

Business Intelligence and Data Analytics Journal (2024-2025)

VidyaVikas Education Society's



VIKAS COLLEGE OF ARTS, SCIENCE & COMMERCE

Affiliated to University of Mumbai
RE-ACCREDITED 'A' GRADE BY NAAC
ISO 9001 : 2008 CERTIFIED

Vikas High School Marg, Kannamwar Nagar No 2, Vikhroli (E), Mumbai – 400083

Dr. R. K. Patra
Principal

Hon' ble: **Shri P. M. Raut**
Chairman. V. V. Edu. Society

This is to certify that, _____
student of T.Y.B.Sc. (Information Technology) (Semester-V) with college enrolled
Roll no. _____ / University Seat _____ has satisfactorily
completed the Project Dissertation work for the Subject Software Project Development
in the program of Information Technology from the UNIVERSITY OF MUMBAI for the
academic year 2024-2025.

Guided By

College Seal

Head Of Department

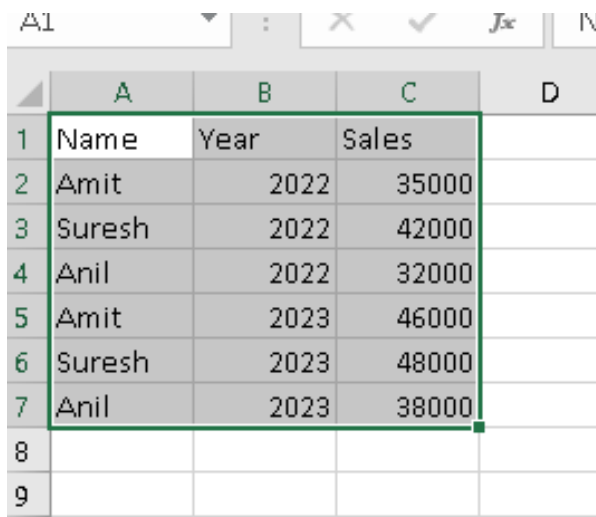
External Examiner

Index

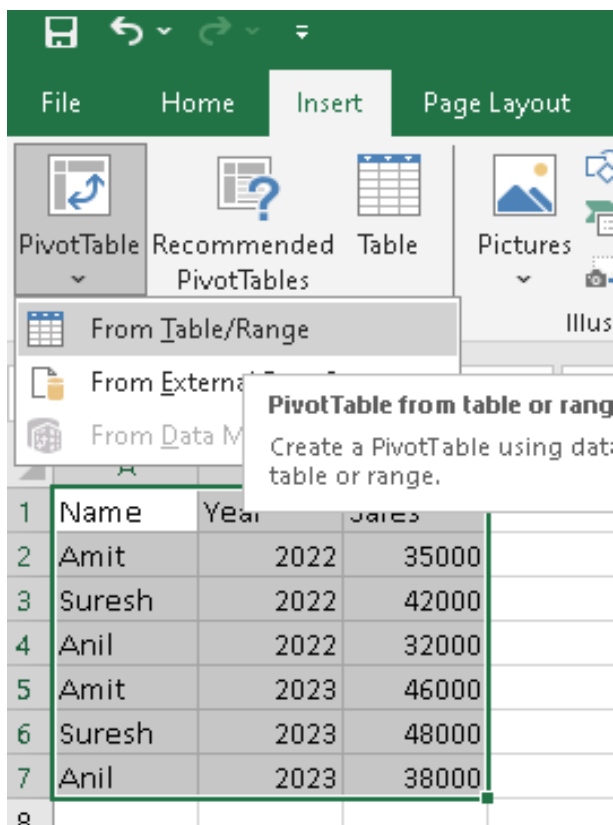
Sr.No	Practical.No	Title	Signature
1	1a	Import the data warehouse data in Microsoft Excel and create the pivot table and pivot chart.	
2	1b	Import the cube in Microsoft Excel and create the Pivot table and Pivot Chart to perform data analysis.	
3	2	Apply the what – if Analysis for data visualization. Design and generate necessary reports based on the data warehouse data. Use Excel.	
4	3	Perform the data classification using classification algorithm using R/Python.	
5	4	Perform the data clustering using clustering algorithm using R/Python.	
6	5	Perform the Linear regression on the given data warehouse data using R/Python.	
7	6	Perform the logistic regression on the given data warehouse data using R/Python.	
8	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).	
9	8a	Perform data visualization using Python on any sales data.	
10	8b	Perform data visualization using PowerBI on any sales data.	
11	9	Create the Data staging area for the selected database using SQL.	

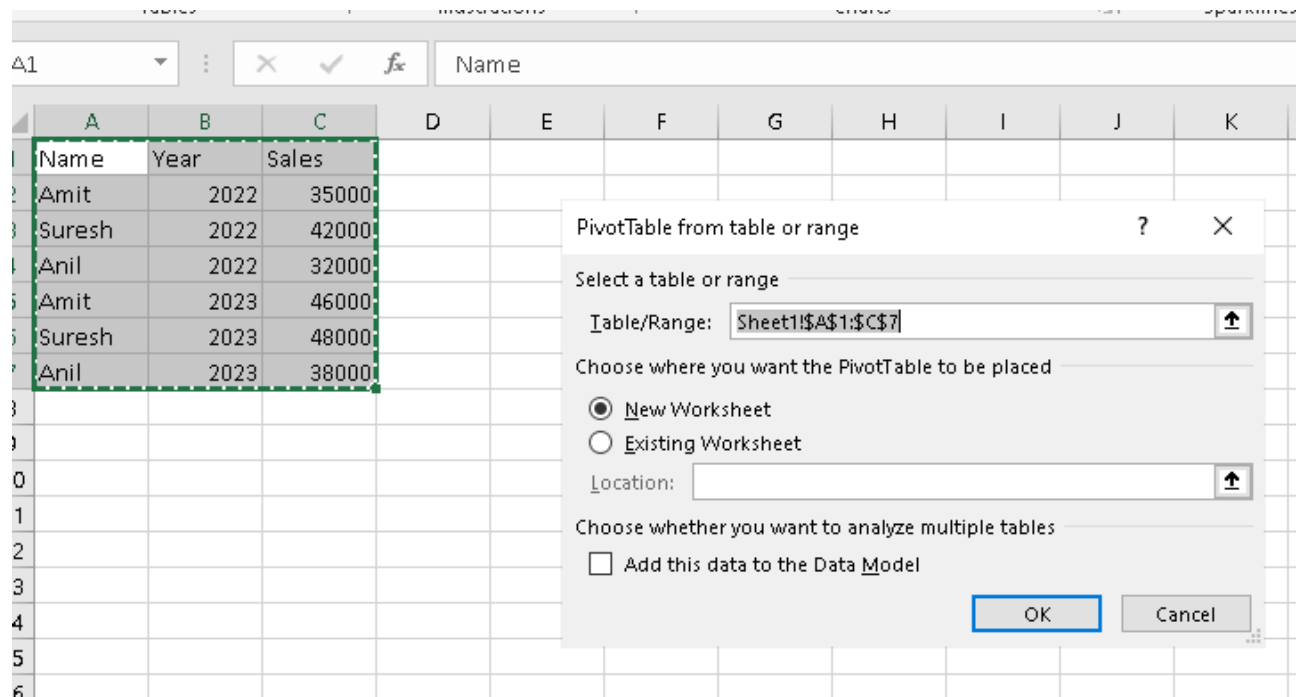
Practical: 1a

Aim: Import the data warehouse data in Microsoft Excel and create the pivot table and pivot chart.

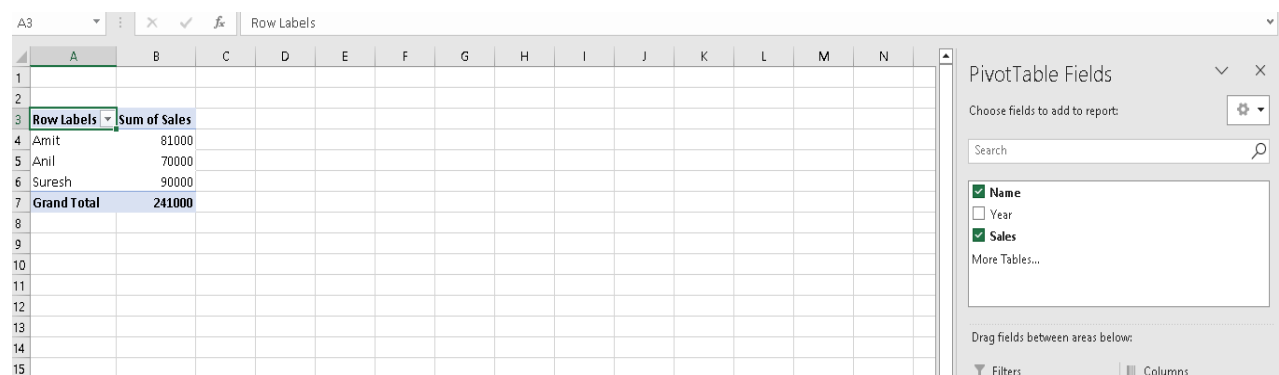


	A	B	C	D
1	Name	Year	Sales	
2	Amit	2022	35000	
3	Suresh	2022	42000	
4	Anil	2022	32000	
5	Amit	2023	46000	
6	Suresh	2023	48000	
7	Anil	2023	38000	
8				
9				





Click on ok button.



Click on Pivott Chart and select any chart.

pract1_2025_Bk.xlsx - Excel

PivotTable Tools: PivotTable Analyze Design

PivotTable Name: PivotTable4 Active Field: Name

Options Field Settings Show Details

PivotTable Active Field Group

Row Labels

	A	B	C	D	E
1					
2					
3	Row Labels	Sum of Sales			
4	Amit	81000			
5	Anil	70000			
6	Suresh	90000			
7	Grand Total	241000			
8					
9					
10					
11					
12					
13					
14					
15					
16					
17					
18					
19					
20					
21					
22					
23					

Insert Chart

All Charts

Recent Templates

Column

Line

Pie

Bar

Area

XY (Scatter)

Stock

Surface

Radar

Treemap

Sunburst

Histogram

Box & Whisker

Waterfall

Combo

Clustered Column

Total

PivotTable Fields

Choose fields to add to report:

Search

☒ Name

☐ Year

☒ Sales

More Tables...

Drag fields between areas below:

Filters

Columns

Rows

Values

Name

Sum of Sales

Defer Layout Update

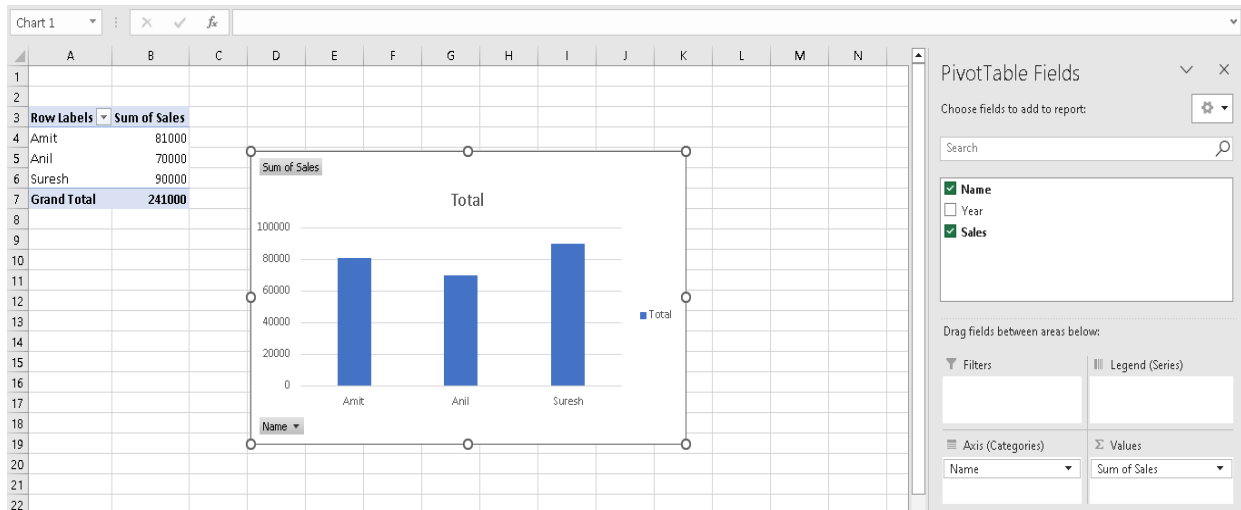
Update

Sheet2 Sheet3 Sheet4 Sheet1

Ready Accessibility: Investigate

Type here to search

19:22 20-03-2025

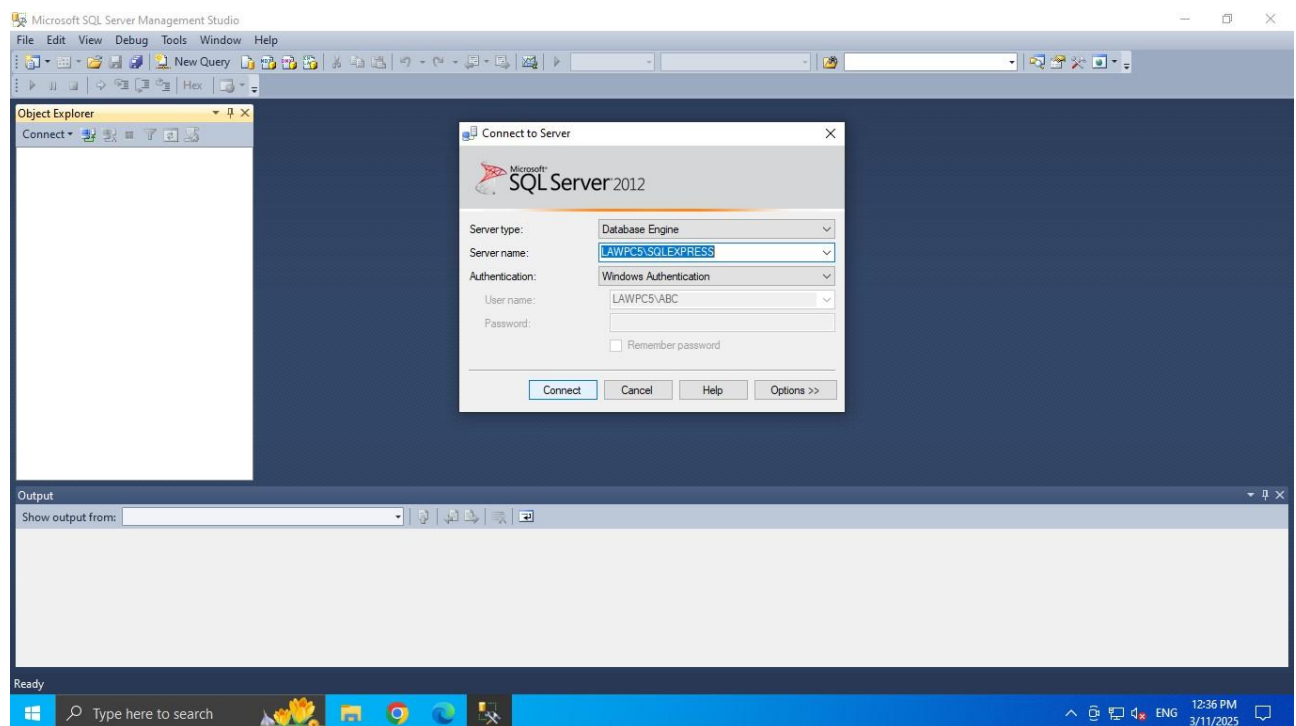


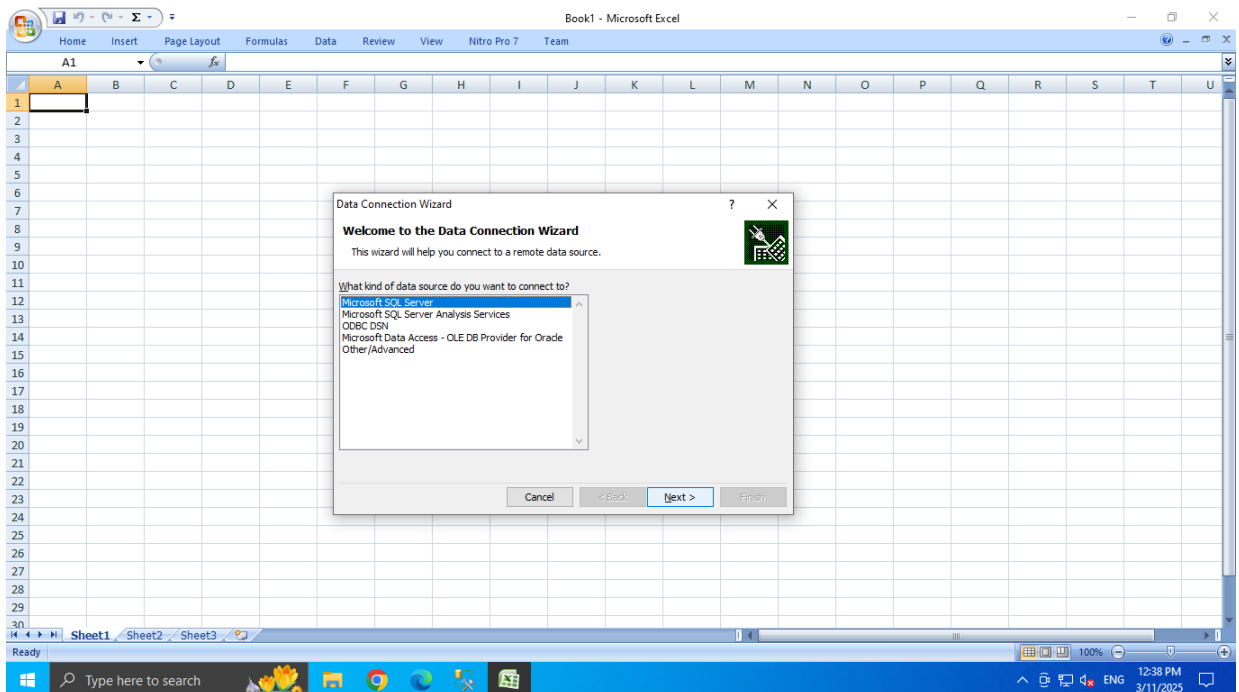
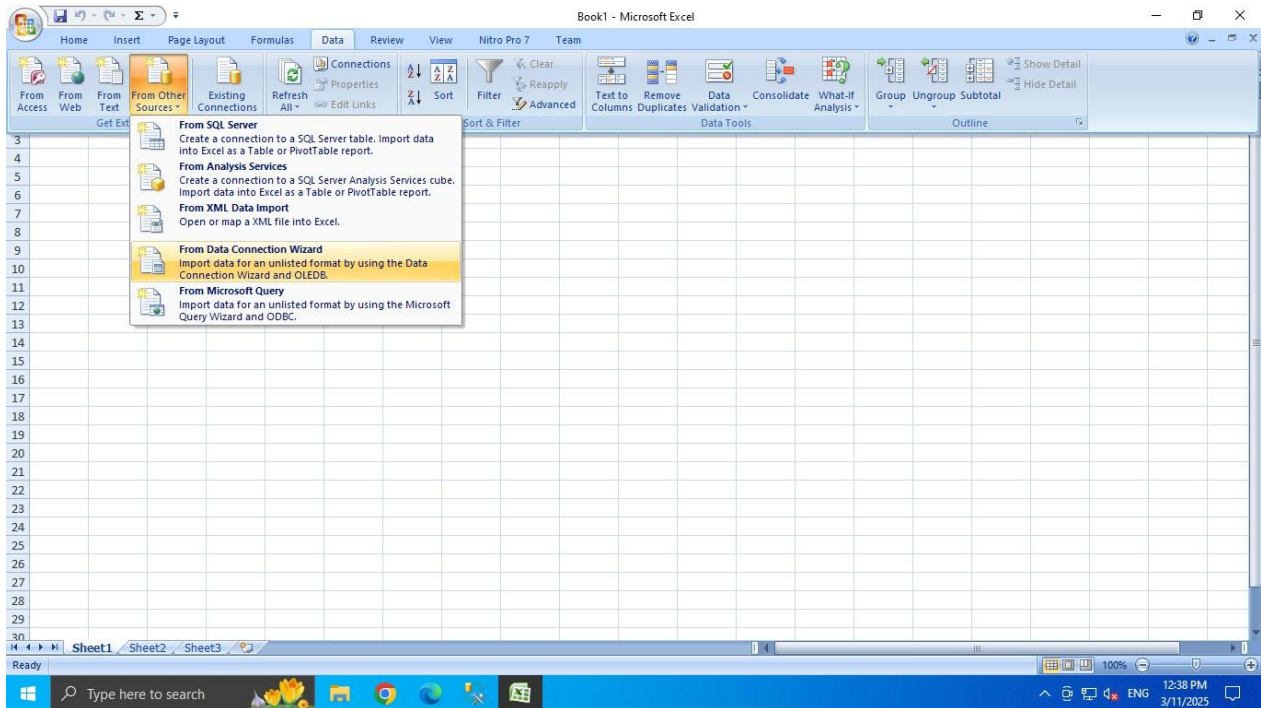
Practical: 1b

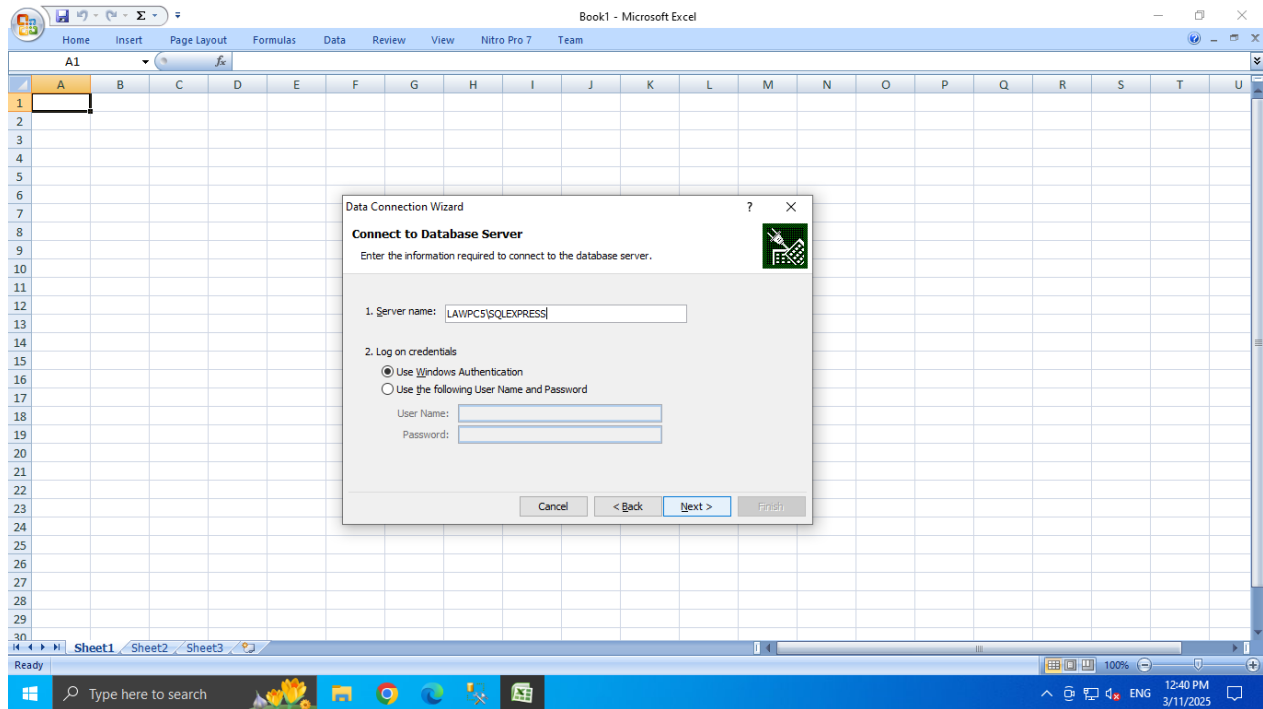
Aim: Import the cube in Microsoft Excel and create the Pivot table and Pivot Chart to perform data analysis.

Step 1: Import Data from Data Warehouse to Excel

1. Open Microsoft Excel.
2. Go to the Data tab.
3. From the drop-down menu, choose **From Other Sources** and select **From Data Connection Wizard**.
4. select **Microsoft From SQL Server** or choose the relevant source for your data warehouse.
5. Connect to database server provide **connection** (Server Name).
6. Click **Next** and then **Finish**.

