

Practical 1

Aim-: Import the legacy data from different sources such as (EXCEL,SQLSERVER,ORACLE) and load in the target system

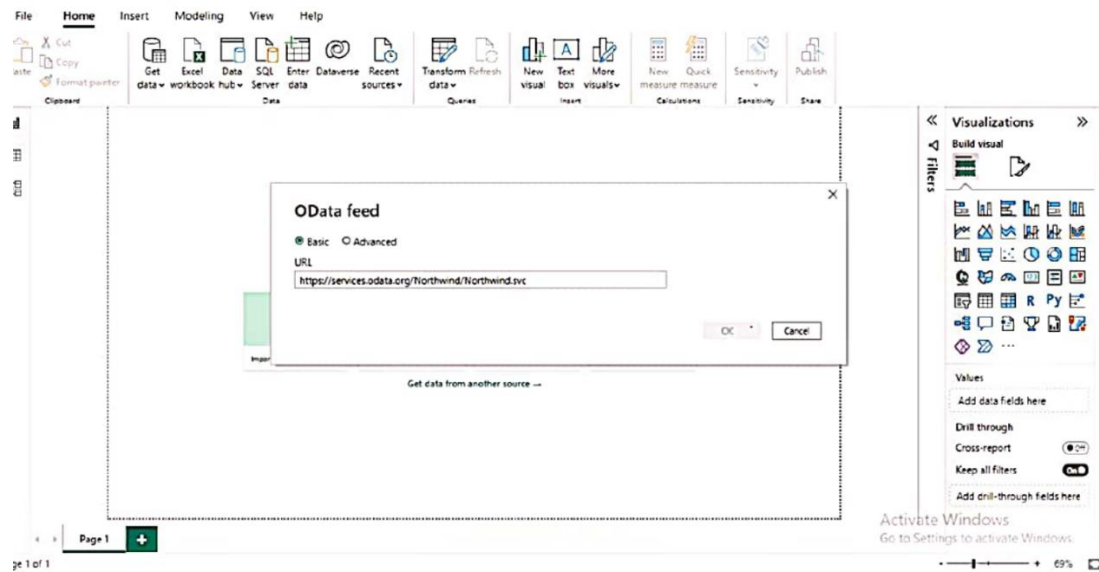
Definition of legacy data-:

- 1) The Legacy data is an electronic data that is only understood by only compatible with or only resides on hardware or software that has become obsolete.
- 2) This means that legacy data includes everything that modern technology is not able to open.

Dataset Used-: Northwind Database

- 1) The Northwind database is a simple database used to demonstrate the features of some of its products.
- 2) The database contains the sales data for Northwind Traders, a fictitious Speciality foods export-import company.

Output-:



Navigator

Navigator

Show All | Show Selected [1]

http://services.odata.org/V3/Northwind/Nort...

- Alphabetical_list_of_products
- Categories
- Category_Sales_for_1997
- Current_Product_Lists
- Customer_and_Suppliers_by_Cities
- CustomerDemographics
- Customers
- Employees
- Invoices
- Order_Details
- Order_Details_Extended
- Order_Subtotals
- Orders**
- Orders_Qries
- Product_Sales_for_1997
- Products
- Products_Above_Average_Prices
- Products_by_Categories
- Regions

