 Create an empty cube
 Generate tables in the data source

Template:
(None)

Description:
Create a cube based on one or more tables in a data source.

< Back Next > Finish >> Cancel

Cube Wizard

Select Measure Group Tables
Select a data source view or diagram and then select the tables that will be used for measure groups.

Data source view:
Adventure Works2012

Measure group tables:
 Employee
 Person
 Document
 EmployeeDepartmentHistory
 SalesPerson
 EmployeePayHistory
 JobCandidate
 PurchaseOrderHeader

Suggest

< Back Next > Finish >> Cancel

Cube Wizard

Select New Dimensions
Select new dimensions to be created, based on available tables.

Dimension
 Person
 Person
 Employee
 Employee

< Back Next > Finish >> Cancel

Cube Wizard

Completing the Wizard
Name the cube, review its structure, and then click Finish to save the cube.

Adventure Works2012

Cube name:
Adventure Works2012

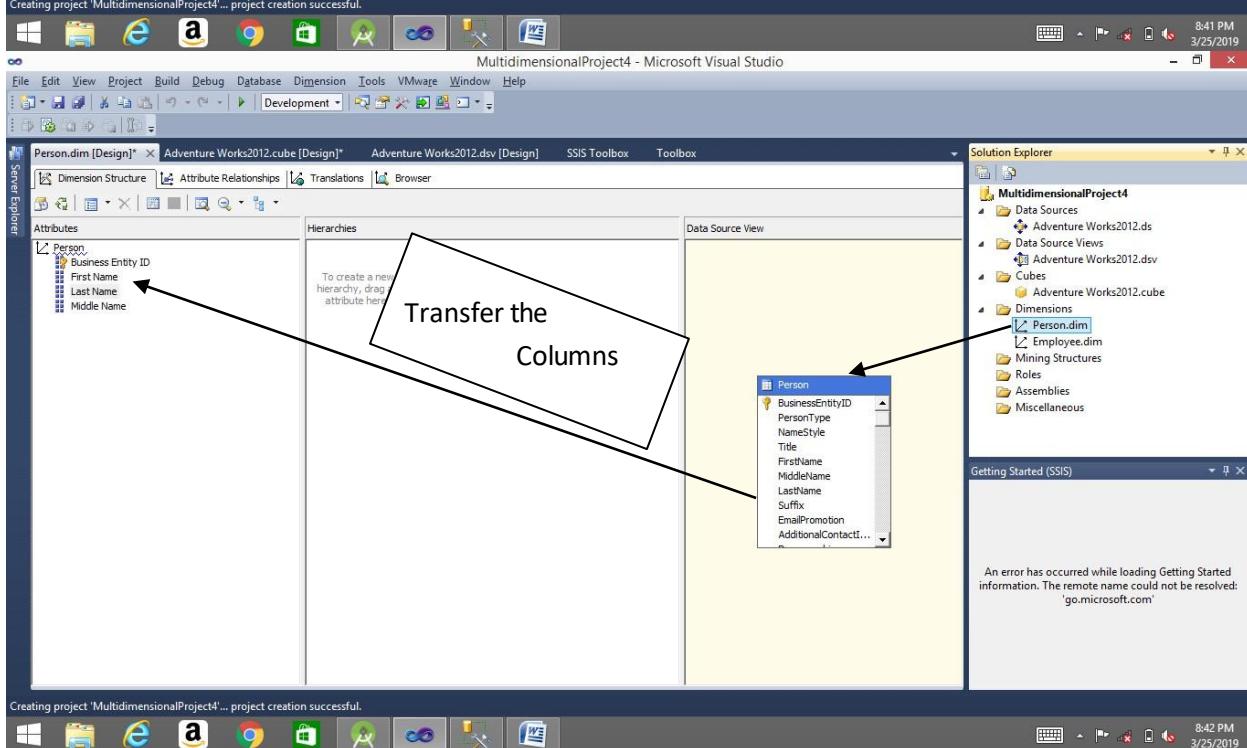
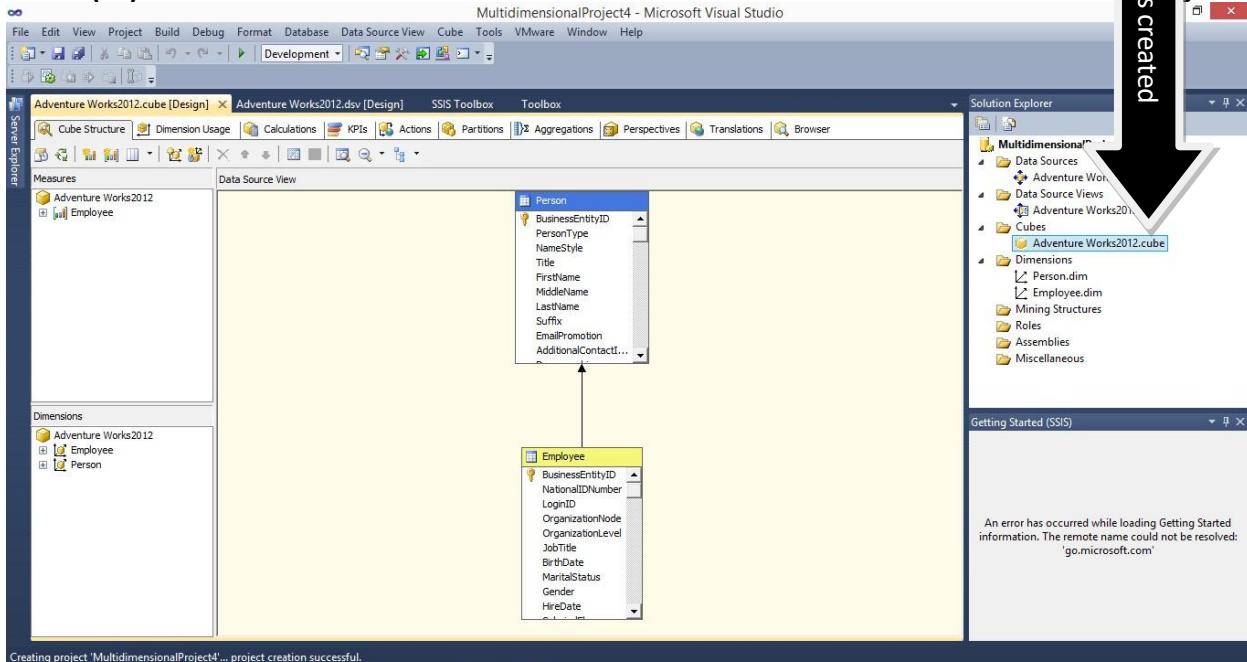
Preview:

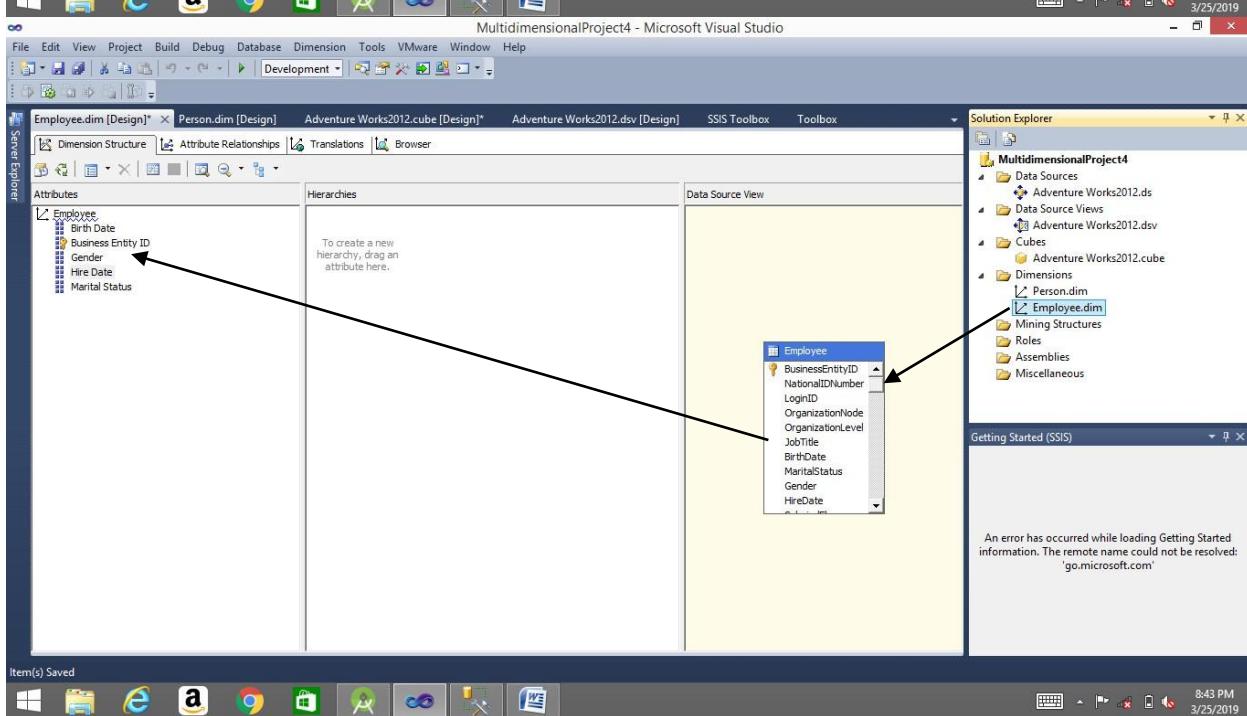
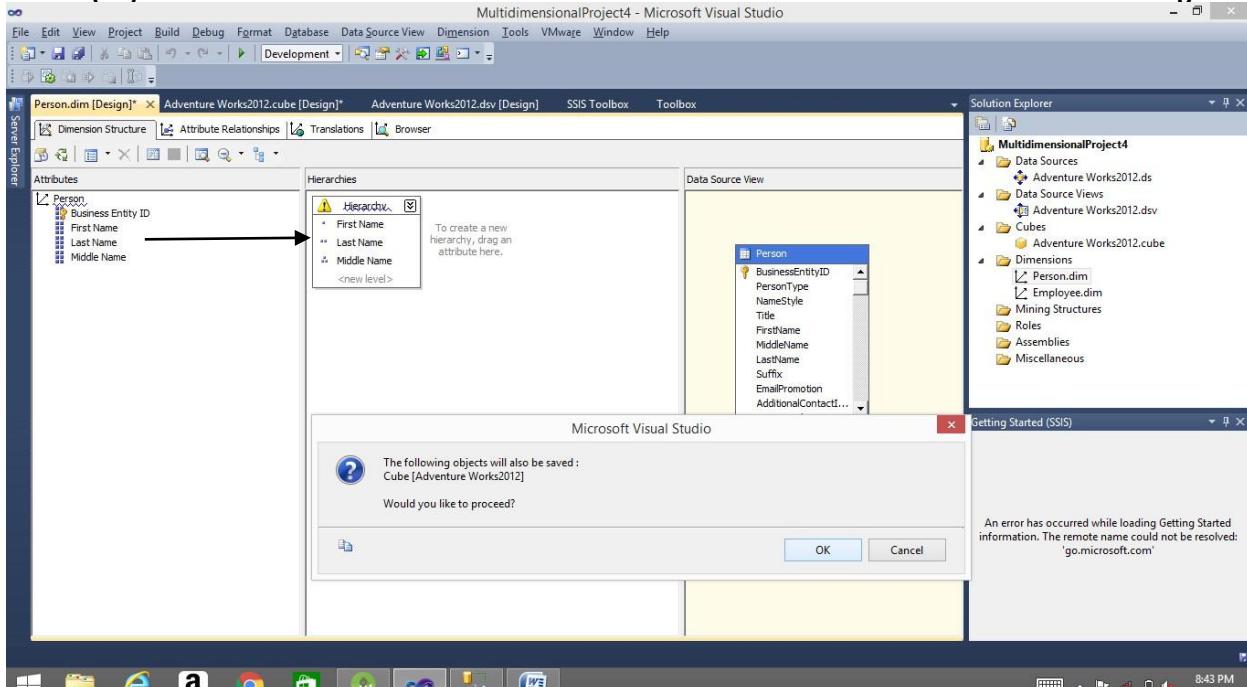
Measure groups
 Employee
 Organization Level
 Vacation Hours
 Sick Leave Hours
 Employee Count