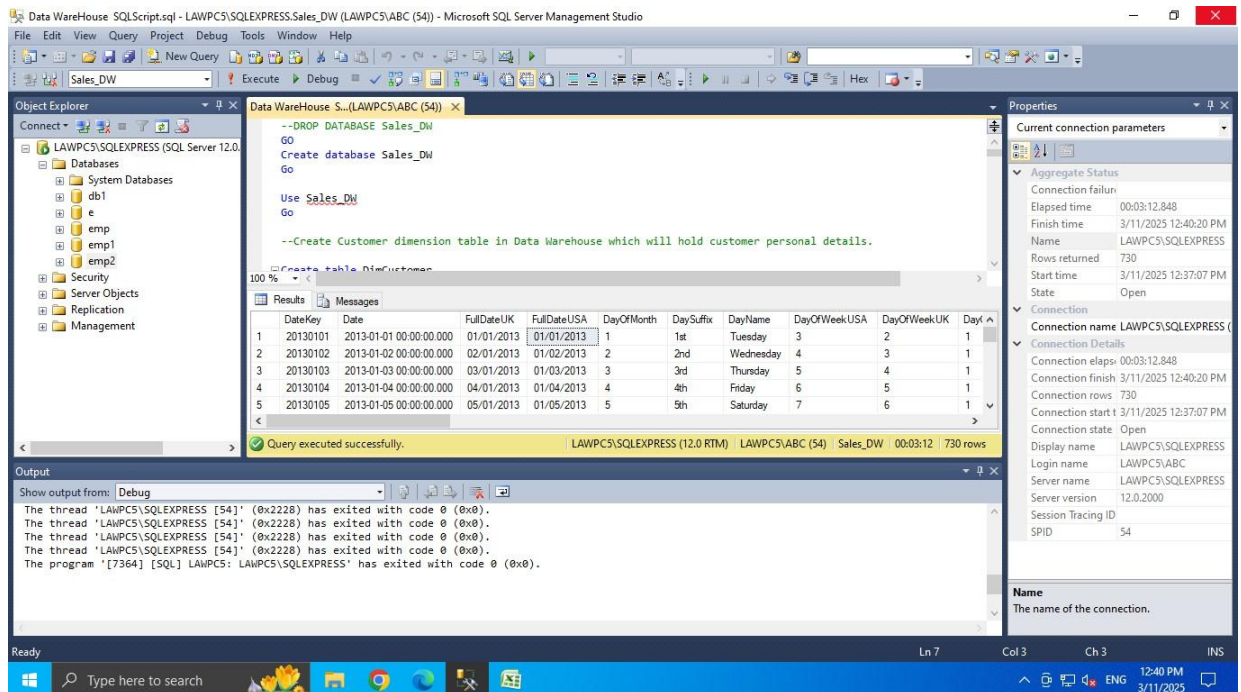
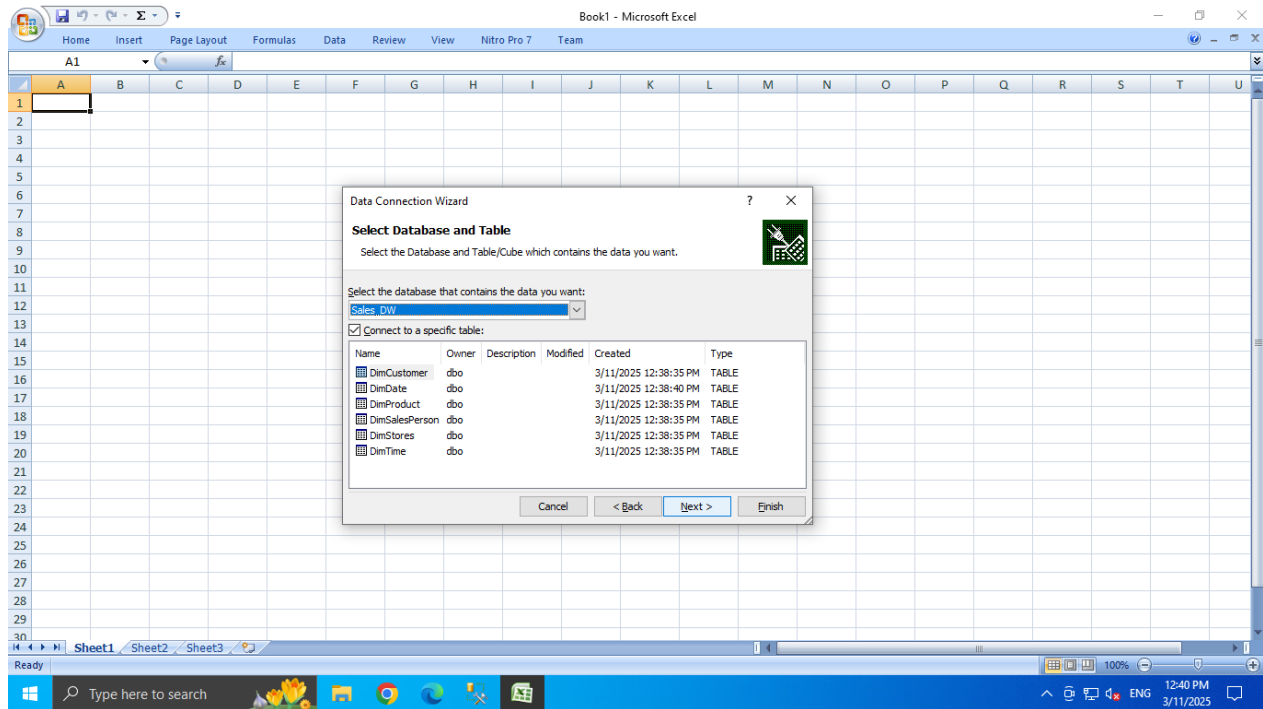


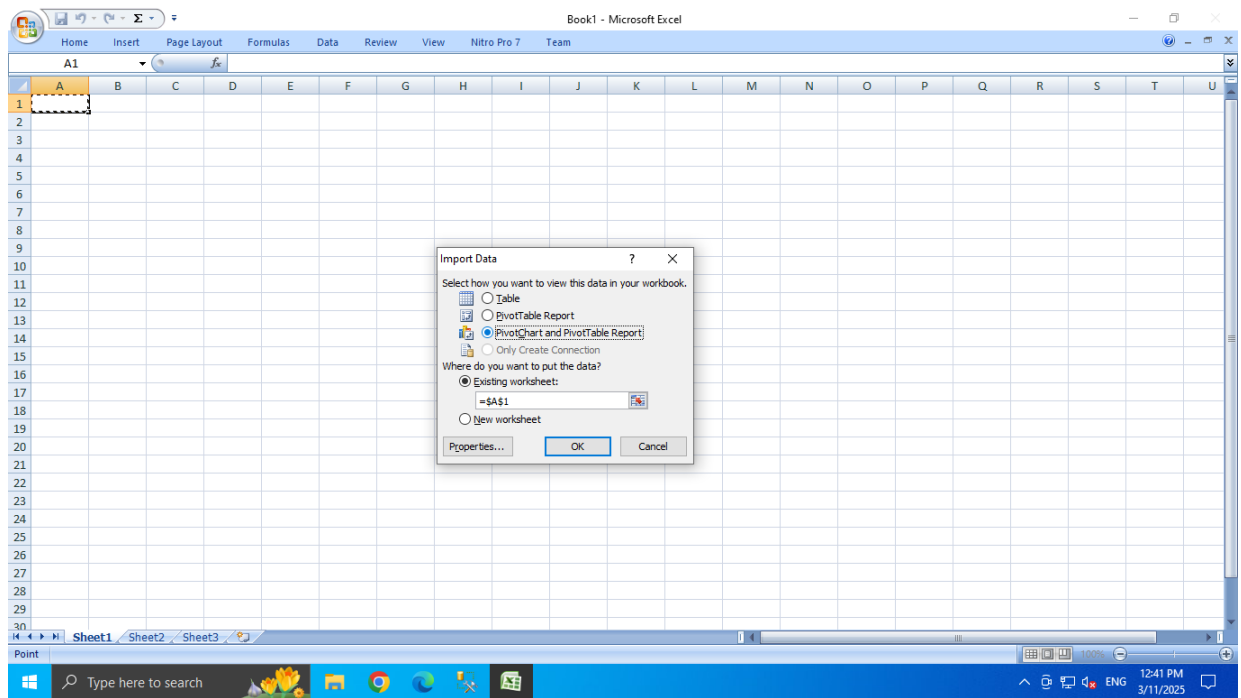
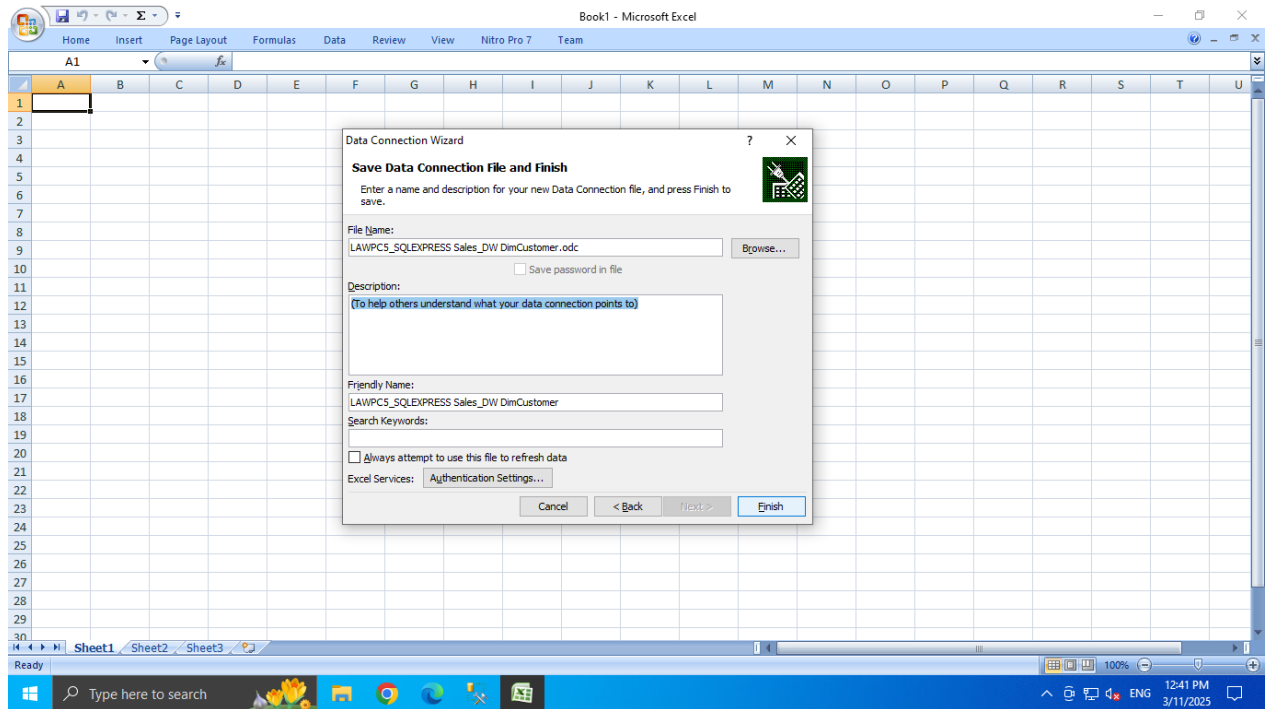


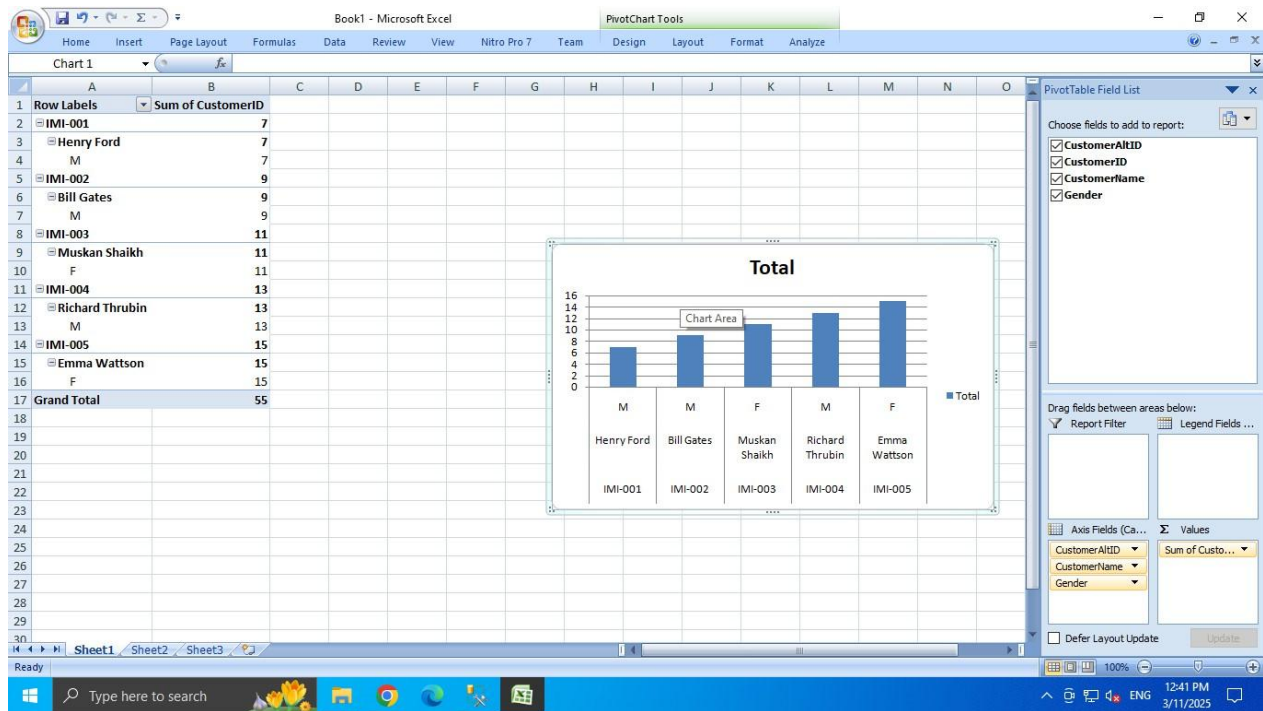


Step 2: Create a Pivot Chart

1. Once you have the PivotTable set up, click anywhere inside the PivotTable.
2. Select database and Table .
3. Save data connections file and finish .
4. Import data to select pivot chart and pivot table report.
5. Choose the chart type you want to visualize your data (e.g., Column, Line, Pie, etc.).
6. Excel will generate the PivotChart based on the data in your PivotTable.