Orders

OrderID	CustomerID	EmployeeID	OrderDate	RequiredDate
10248	VINET	5	7/4/1996 12:00:00 AM	8/1/1996
10249	TOMSP	6	7/5/1996 12:00:00 AM	8/16/1996
10250	HANAR	4	7/8/1996 12:00:00 AM	8/5/1996
10251	VICTE	3	7/8/1996 12:00:00 AM	8/6/1996
10252	SUPED	4	7/9/1996 12:00:00 AM	8/6/1996
10253	HANAR	3	7/10/1996 12:00:00 AM	7/24/1996
10254	CHOPS	5	7/11/1996 12:00:00 AM	8/8/1996
10255	RICSU	9	7/12/1996 12:00:00 AM	8/9/1996
10256	WELLI	3	7/13/1996 12:00:00 AM	8/12/1996
10257	HILAA	4	7/16/1996 12:00:00 AM	8/13/1996
10258	ERNST	1	7/17/1996 12:00:00 AM	8/14/1996
10259	CENTC	4	7/18/1996 12:00:00 AM	8/15/1996
10260	OTTIK	4	7/19/1996 12:00:00 AM	8/16/1996
10261	QUEDE	4	7/19/1996 12:00:00 AM	8/16/1996
10262	RATTC	8	7/22/1996 12:00:00 AM	8/18/1996
10263	FRNSH	9	7/23/1996 12:00:00 AM	8/20/1996
10264	POLAO	6	7/24/1996 12:00:00 AM	8/21/1996
10265	BLONP	2	7/25/1996 12:00:00 AM	8/22/1996
10266	WARTH	3	7/26/1996 12:00:00 AM	9/6/1996
10267	FRANK	4	7/26/1996 12:00:00 AM	8/26/1996
10268	GROSR	8	7/30/1996 12:00:00 AM	8/27/1996
10269	WHITC	5	7/31/1996 12:00:00 AM	8/14/1996
10270	WARTH	1	8/1/1996 12:00:00 AM	8/29/1996

OK Cancel

Practical 2

Aim-: Perform the Extraction Transformation and Loading (ETL) Process to Construct the database in the SQL Server

Definition of Extraction Transformation and Loading-:

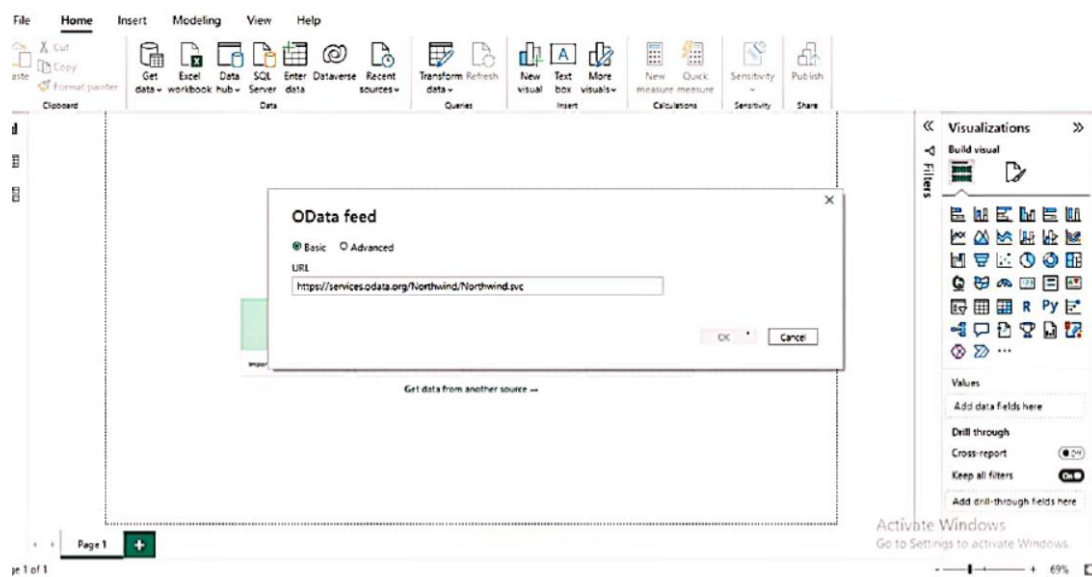
- 1) In Computing, extract, transform, load is a 3 Phases process where data is extracted, transformed & loaded into an output data container.
- 2) It is a tool performed on a given set of data in order to accomplish a specific business goal.
- 3) Power B.I. Extract Transform load (ETL) Dataflow is a cloud-based platform that users in data preparation.
- 4) Users can also utilize Power B.I. Dataflows to ingest, convert, and load data into Microsoft Dataverse environments, Power B.I. Workspace.
- 5) Without ETL Tools, the analytical reports and dashboards look old because of outdated data.
- 6) ETL helps update them so your reports are current.
- 7) ETL includes the following processes :-
The data is extracted from a data source, then transformed, validated standardized, corrected

loaded into a quality checked and ultimately data repository, such as data mart or data warehouse, where it is streamlined for analysis and reporting.

Output:-

Extract:-

Extract



Navigator

Show All | Show Selected [1]

http://services.odata.org/V3/Northwind/Nort...