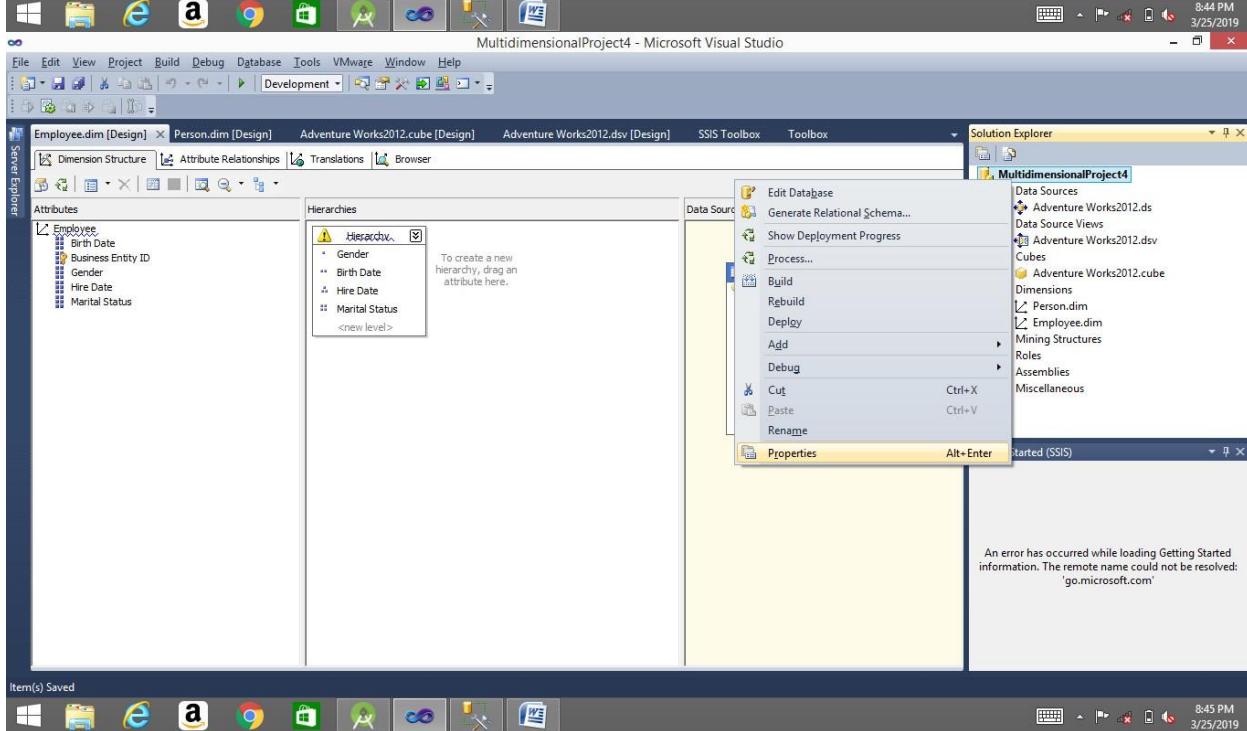
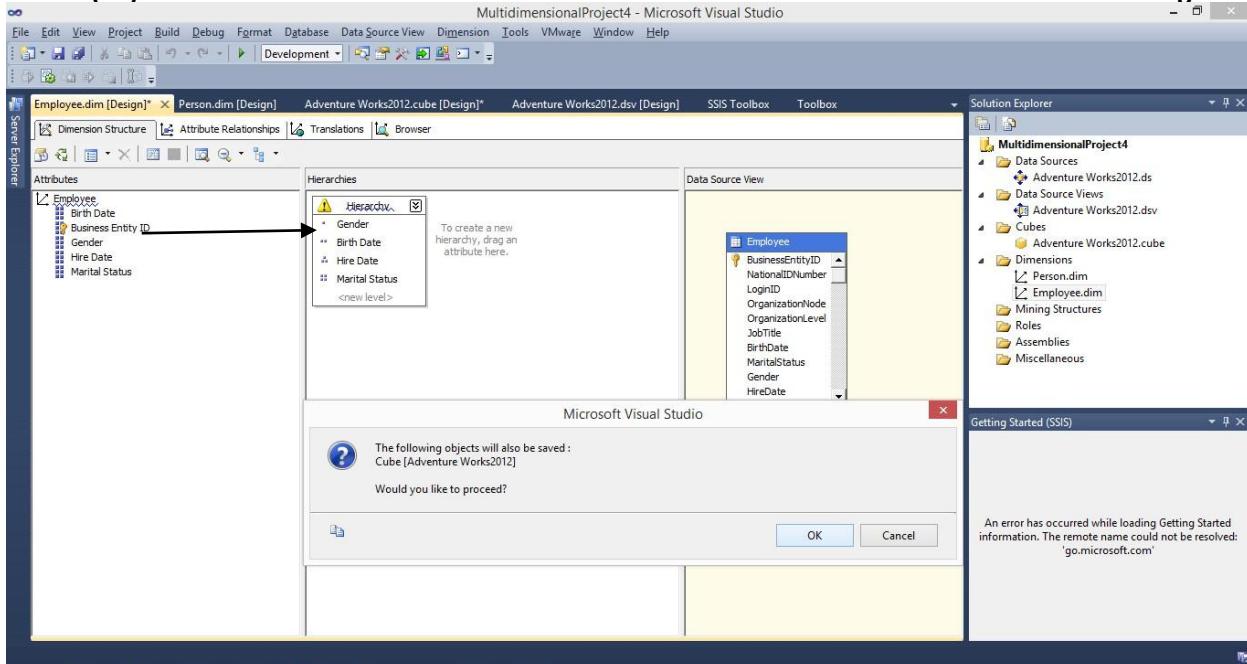
Dimensions
 Person
 Employee