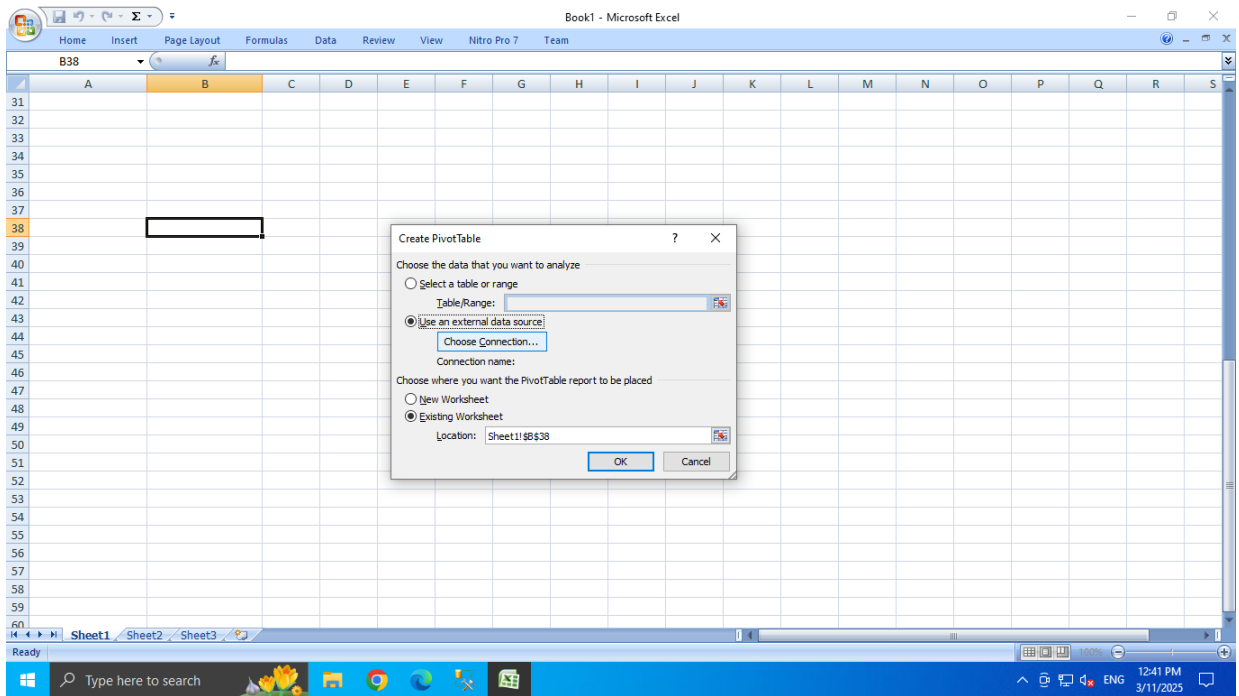
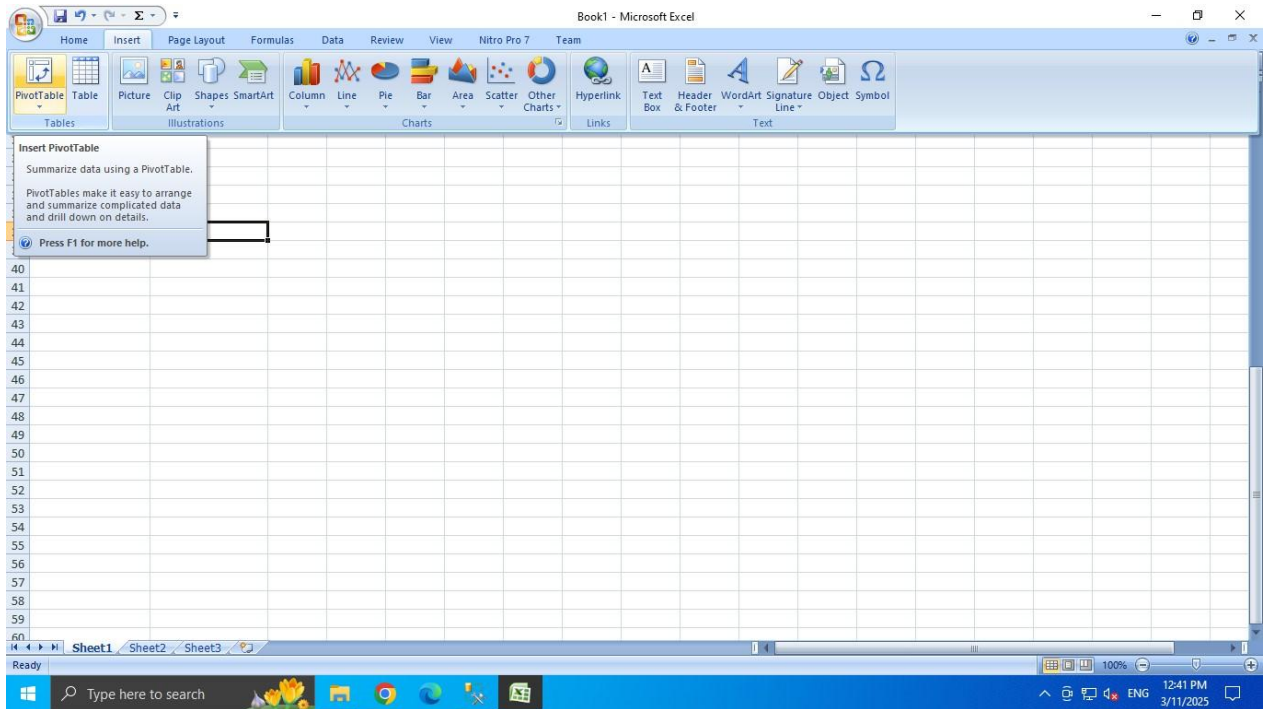


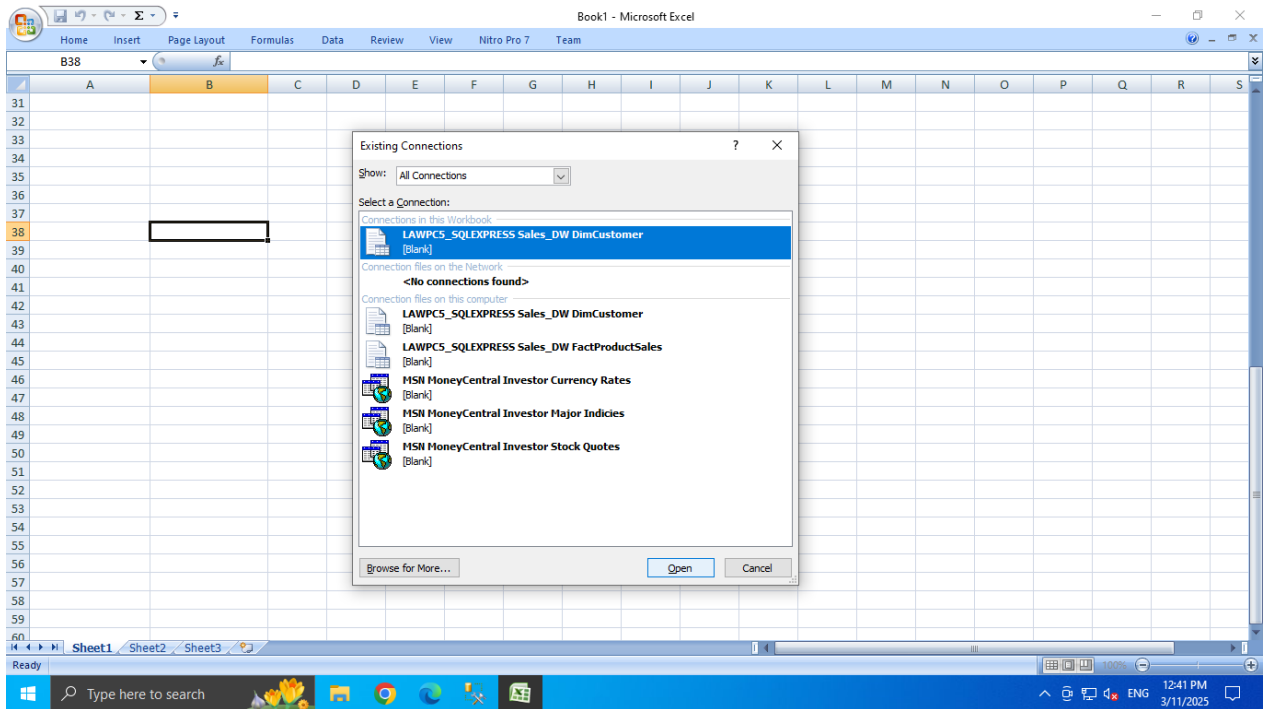




Step 3: Create a Pivot Table

1. Once the data is imported, click anywhere inside the data range.
2. Go to the **Insert tab** on the Ribbon.
3. Click on **PivotTable** in the Tables group.
4. You can also choose to place the PivotTable in a new worksheet or an existing one.
5. Click **OK** to create the PivotTable.
6. In the **PivotTable Field List** pane, drag fields into the **Rows**, **Columns**, **Values**, and **Filters** areas to organize the data. For example:
 - a. **Rows:** Place dimensions like product names or regions.
 - b. **Columns:** Place time periods (e.g., months or years).
 - c. **Values:** Put numerical values like sales or revenue.





Book1 - Microsoft Excel

PivotTable Tools: Options, Design

PivotTable Field List

Choose fields to add to report:

- ☒ CustomerAltID
- ☒ CustomerID
- ☒ CustomerName
- ☒ Gender

Drag fields between areas below:

Report Filter

Column Labels

Row Labels

Values

CustomerAltID

CustomerName

Gender

CustomerID

Sum of CustomerID

Count of CustomerName

Count of Gender

Defer Layout Update

Update

Row Labels	Sum of CustomerID	Count of CustomerName	Count of Gender
IMI-001	7	2	2
Henry Ford	7	2	2
M	7	2	2
1	1	1	1
6	6	1	1
IMI-002	9	2	2
Bill Gates	9	2	2
M	9	2	2
2	2	1	1
7	7	1	1
IMI-003	11	2	2
Muskan Shaikh	11	2	2
F	11	2	2
3	3	1	1
8	8	1	1
IMI-004	13	2	2
Richard Thruvin	13	2	2
M	13	2	2
4	4	1	1
9	9	1	1
IMI-005	15	2	2
Emma Wattson	15	2	2
F	15	2	2
5	5	1	1
10	10	1	1
Grand Total	55	10	10

Practical: 2

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

Step 1: Apply What-If Analysis in Excel

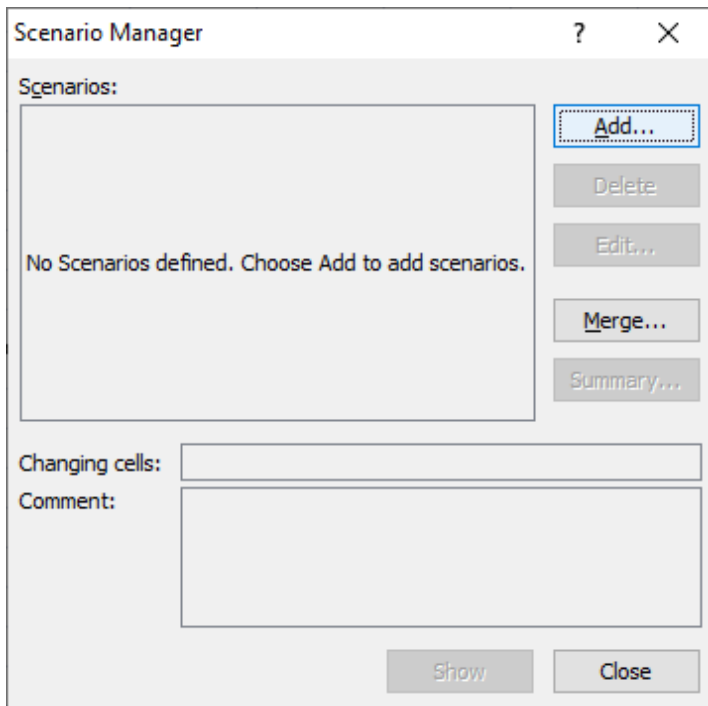
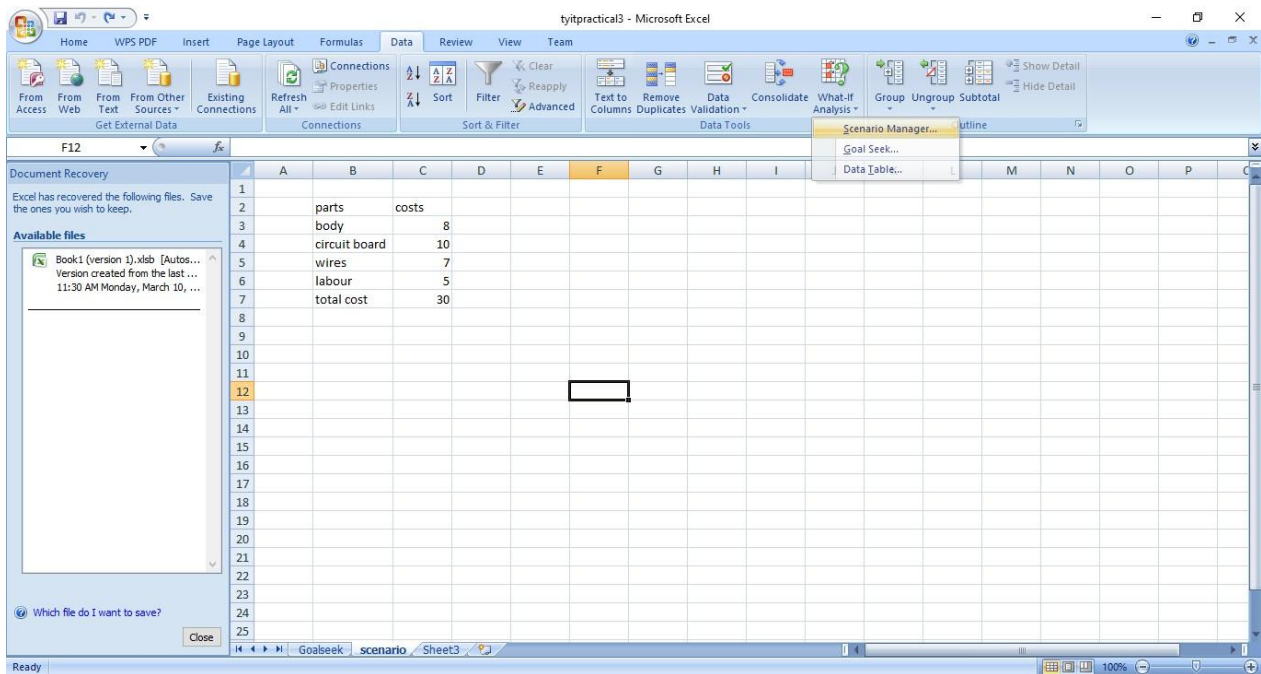
Microsoft Excel offers several tools for **What-If Analysis** that help simulate different scenarios and outcomes based on changes to input data. These tools are:

1. **Scenario Manager**
2. **Goal Seek**

2.1. Scenario Manager


The Scenario Manager lets you create and compare different sets of input values to see how changes affect the output.