- Alphabetical_list_of_products
- Categories
- Category_Sales_for_1997
- Current_Product_Lists
- Customer_and_Suppliers_by_Cities
- CustomerDemographics
- Customers
- Employees
- Invoices
- Order_Details
- Order_Details_Extended
- Order_Subtotals
- Orders**
- Orders_Qries
- Product_Sales_for_1997
- Products
- Products_Above_Average_Prices
- Products_by_Categories
- Regions

Orders

OrderID	CustomerID	EmployeeID	OrderDate	RequiredDate
10248	VINET	5	7/4/1996 12:00:00 AM	8/2/1996
10249	TONSP	6	7/5/1996 12:00:00 AM	8/16/1996
10250	HANAR	4	7/8/1996 12:00:00 AM	8/5/1996
10251	VICTE	3	7/8/1996 12:00:00 AM	8/5/1996
10252	SUPRD	4	7/9/1996 12:00:00 AM	8/6/1996
10253	HANAR	5	7/10/1996 12:00:00 AM	7/24/1996
10254	CHOPS	5	7/11/1996 12:00:00 AM	8/8/1996
10255	RICSU	9	7/12/1996 12:00:00 AM	8/9/1996
10256	WELLI	5	7/15/1996 12:00:00 AM	8/12/1996
10257	HILAA	4	7/16/1996 12:00:00 AM	8/13/1996
10258	FRNSH	1	7/17/1996 12:00:00 AM	8/14/1996
10259	CENIC	4	7/18/1996 12:00:00 AM	8/15/1996
10260	OTTM	4	7/19/1996 12:00:00 AM	8/16/1996
10261	QUEDE	4	7/19/1996 12:00:00 AM	8/16/1996
10262	RATTC	8	7/22/1996 12:00:00 AM	8/19/1996
10263	FRNSH	9	7/23/1996 12:00:00 AM	8/20/1996
10264	POLKO	6	7/24/1996 12:00:00 AM	8/21/1996
10265	BLOAP	2	7/25/1996 12:00:00 AM	8/22/1996
10266	WARTH	3	7/26/1996 12:00:00 AM	8/6/1996
10267	FRANK	4	7/29/1996 12:00:00 AM	8/26/1996
10268	GROSR	8	7/30/1996 12:00:00 AM	8/27/1996
10269	WHITC	5	7/31/1996 12:00:00 AM	8/14/1996
10270	WARTH	1	8/1/1996 12:00:00 AM	8/29/1996

OK Cancel

Transform-:

Untitled - Power Query Editor

File Home Transform Add Column View Tools Help

Queries [2]

Products

Orders

Table.RemoveColumns(Products_Table, {"SupplierID", "CategoryID", "UnitPrice", "UnitsInOrd