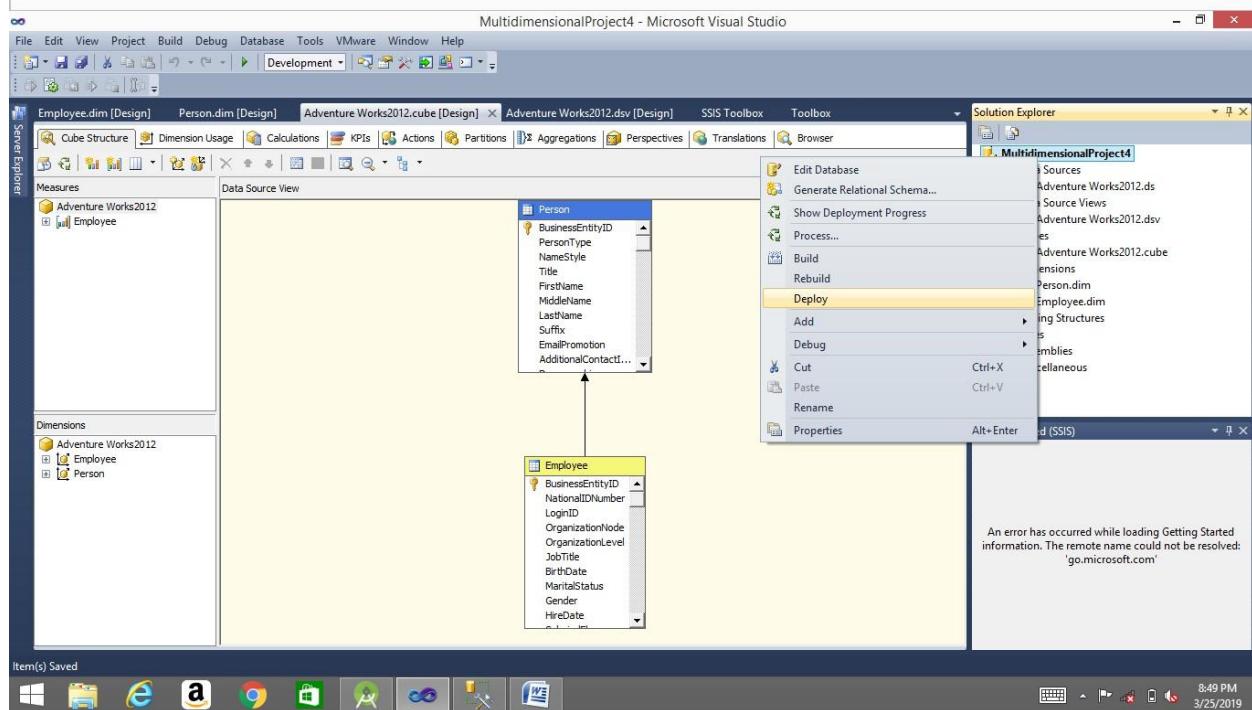
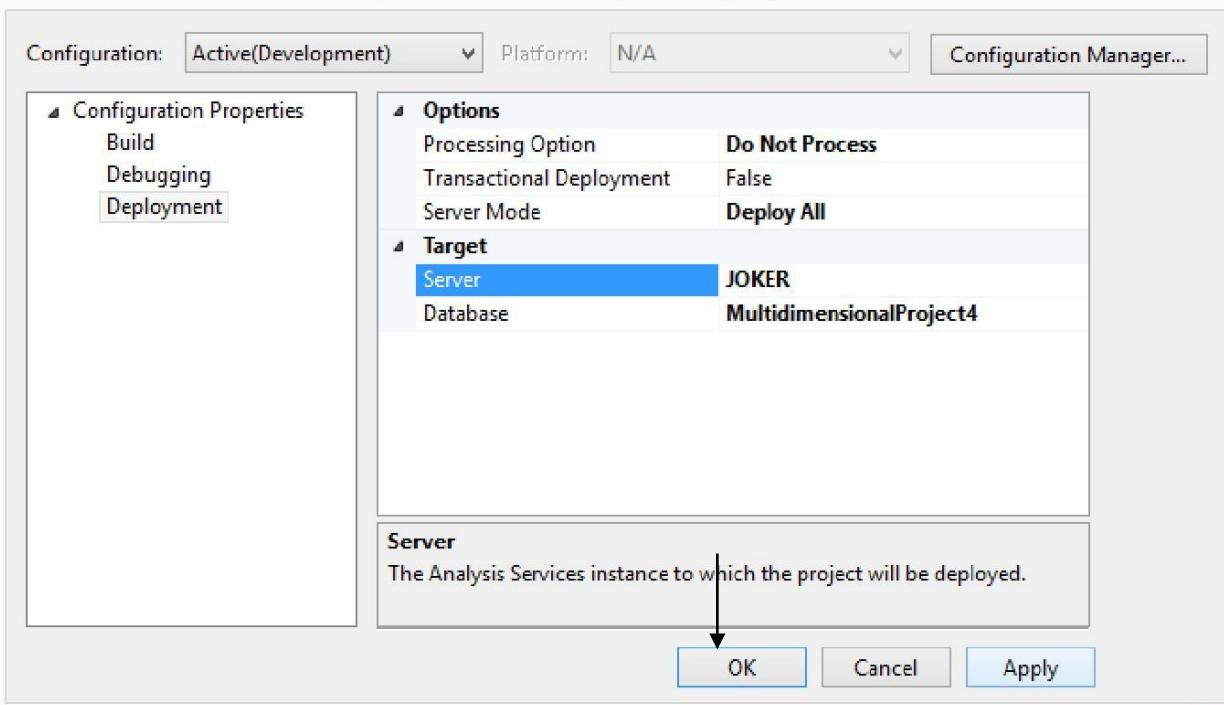
< Back Next > Finish Cancel