- **Step-by-Step:**
 1. Go to the **Data tab** on the Ribbon.
 2. Click on **What-If Analysis**, then select **Scenario Manager**.
 3. In the **Scenario Manager** dialog box, click **Add** to create a new scenario.
 4. Name the scenario (e.g., “Best Case”, “Worst Case”, “Most Likely Case”).
 5. Define the changing cells (these are the input values you want to experiment with, such as sales volume, price, etc.).
 6. Enter the values for each scenario you want to test.
 7. After setting up all scenarios, click **OK**.
 8. To view the results of each scenario, click **Show** in the Scenario Manager dialog box.



Edit Scenario ? X

Scenario name:

Changing cells:
 

Ctrl+click cells to select non-adjacent changing cells.

Comment:

Protection

☒ Prevent changes
☐ Hide


Scenario Values ? X

Enter values for each of the changing cells.

1: labour

Edit Scenario ? X

Scenario name:

Changing cells:
 

Ctrl+click cells to select non-adjacent changing cells.

Comment:

Protection

☒ Prevent changes
☐ Hide

Scenario Values ? X

Enter values for each of the changing cells.

1: circuit_board 20

Add OK Cancel

Edit Scenario ? X

Scenario name:
future

Changing cells:
\$C\$5

Ctrl+click cells to select non-adjacent changing cells.

Comment:
Created by ADMIN on 3/11/2025

Protection
☒ Prevent changes
☐ Hide

OK Cancel

Scenario Values ? X

Enter values for each of the changing cells.

1: wires 11

Add OK Cancel

Scenario Manager

Scenarios:

present
past
future

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

Changing cells: wires

Comment: Created by ADMIN on 3/11/2025

Show Close

Scenario Summary

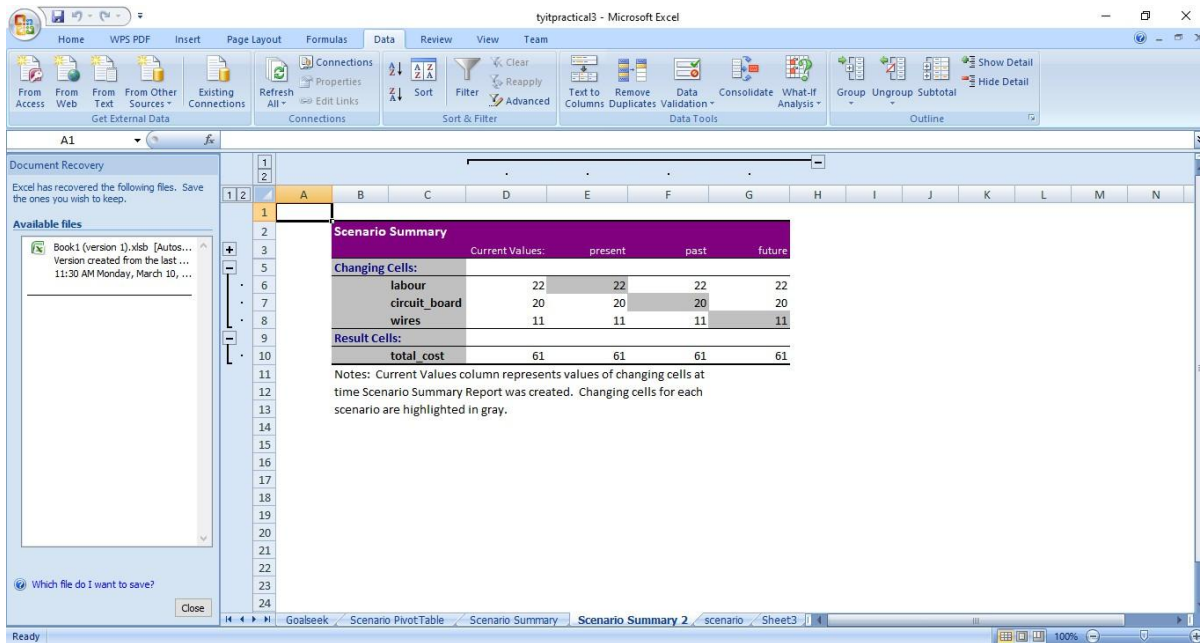
Report type

☒ Scenario summary
☐ Scenario PivotTable report

Result cells:

=\$C\$7

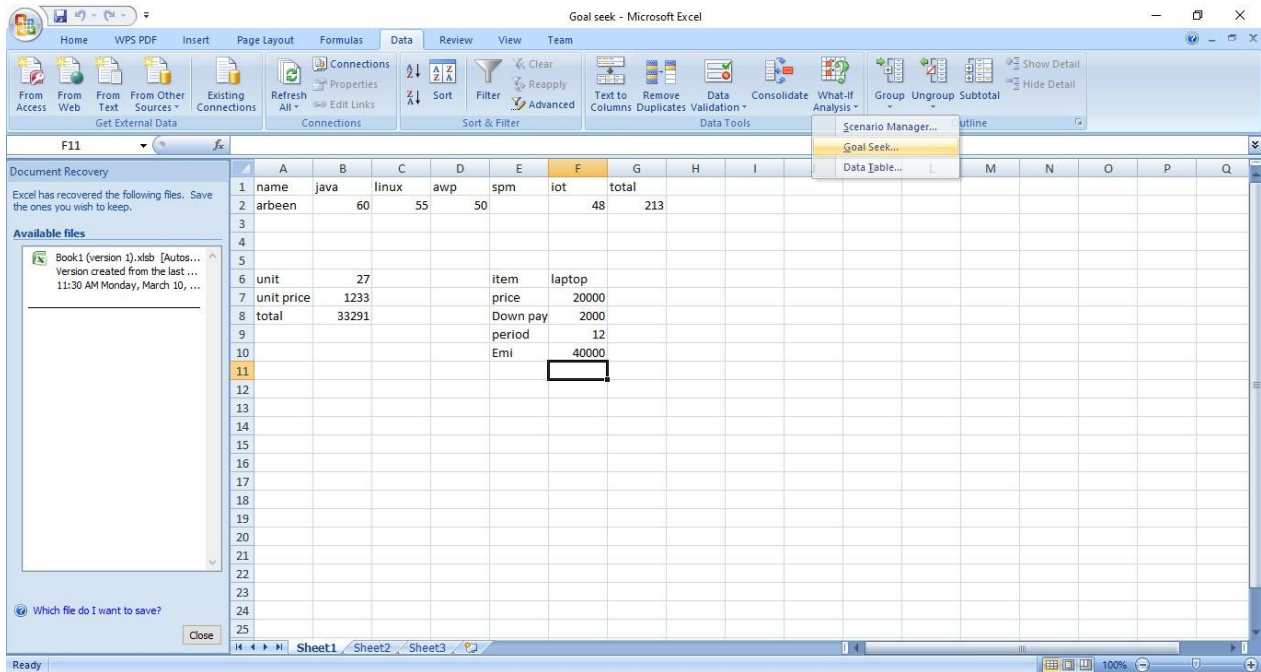
OK Cancel



2.2. Goal Seek

Goal Seek helps you find the input value needed to reach a specific result.

- **Step-by-Step:**
 1. Go to the **Data** tab on the Ribbon.
 2. Click on **What-If Analysis**, then select **Goal Seek**.
 3. In the **Goal Seek** dialog box, specify the following:
 - **Set cell:** The cell that contains the formula you want to solve.
 - **To value:** The result you want to achieve.
 - **By changing cell:** The input value that will be adjusted to reach the desired result.
 4. Click **OK**. Excel will calculate the value needed to achieve the target.



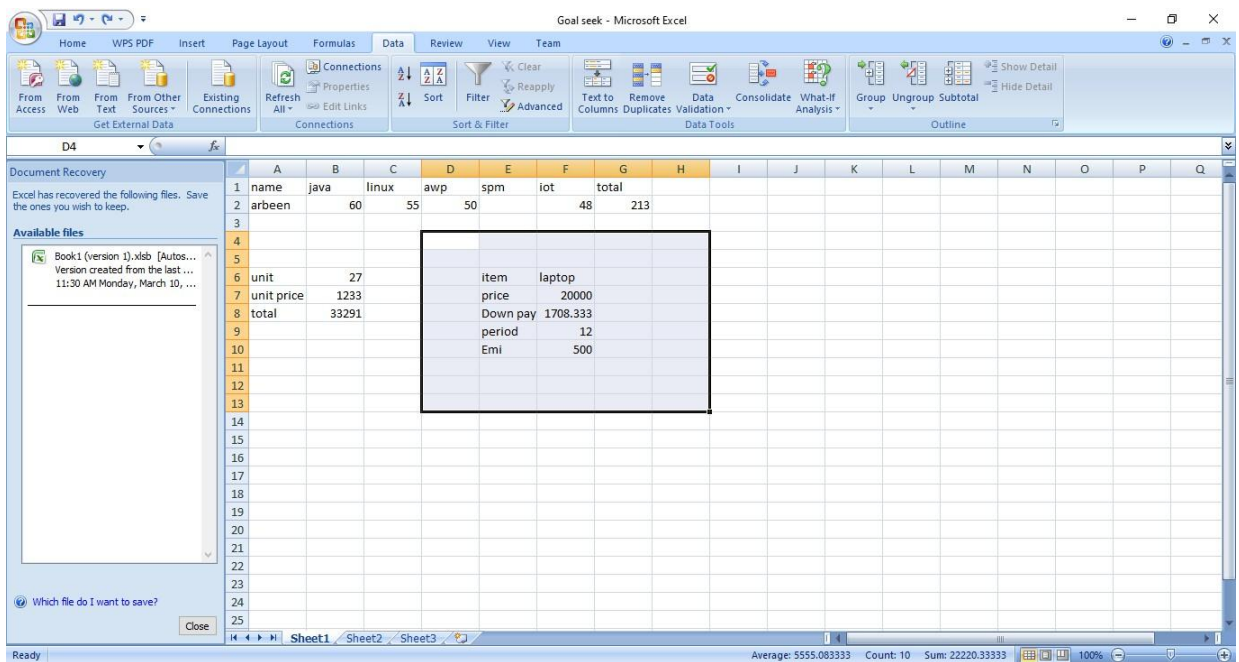
Goal Seek

Set cell:

To value:

By changing cell:

OK Cancel



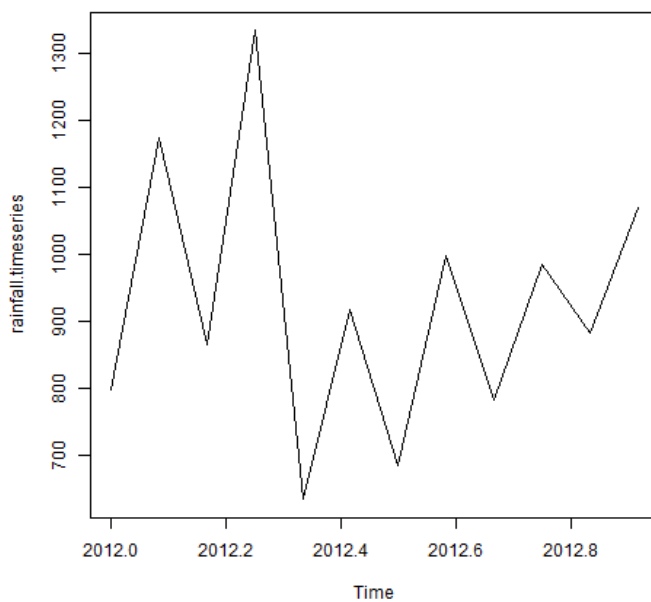
Practical: 3

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

Code:

```
rainfall <-c(799,1174.8,865.1,1334.6,635.4,918.5,685.5,998.6,784.2,985,882.8,1071)
rainfall.timeseries <- ts(rainfall,start = c(2012,1),frequency = 12)
print(rainfall.timeseries)
png(file = "rainfall.png")
plot(rainfall.timeseries)
dev.off()
```

Output:



Practical: 4

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

Code:

```
library(party)

print(head(readingSkills))

input.dat <- readingSkills[c(1:105),]

png(file = "decision_tree.png")

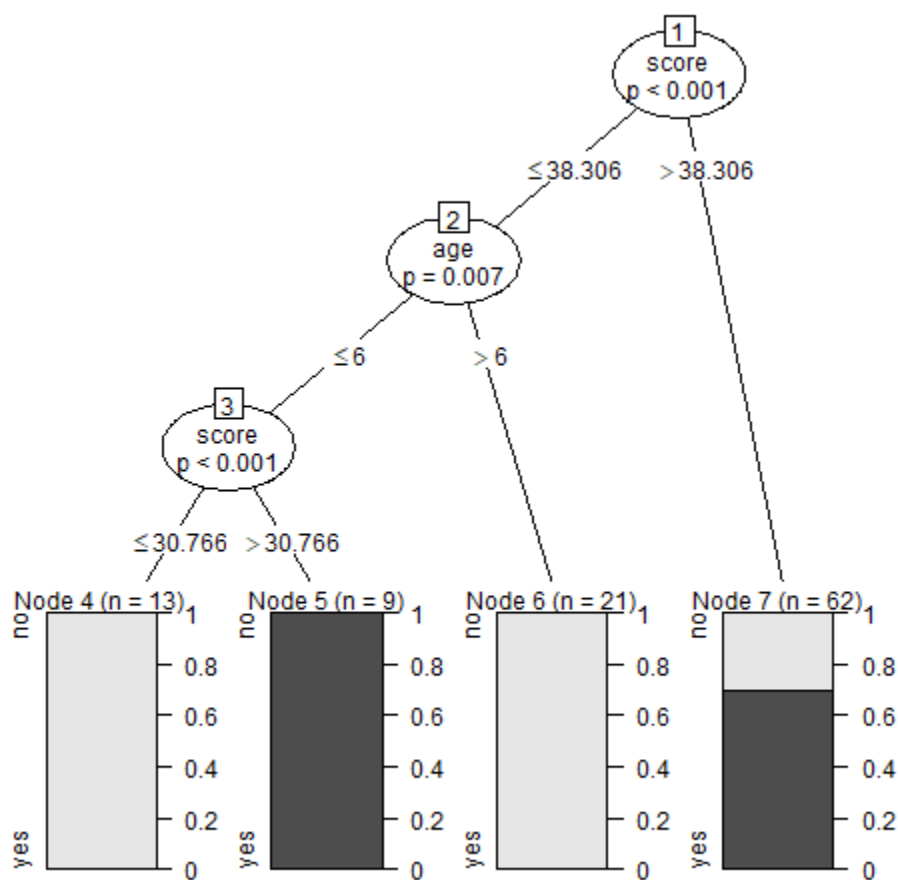
output.tree <- ctree(nativeSpeaker ~ age + shoeSize + score, data = input.dat)

plot(output.tree)

dev.off()
```