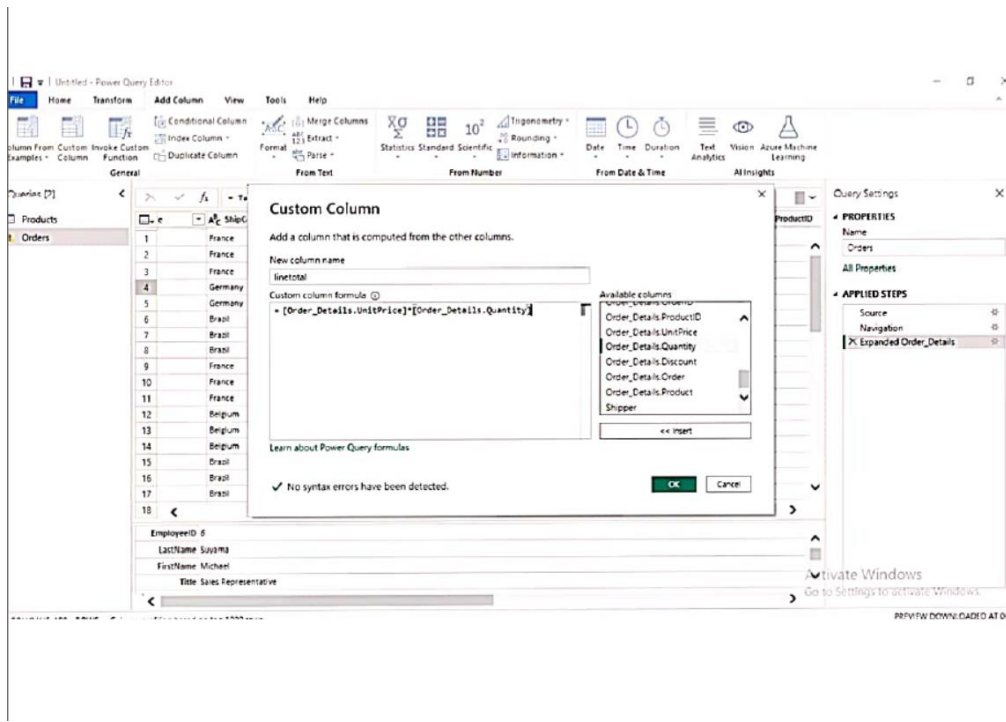
ProductID	ProductName	QuantityPerUnit	UnitPrice	UnitsInOrder
1	Chai	10 boxes x 20 bags	39	Re
2	Chang	24 - 12 oz bottles	27	Re
3	Aniseed Syrup	12 - 500 ml bottles	28	Re
4	Chef Anton's Cajun Seasoning	48 - 6 oz jars	53	Re
5	Chef Anton's Gumbo Mix	36 boxes	0	Re
6	Grandma's Boysenberry Spread	12 - 8 oz jars	120	Re
7	Uncle Bob's Organic Dried Pears	12 - 1 lb pags	15	Re
8	Northwoods Cranberry Sauce	12 - 12 oz jars	6	Record
9	Mishi Kobe Niku	18 - 500 g pags	29	Record
10	Kuon	12 - 200 ml jars	21	Record
11	Queso Cabrales	1 kg pkg	22	Record
12	Queso Manchego La Pastora	10 - 500 g pags	85	Record
13	Korhu	24 boxes	24	Record
14	Tofu	40 - 100 g pags	35	Record
15	Garden of Eatin'	24 - 250 ml bottles	39	Record
16	Pavlova	12 - 500 g boxes	29	Record
17	Alice Mutton	20 - 1 kg tins	0	Record
18	Camu Camu	18 kg pkg	42	Record
19	Treatime Chocolate Biscuits	10 boxes x 12 pieces	25	Record
20	Sir Rodney's Marmalade	30 g/tt boxes	40	Record
21	Sir Rodney's Scones	24 pkts, x 4 pieces	3	Record
22	Gustaf's Knackebrod	24 - 500 g pags	104	Record
23	Tunnbrod	12 - 250 g pags	61	Record