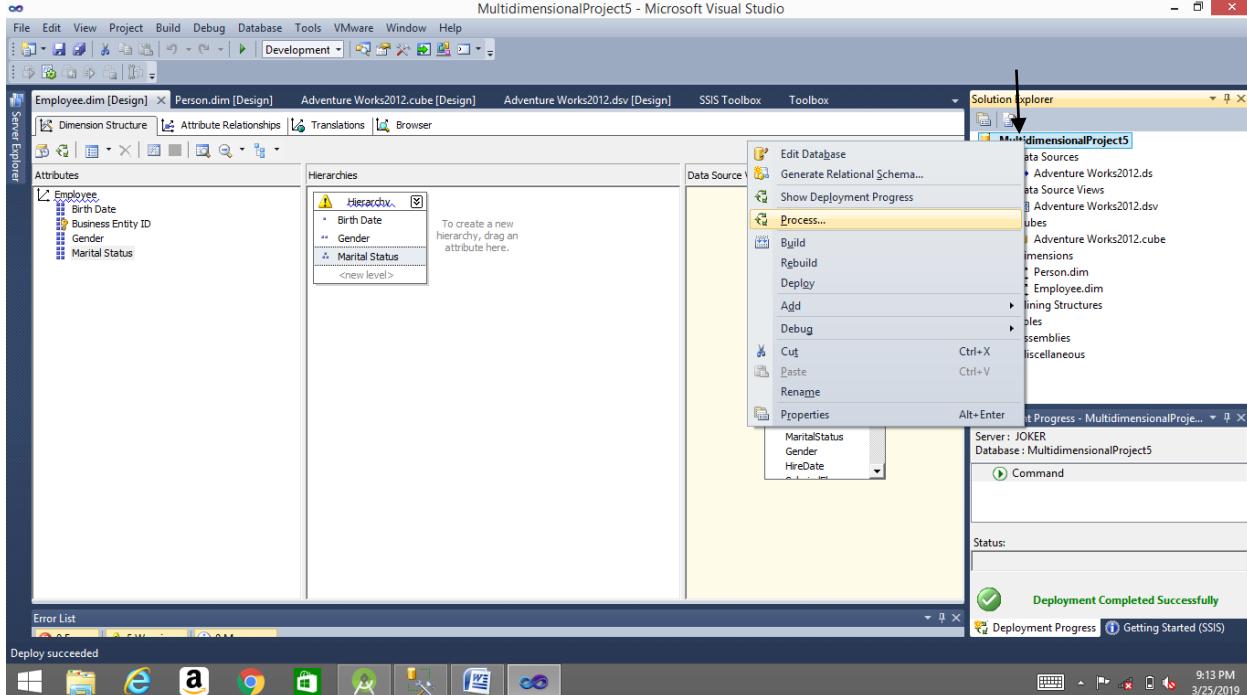
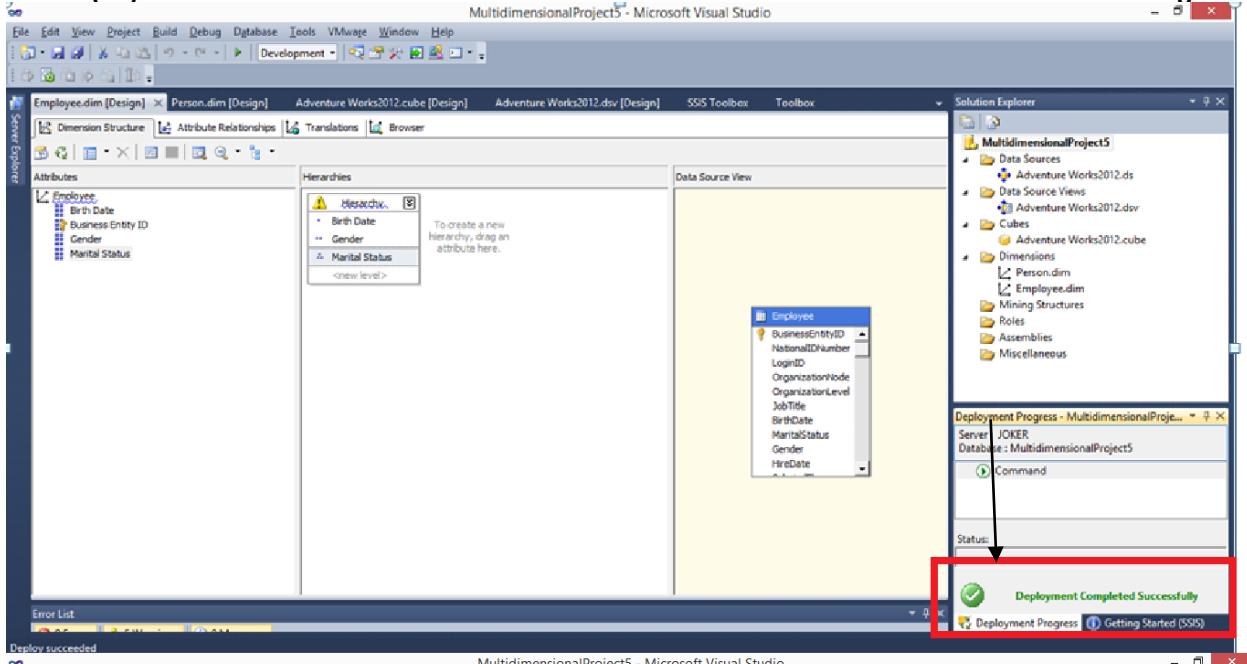
T.Y.BSc(I.T) SEM-VI