Output:

```
nativeSpeaker age shoeSize  score
1          yes   5 24.83189 32.29385
2          yes   6 25.95238 36.63105
3          no  11 30.42170 49.60593
4          yes   7 28.66450 40.28456
5          yes  11 31.88207 55.46085
6          yes  10 30.07843 52.83124
> |
```



Practical: 5

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

Code:

```
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 to create a linear regression model
relation <- lm(y ~ x)

# Print the model summary
print(summary(relation))

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

# Visualizing the Regression Graphically
png(file = "linearregression.png")
plot(x, y, col = "blue", main = "Height & Weight Regression",
      xlab = "Height in cm", ylab = "Weight in Kg", pch = 16)
abline(relation, col = "red")
dev.off()
```

Output:

```
> source("C:/Users/Lenovo/OneDrive/Desktop/SEM 6 practicals/BI PRACTICALS/p7.R")
```

```
Call:
lm(formula = y ~ x)
```

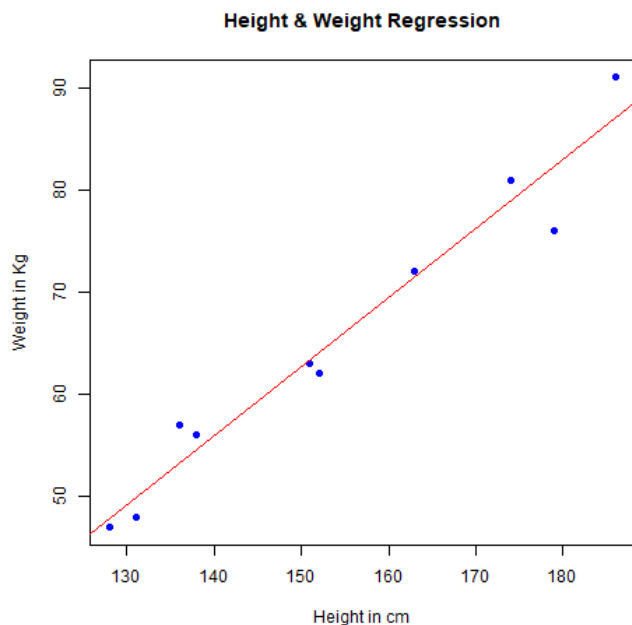
```
Residuals:
    Min       1Q   Median       3Q      Max
-6.3002 -1.6629  0.0412  1.8944  3.9775
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -38.45509     8.04901  -4.778  0.00139 **
x             0.67461     0.05191  12.997 1.16e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 3.253 on 8 degrees of freedom
Multiple R-squared:  0.9548,    Adjusted R-squared:  0.9491
F-statistic: 168.9 on 1 and 8 DF,  p-value: 1.164e-06
```

```
1
76.22869
```

```
> |
```



Practical :6

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

Code:

```
# Load necessary libraries
library(dplyr)
library(titanic)
library(pROC)

# Load Titanic dataset
data("titanic_train")

# Check dataset structure
head(titanic_train)

# Data Cleaning: Removing rows with missing values
titanic_clean <- titanic_train %>%
  filter(!is.na(Age), !is.na(Embarked), !is.na(Sex), !is.na(Pclass))

# Convert categorical variables to factors
titanic_clean$Survived <- as.factor(titanic_clean$Survived)
titanic_clean$Pclass <- as.factor(titanic_clean$Pclass)
titanic_clean$Sex <- as.factor(titanic_clean$Sex)
titanic_clean$Embarked <- as.factor(titanic_clean$Embarked)

# Build Logistic Regression Model
model <- glm(Survived ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked,
  data = titanic_clean, family = binomial)

# Print model summary
summary(model)
```

```

# Predict probabilities
predictions <- predict(model, type = "response")

# Convert probabilities to binary classification
predictions_class <- ifelse(predictions > 0.5, 1, 0)

# Evaluate Model Accuracy
confusion_matrix <- table(Predicted = predictions_class, Actual =
titanic_clean$Survived)
print(confusion_matrix)

accuracy <- mean(predictions_class == as.numeric(titanic_clean$Survived) - 1)
print(paste("Accuracy:", accuracy))

# ROC Curve
roc_curve <- roc(as.numeric(titanic_clean$Survived) - 1, predictions)
plot(roc_curve, main = "ROC Curve")

```

Output:

```

> # Load necessary libraries
> library(dplyr)
> library(titanic)
> library(pROC)
>
> # Load Titanic dataset
> data("titanic_train")
>
> # Check dataset structure
> head(titanic_train)

```

PassengerId	Survived	Pclass	Name	Sex	Age
1	1	0	3	Braund, Mr. Owen Harris	male 22
2	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female 38
3	3	1	3	Heikkinen, Miss. Laina	female 26
4	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female 35
5	5	0	3	Allen, Mr. William Henry	male 35
6	6	0	3	Moran, Mr. James	male NA

```

SibSp Parch Ticket Fare Cabin Embarked
1 1 0 A/5 21171 7.2500 S
2 1 0 PC 17599 71.2833 C85 C
3 0 0 STON/O2. 3101282 7.9250 S

```

```

4  1  0      113803 53.1000 C123   S
5  0  0      373450 8.0500      S
6  0  0      330877 8.4583      Q
>
> # Data Cleaning: Removing rows with missing values
> titanic_clean <- titanic_train %>%
+   filter(!is.na(Age), !is.na(Embarked), !is.na(Sex), !is.na(Pclass))
>
> # Convert categorical variables to factors
> titanic_clean$Survived <- as.factor(titanic_clean$Survived)
> titanic_clean$Pclass <- as.factor(titanic_clean$Pclass)
> titanic_clean$Sex <- as.factor(titanic_clean$Sex)
> titanic_clean$Embarked <- as.factor(titanic_clean$Embarked)
>
> # Build Logistic Regression Model
> model <- glm(Survived ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked,
+             data = titanic_clean, family = binomial)
>
> # Print model summary
> summary(model)

```

Call:

```
glm(formula = Survived ~ Pclass + Sex + Age + SibSp + Parch +
    Fare + Embarked, family = binomial, data = titanic_clean)
```

Coefficients:

```

              Estimate Std. Error z value Pr(>|z|)
(Intercept) 16.691979 607.920015  0.027 0.978095
Pclass2     -1.189637  0.329197 -3.614 0.000302 ***
Pclass3     -2.395220  0.343356 -6.976 3.04e-12 ***
Sexmale     -2.637859  0.223006 -11.829 < 2e-16 ***
Age         -0.043308  0.008322 -5.204 1.95e-07 ***
SibSp       -0.362925  0.129290 -2.807 0.005000 **
Parch       -0.060365  0.123944 -0.487 0.626233
Fare         0.001451  0.002595  0.559 0.576143
EmbarkedC   -12.259048 607.919885 -0.020 0.983911
EmbarkedQ   -13.082427 607.920088 -0.022 0.982831
EmbarkedS   -12.661895 607.919868 -0.021 0.983383
---

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

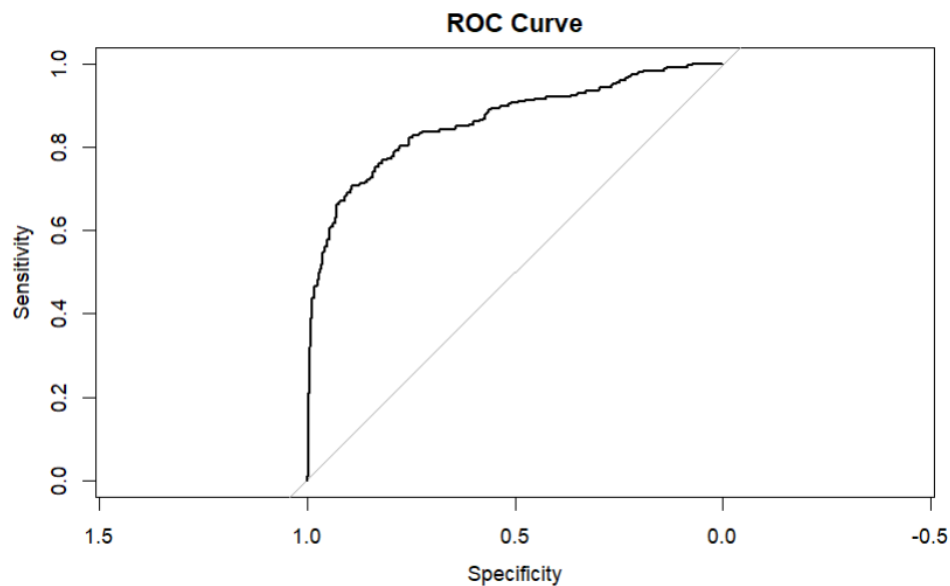
```

Null deviance: 964.52 on 713 degrees of freedom
Residual deviance: 632.34 on 703 degrees of freedom
AIC: 654.34

```

Number of Fisher Scoring iterations: 13

```
>
> # Predict probabilities
> predictions <- predict(model, type = "response")
>
> # Convert probabilities to binary classification
> predictions_class <- ifelse(predictions > 0.5, 1, 0)
>
> # Evaluate Model Accuracy
> confusion_matrix <- table(Predicted = predictions_class, Actual = titanic_clean$Survived)
> print(confusion_matrix)
      Actual
Predicted 0  1
      0 365 83
      1  59 207
>
> accuracy <- mean(predictions_class == as.numeric(titanic_clean$Survived) - 1)
> print(paste("Accuracy:", accuracy))
[1] "Accuracy: 0.801120448179272"
>
> # ROC Curve
> roc_curve <- roc(as.numeric(titanic_clean$Survived) - 1, predictions)
Setting levels: control = 0, case = 1
Setting direction: controls < cases
> plot(roc_curve, main = "ROC Curve")
```



Practical: 7

Aim: 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).

Code:

```
import pandas as pd
file_path = 'data.csv'
data = pd.read_csv(file_path)
print("First 5 rows of the dataset:")
print(data.head())
print("\nDataset Information:")
print(data.info())
print("\nSummary Statistical:")
print(data.describe())
if 'Category' in data.columns:
    print("\nUnique values in 'Category' column:")
    print(data['Category'].value_counts())
```

	A	B	C
1	Category	Sales	Profit
2	Electronics	300	40
3	Furniture	400	50
4	Clothing	350	60
5	Electronics	250	80
6	Clothing	200	70
7	Furniture	500	100

Output:

```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

warnings.warn(
PS C:\Users\Rashmi Pandey\Desktop\Logistic> & "C:/Program Files/Python313/python.exe" "c:/Users/Rashmi Pandey/Desktop/Logistic/data.py"
First 5 rows of the dataset:
   Category  Sales  Profit
0  Electronics    300     40
1  Furniture     400     50
2   Clothing     350     60
3  Electronics    250     80
4   Clothing     200     70

Dataset Information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  ---
0   Category  6 non-null      object
1   Sales      6 non-null      int64
2   Profit     6 non-null      int64
dtypes: int64(2), object(1)
memory usage: 276.0+ bytes
None

Summary Statistical:
      Sales      Profit
count  6.000000  6.000000
mean   333.333333  66.666667
std    108.012345  21.602469
min     200.000000  40.000000
25%    262.500000  52.500000
50%    325.000000  65.000000
75%    387.500000  77.500000
max     500.000000 100.000000

Unique values in 'Category' column:
Category
Electronics    2
Furniture      2
Clothing       2
Name: count, dtype: int64
```


Practical: 8a

Aim: Perform data visualization using python on any sales data.

Code:

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
data = {
    'Data': pd.date_range(start='2023-01-01', periods=100, freq='D'),
    'Product': ['Product A', 'Product B', 'Product C', 'Product D']* 25,
    'Category': ['Electronics', 'Furniture', 'Clothing', 'Books']* 25,
    'Region': ['North', 'South', 'East', 'West']* 25,
    'Sales': [x * 10 for x in range(100)],
    'Profit': [x * 2 for x in range(100)],
}
df = pd.DataFrame(data)
print(df.head())
plt.figure(figsize=(10, 5))
sns.lineplot(data=df, x='Data', y='Sales', marker='o')
plt.title('Sales Trend Over Time')
plt.xlabel('Data')
plt.ylabel('Sales')
plt.grid(True)
plt.show()
plt.figure(figsize=(8, 5))
sns.barplot(data=df, x='Category', y='Sales', ci=None, palette='viridis')
plt.title('Sales by Category')
plt.xlabel('Category')
plt.ylabel('Total Sales')
plt.show()
plt.figure(figsize=(8, 5))
sns.boxplot(data=df, x='Region', y='Profit', palette='coolwarm')
plt.title('Profit Distribution by Region')
plt.xlabel('Region')
plt.ylabel('Profit')
plt.show()
plt.figure(figsize=(8, 5))
product_sales = df.groupby('Product')['Sales'].sum().reset_index()
sns.barplot(data=product_sales, x='Product', y='Sales', palette='magma')
plt.title('Total Sales by Product')
plt.xlabel('Product')
```

```
plt.ylabel('Sales')
plt.show()
```