Properties

Applied Steps

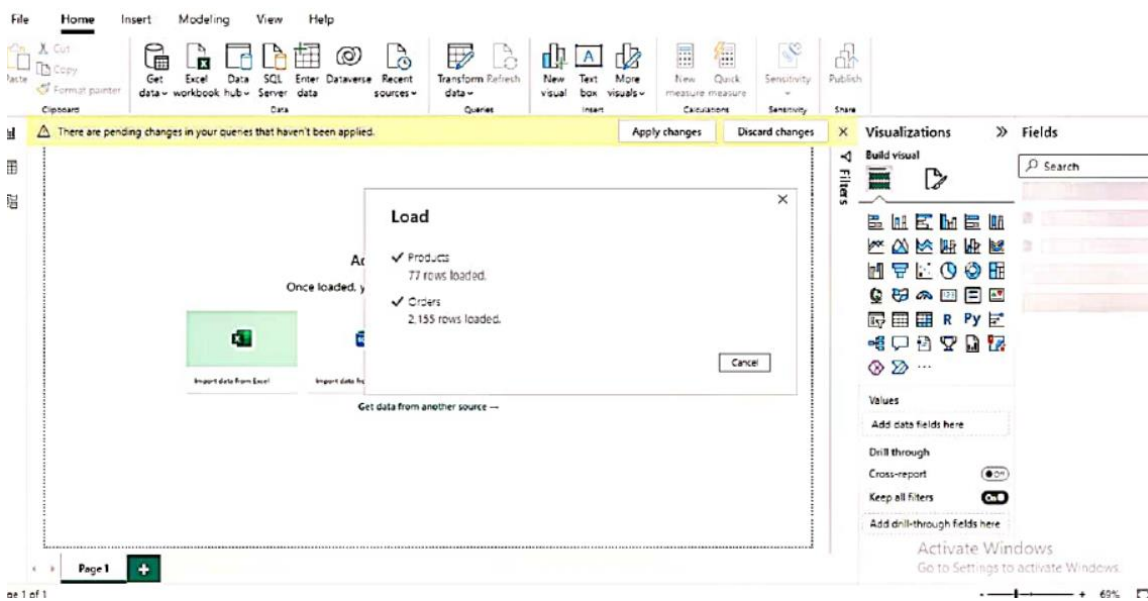
Removed Columns

PREVIEW DOWNLOADED AT 01



Load-:

Load



Practical 3

Aim-: Perform Data visualization from ETL (Using Power B.I)

Definition of Data Visualization-:

- 1) Data visualisation helps you turn all the granular data into easily understood, visually compelling and useful business information.
- 2) By using Power B.I. as a tool for Data visualization it helps us to see our KPIs more clearly, it unifies the data.
- 3) Data visualization brings data to life, making you the master storyteller of the insights hidden within your data.
- (4) Data Visualization helps users develop powerful business insight quickly and effectively.

Definition of ETL Process-:

- 1) Tableau Prep tool is a tool which excellent for Data Visualization and has some ETL capabilities in the Tableau.
- 2) ETL Tool is anacronym for Extract, Transform and load
- 3) In ETL Tool the data is collected from different heterogenous sources and transformed into information and loaded into data warehouses.

The screenshot displays the Microsoft Power BI Desktop interface. The main canvas shows a sunburst chart titled "COUNT OF PRODUCTID AND ... BY CUSTOMERID AN...". The chart is a circular sunburst with multiple segments, each labeled with "69 (1.25%)". A legend on the right side of the chart lists various customer IDs with corresponding colored circles: (Blank), ALFKI, ANATR, ANTON, AROUT, BERGS, BLAUS, BLONP, BOLID, BONAP, BOTTM, BSBEV, CACTU, CENTC, CHOPS, COMMI, and CONSH.

The right-hand pane is the "Visualizations" pane, which is currently set to "Build visual". It shows a list of fields on the right, including "Current_Product_Lists", "Customer_and_Suppli...", and "Customers". The "Customers" field is expanded, showing a list of fields: Address, City, CompanyName, ContactName, ContactTitle, Country, CustomerID (checked), Fax, Phone, and PostalCode. The "Visualizations" pane also shows a list of visual types on the left, including Bar chart, Line chart, Pie chart, and Sunburst chart. The "Sunburst chart" is selected. The "Filters" pane on the left shows a search bar and a list of filters: "CompanyName is (All)", "Count of CategoryID is (All)", "Count of ProductID is (All)", and "CustomerID is (All)".

Practical 4

Aim-:

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

Definition of Pivot Table-:

- 1) Pivot table is a simple tool used to create a Summarized report from a large set of databases.
- 2) Pivot Table is an interactive way to quickly Summarize large amounts of data.
- 3) Pivot tables in Power BI are used to turn rows into columns and unpivot columns into rows.

Definition of Pivot Charts-:

- 1) Pivot Charts are an in-built helps you summarize sale program tool that selected rows and spreadsheet columns of data in a spreadsheet.
- 2) It gives you the big picture of your raw data using various types of graphs and layouts.

Definition of Data Analysis-:

- 1) Data analysis is a transforming, and monitoring process of inspecting, valuable insights.
- 2) Data insights helps in making the required decision for the growth of the business & Company.
- 3) It involves subject data to obtain precise informed decisions or expand knowledge which can help businesses make on various subjects.
- 4) Data analysis has multiple facets and approaches.

Definition of Cube-:

- a) **ROLAP-:** It stands for Relational Online Analytical Processing and stores data in column, Rows.
- b) **MOLAP-:** It stands for Multidimensional online Analytical Processing and accesses data through various combinations.
- c) **HOLAP-:** It for Hybrid Online Analytical & is combination of both.