Object list:

Object Name	Type	Process Options	Settings
MultidimensionalProject5	Database	Process Full	

Batch Settings Summary

Processing order:
Parallel

Transaction mode:
(Default)

Dimension errors:
(Default)

Dimension key error log path:
(Default)

Process affected objects:
Do not process

Remove Impact Analysis...

Run... Close

 9:13 PM
3/25/2019

Process Progress

Command

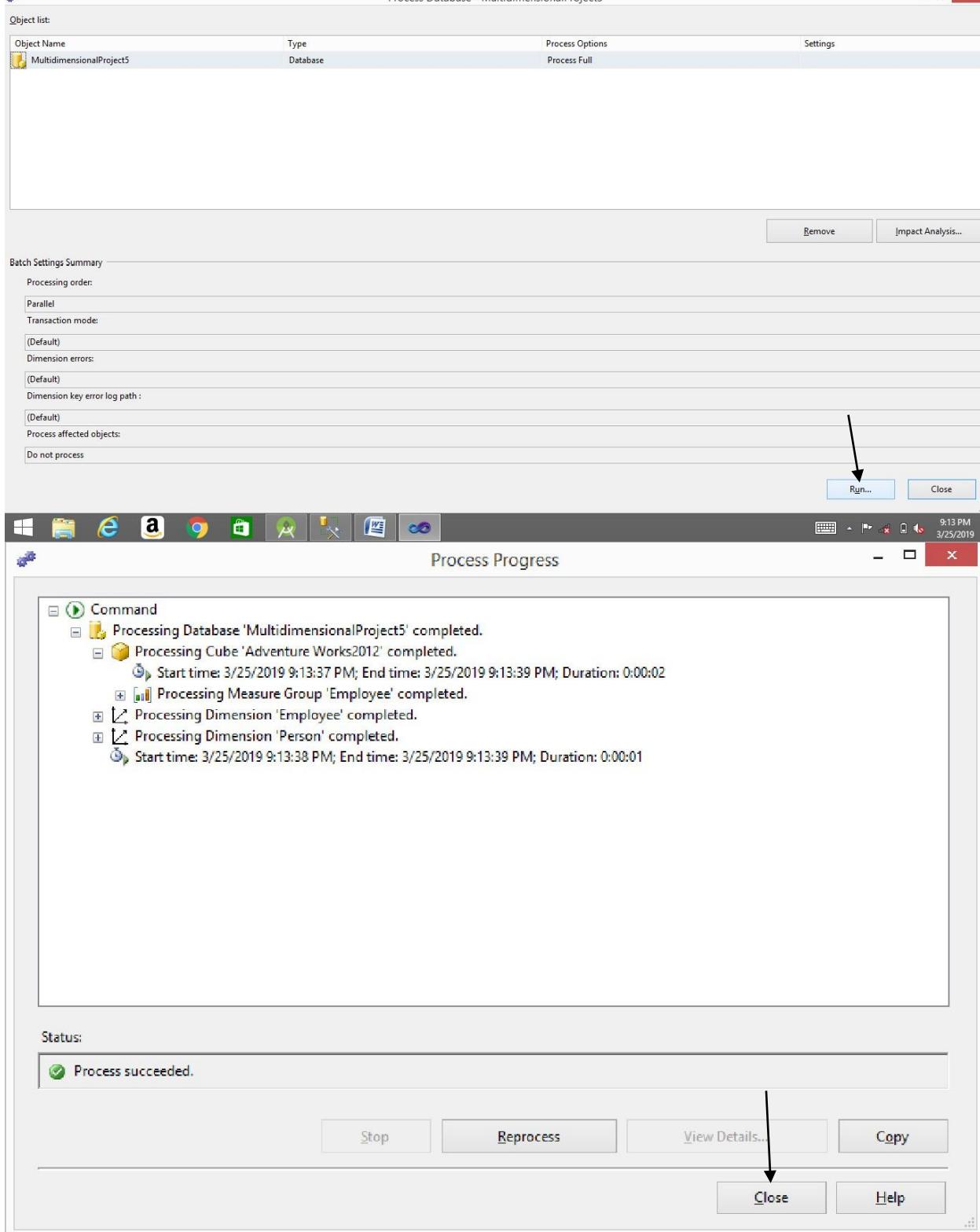
- Processing Database 'MultidimensionalProject5' completed.
- Processing Cube 'Adventure Works2012' completed.
 - Start time: 3/25/2019 9:13:37 PM; End time: 3/25/2019 9:13:39 PM; Duration: 0:00:02
 - Processing Measure Group 'Employee' completed.
 - Processing Dimension 'Employee' completed.
 - Processing Dimension 'Person' completed.
 - Start time: 3/25/2019 9:13:38 PM; End time: 3/25/2019 9:13:39 PM; Duration: 0:00:01