Output:

	Data	Product	Category	Region	Sales	Profit
0	2023-01-01	Product A	Electronics	North	0	0
1	2023-01-02	Product B	Furniture	South	10	2
2	2023-01-03	Product C	Clothing	East	20	4
3	2023-01-04	Product D	Books	West	30	6
4	2023-01-05	Product A	Electronics	North	40	8

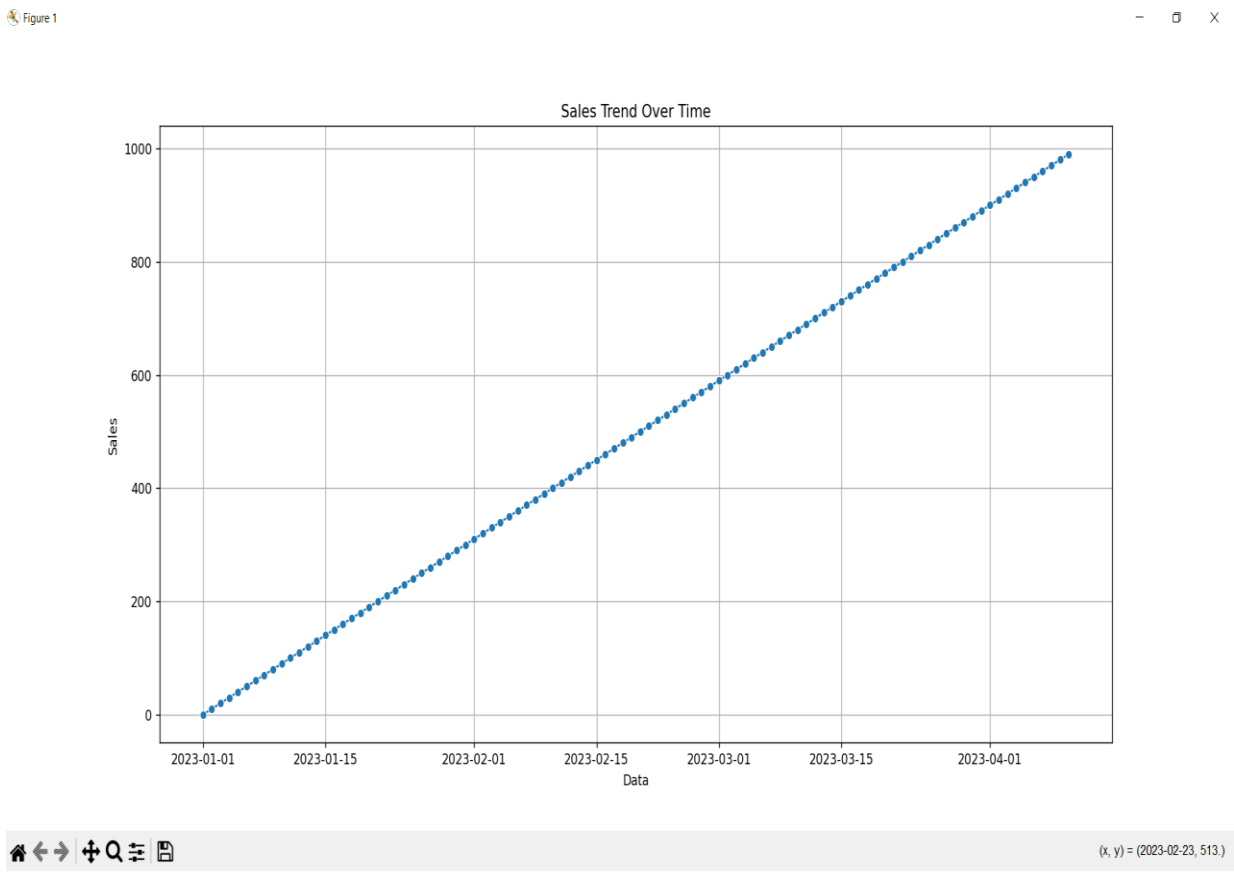
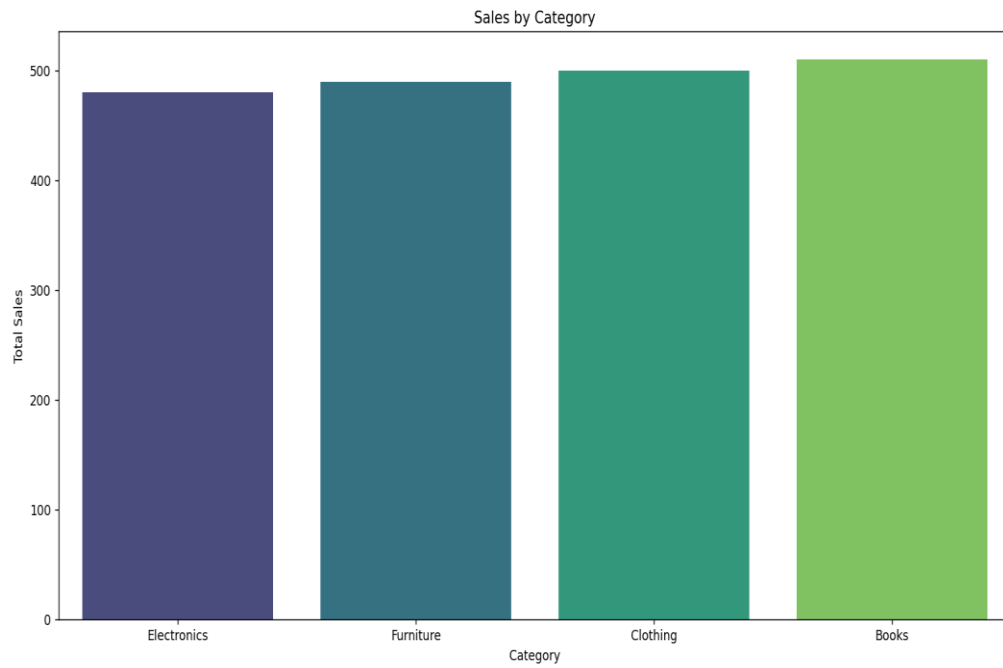
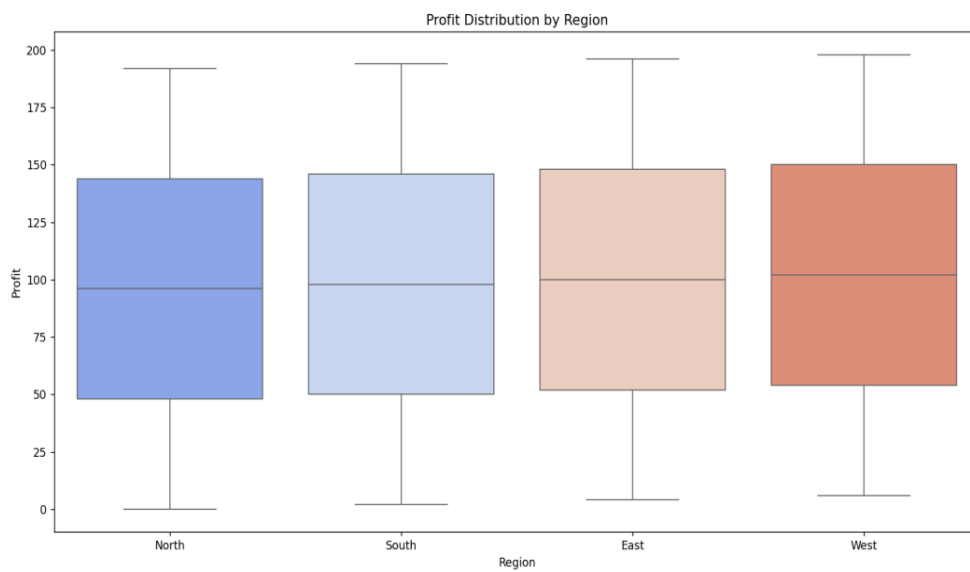


Figure 1



(x, y) = (Furniture, 445.3)

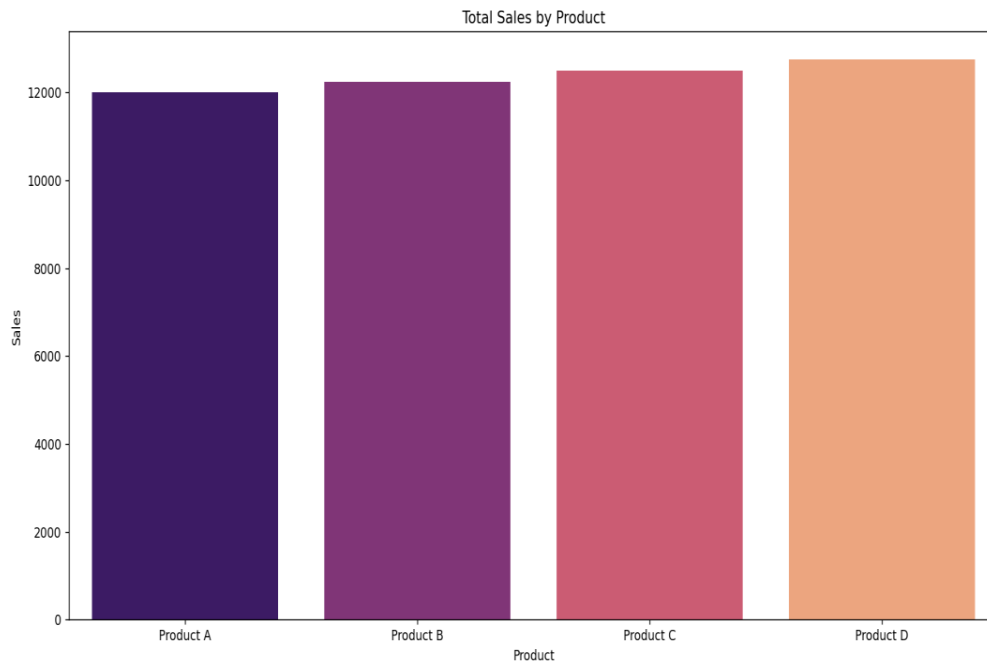
Figure 1



(x, y) = (East, 28.4)

Figure 1

— □ ×



🏠 ⬅ ➡ 🔍 📄

(x, y) = (Product A, 7.90e+03)

Practical: 8b

Aim: Perform data visualization using PowerBI on any sales data.

Steps:-

Step 1:- Download & Install Power BI

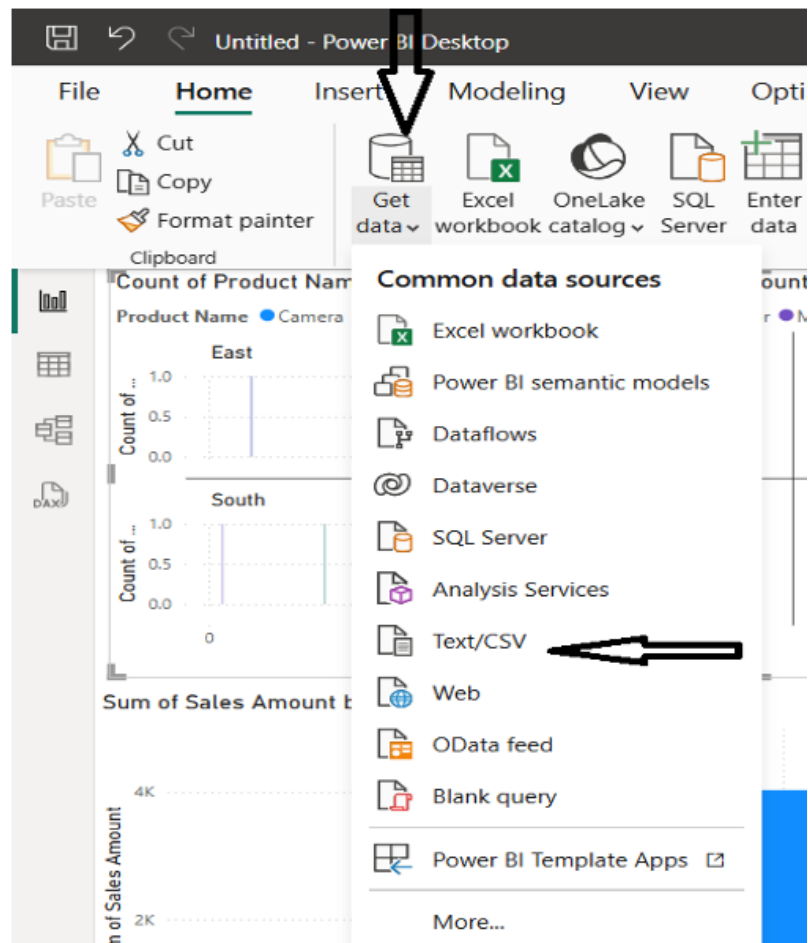
Step 2:- create a sample data

	A	B	C	D	E	F	G	H	I	J	K
1	Order ID	Date	Product Name	Category	Sales Amount	Quantity Sold	Customer Name	Customer Region	Payment Method	Discount (%)	Profit
2	ORD1001	01-01-2024	Laptop	Electronics	1200	2	John Doe	New York	Credit Card	10	300
3	ORD1002	02-01-2024	Smartphone	Electronics	800	1	Alice Smith	California	Debit Card	5	150
4	ORD1003	03-01-2024	Shoes	Fashion	100	3	Robert Brown	Texas	Cash	15	30
5	ORD1004	04-01-2024	Washing Machine	Home Appliances	500	1	Emily Davis	Florida	Net Banking	8	80
6	ORD1005	05-01-2024	Tablet	Electronics	300	2	Michael Lee	Illinois	UPI	10	60
7	ORD1006	06-01-2024	Headphones	Accessories	50	4	Sophia Wilson	New York	Credit Card	5	15
8	ORD1007	07-01-2024	Refrigerator	Home Appliances	900	1	Daniel Scott	Texas	Debit Card	12	180

Step 3 :- Import the Sales Dataset

- Click "Home" > "Get Data".
- Choose your data source:
 - o Excel (XLSX/CSV)
 - o SQL Server
 - o Online Services (Google Sheets, SharePoint, etc.)
- Browse and select the dataset.

- Click "Load" to import.



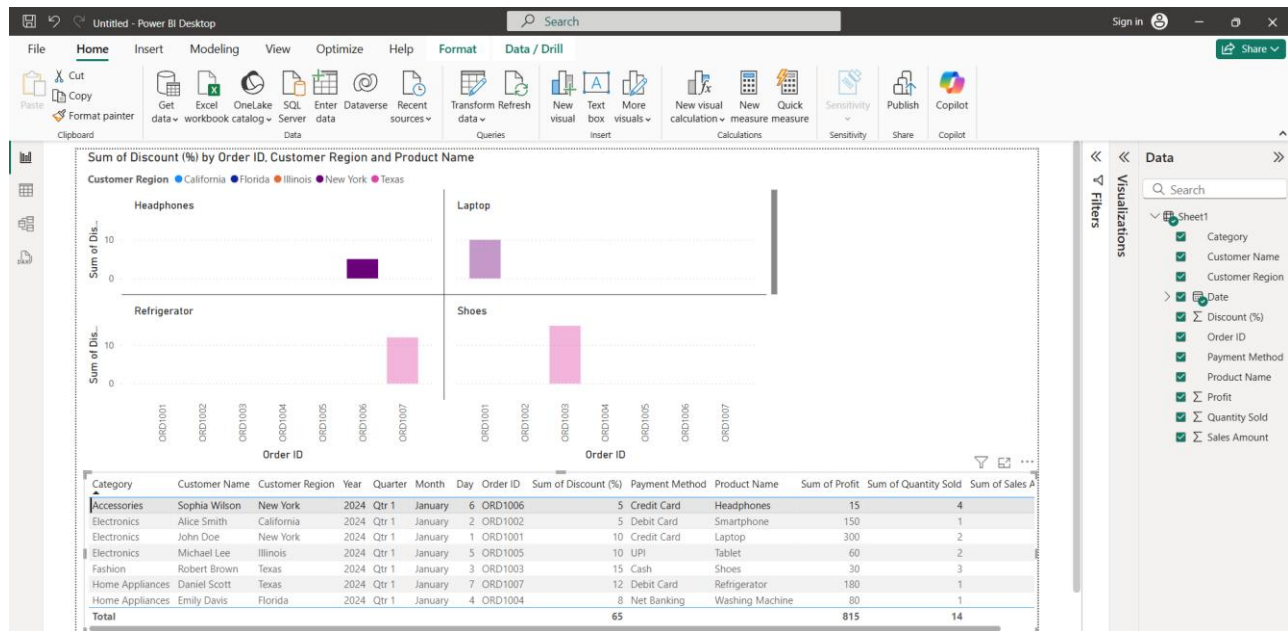
Step 4 :- Clean & Transform Data (Power Query)

Step 5 :- Create Data Visualizations

- Click on "Line Chart" in the "Visualizations" pane.
- Drag Order Date to the X-axis.
- Drag Sales Amount to the Y-axis.
- Customize the chart (format labels, add title, etc.).

Step 6:- Add Filters & Interactivity

Output:-



Practical: 9

Aim: Create the Data staging area for the selected database using SQL

STEPS:-

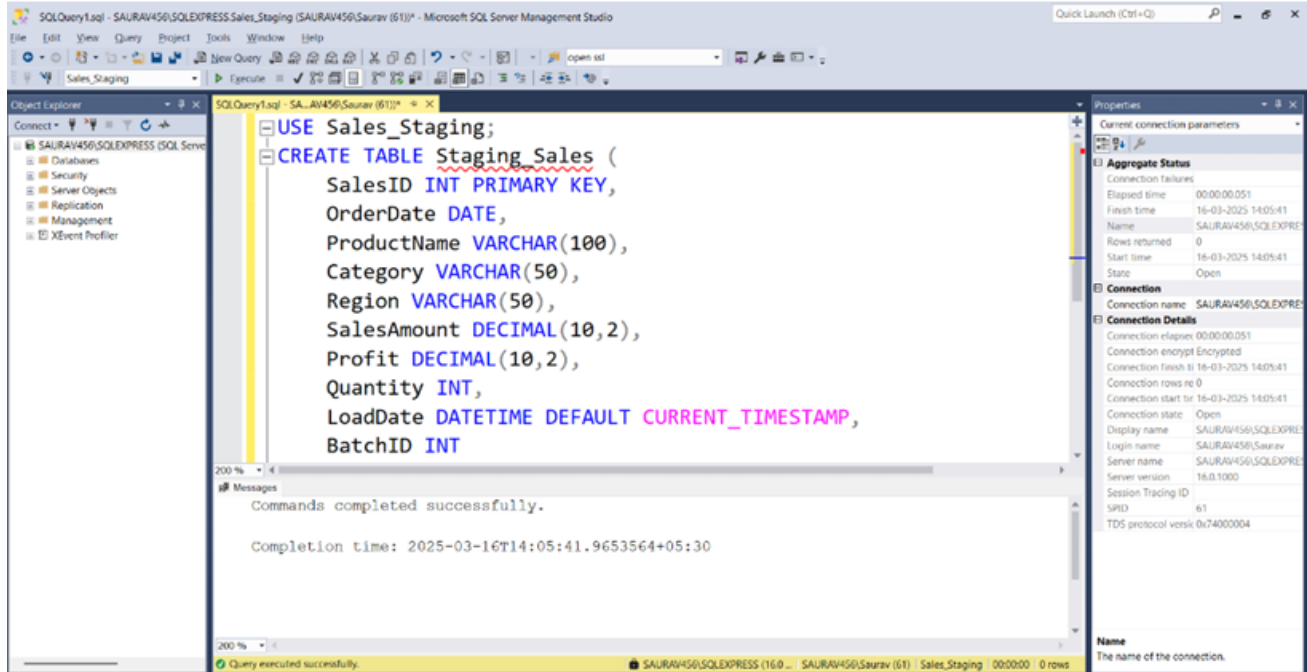
Step 1:- Create a Staging Database First, create a staging database to store raw sales data.

```
CREATE DATABASE Sales_Staging;
```

```
USE Sales_Staging;
```

Step 2:- Create Staging Tables Create tables that match the structure of raw sales data but include additional fields like load date and batch ID.

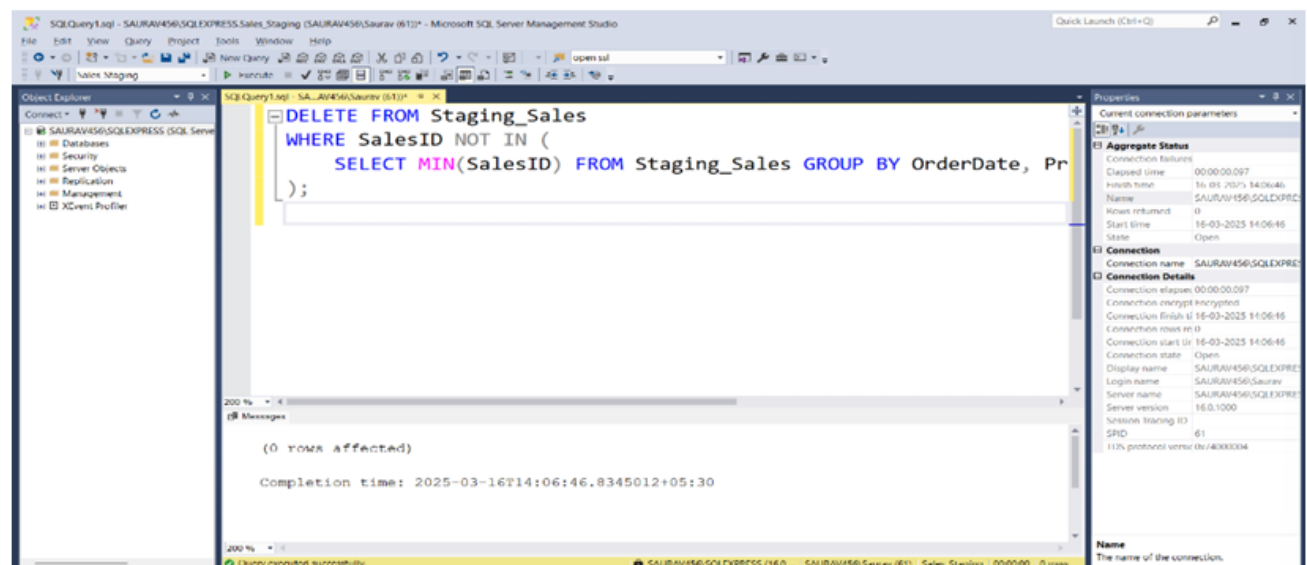
```
CREATE TABLE Staging_Sales (  
SalesID INT PRIMARY KEY,  
OrderDate DATE,  
ProductName VARCHAR(100),  
Category VARCHAR(50),  
Region VARCHAR(50),  
SalesAmount DECIMAL(10,2),  
Profit DECIMAL(10,2),  
Quantity INT,  
LoadDate DATETIME DEFAULT CURRENT_TIMESTAMP,  
BatchID INT );
```

Step 4:- Perform Data Cleansing & Transformation

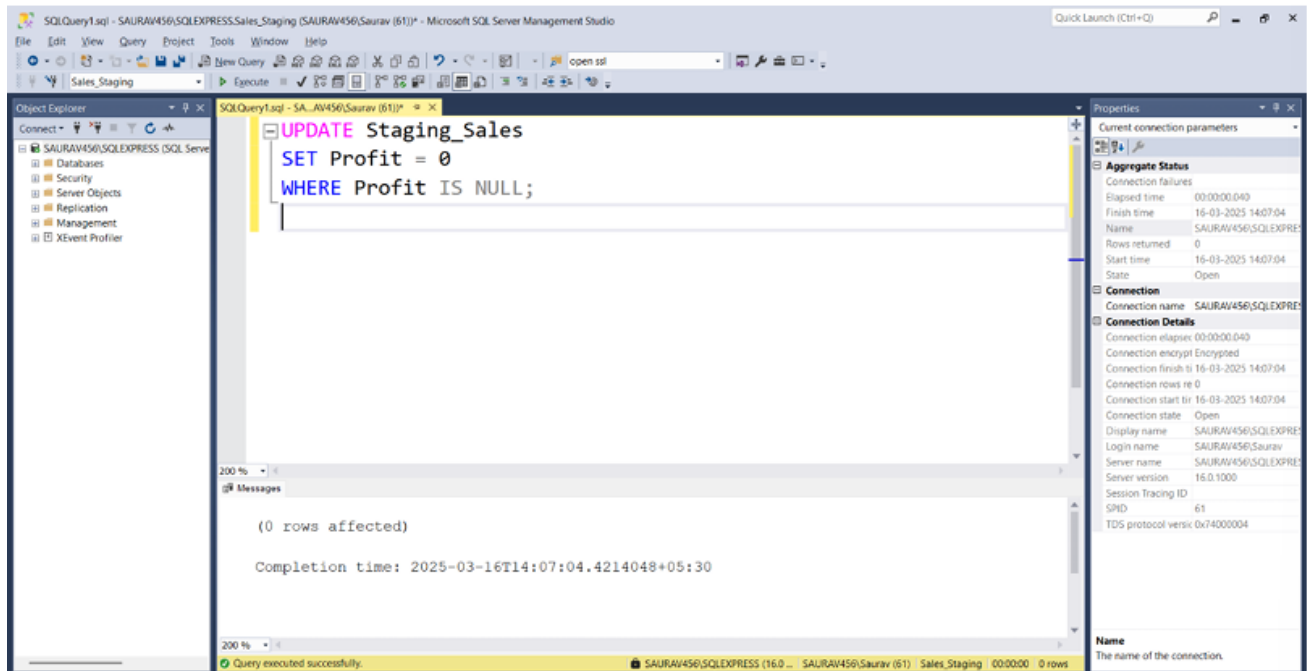
- Remove Duplicates DELETE FROM Staging_Sales

WHERE SalesID NOT IN (
SELECT MIN(SalesID) FROM Staging_Sales GROUP BY OrderDate, ProductName, Region);



- **Handle Null Values**

```
UPDATE Staging_Sales
SET Profit = 0
WHERE Profit IS NULL;
```



Output:

Results		Messages								
	SalesID	OrderDate	ProductName	Category	Region	SalesAmount	Profit	Quantity	LoadDate	BatchID
1	1	2024-01-01	Laptop	Electronics	North	1200.00	200.00	3	2025-03-21 13:49:51.297	101
2	2	2024-01-02	Smartphone	Electronics	South	800.00	150.00	2	2025-03-21 13:49:51.297	101
3	3	2024-01-03	Tablet	Electronics	East	600.00	100.00	5	2025-03-21 13:49:51.297	101