Status:

Process succeeded.

Stop Reprocess View Details... Copy

Close Help



MultidimensionalProject3 - Microsoft Visual Studio

File Edit View Project Build Debug Database Cube Tools VMware Window Help

Employee.dim [Design] Person.dim [Design] Adventure Works2012.cube [Design] Adventure Works2012.ds [Design] Toolbox

Cube Structure Dimension Usage Calculations KPIs Actions Partitions Aggregations Perspectives Translations Browser

Language: Default

Adventure Works2012

Dimension Hierarchy Operator Filter Expression Parameter

<Select dimension>

Metadata

Measure Group: <All>

- Adventure Works2012
 - Measures
 - KPIs
 - Employee
 - Person

Calculated Members

Drag levels or measures here to add to the query.

Error List

Deploy succeeded

9:16 PM 3/25/2019

Solution Explorer

MultidimensionalProject3

- Data Sources
 - Adventure Works2012.ds
- Data Source Views
 - Adventure Works2012.ds
- Cubes
 - Adventure Works2012.cube
- Dimensions
 - Person.dim
 - Employee.dim
- Mining Structures
- Roles
- Assemblies
- Miscellaneous

Deployment Progress - MultidimensionalProj... Server: JOKER Database: MultidimensionalProject3

Command

Status:

Deployment Completed Successfully

Deployment Progress Getting Started (SSIS)

MultidimensionalProject3 - Microsoft Visual Studio

File Edit View Project Build Debug Database Cube Tools VMware Window Help

Employee.dim [Design] Person.dim [Design] Adventure Works2012.cube [Design] Adventure Works2012.ds [Design] Toolbox

Cube Structure Dimension Usage Calculations KPIs Actions Partitions Aggregations Perspectives Translations Browser

Language: Default

Adventure Works2012

Dimension Hierarchy Operator Filter Expression Parameter

<Select dimension>

Metadata

Measure Group: <All>

- Adventure Works2012
 - Employee Count
 - Organization Level
 - Sick Leave Hours
 - Vacation Hours
 - KPIs
 - Employee
 - Birth Date
 - Business Entity ID
 - Gender
 - Hire Date
 - Marital Status
 - Hierarchy

Calculated Members

Birth Date	Gender	Marital Status	Hire Date	Employee Count	Organization Level	Sick Leave Hours	Vacation Hours
1951-10-17	M	M	2011-01...	1	2	27	14
1952-03-02	M	M	2010-02...	1	4	53	66
1952-05-12	M	M	2010-01...	1	4	49	58
1952-09-27	F	M	2008-01...	1	3	22	5
1953-04-30	M	M	2010-01...	1	4	22	5
1954-04-24	F	M	2010-03...	1	4	65	91
1955-01-30	M	S	2010-01...	1	4	35	30
1956-01-16	F	S	2007-12...	1	3	61	82
1956-03-26	M	S	2008-01...	1	4	64	88
1956-03-29	F	M	2008-03...	1	4	63	87
1956-04-01	M	M	2008-02...	1	3	59	79
1956-04-04	F	M	2008-03...	1	4	63	86
1956-06-04	F	M	2008-01...	1	4	61	83
1956-07-11	M	S	2008-02...	1	4	62	85
1956-08-07	M	M	2008-03...	1	3	60	81
1956-09-20	M	M	2009-02...	1	4	62	84

Error List

Deploy succeeded

9:18 PM 3/25/2019

Solution Explorer

MultidimensionalProject3

- Data Sources
 - Adventure Works2012.ds
- Data Source Views
 - Adventure Works2012.ds
- Cubes
 - Adventure Works2012.cube
- Dimensions
 - Person.dim
 - Employee.dim
- Mining Structures
- Roles
- Assemblies
- Miscellaneous

Deployment Progress - MultidimensionalProj... Server: JOKER Database: MultidimensionalProject3

Command

Status:

Deployment Completed Successfully

Deployment Progress Getting Started (SSIS)