Output-:

4a)

PivotTable Fields

Choose fields to add to report:

- ☐ BusinessType
- ☐ Construction
- ☐ Earthquake
- ☐ Expiry
- ☐ Flood
- ☐ InsuredValue
- ☐ Location
- ☐ Policy
- ☐ Region

Drag fields between areas below:

Y FILTERS II COLUMNS

≡ ROWS Σ VALUES

Create PivotTable

Choose the data that you want to analyze

☒ Select a table or range

Table Range: [Table Range]

☐ Use an external data source

Choose where you want the PivotTable report to be placed

☐ New Worksheet

☒ Existing Worksheet

Location: Sheet1!\$A\$1

☐ Add this data to the Data Model

PivotTable Fields

Choose fields to add to report:

- ☐ Construction
- ☐ Expiry
- ☐ Flood
- ☒ InsuredValue
- ☒ Location
- ☐ Policy
- ☐ Region
- ☐ State

Drag fields between areas below:

Y FILTERS II COLUMNS

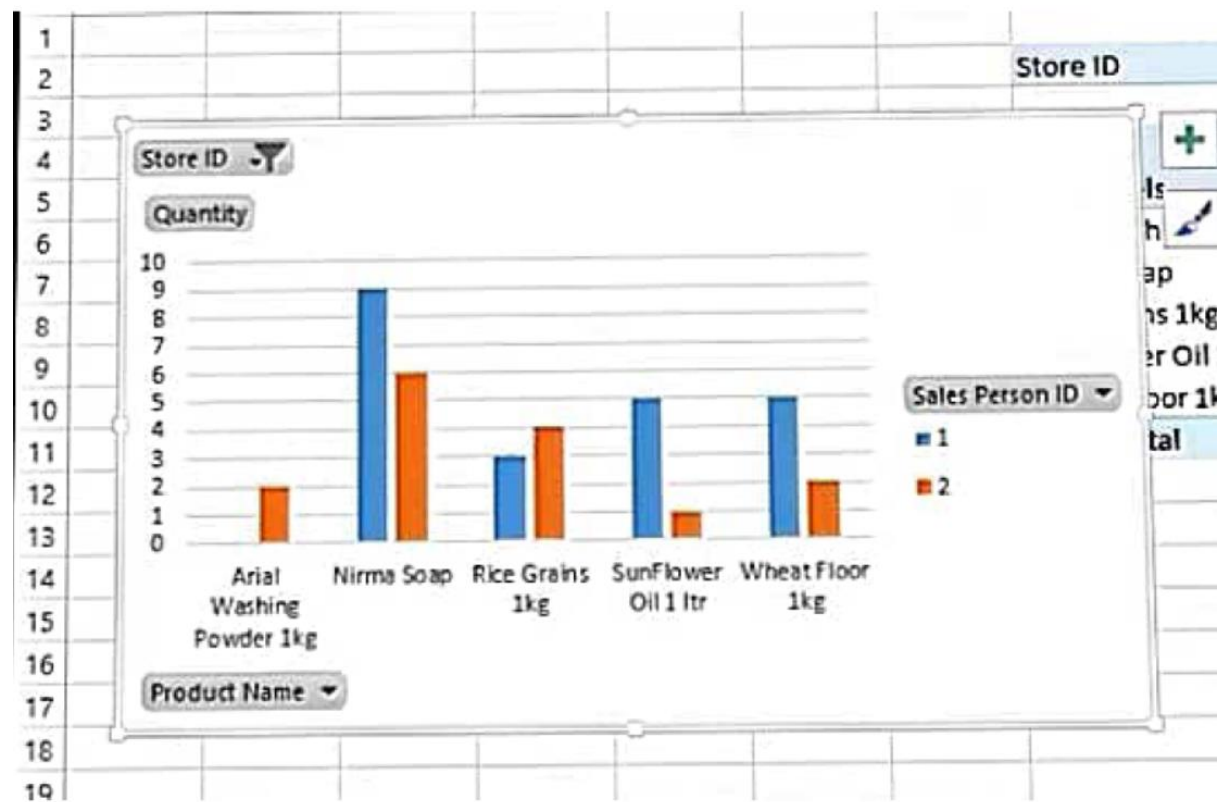
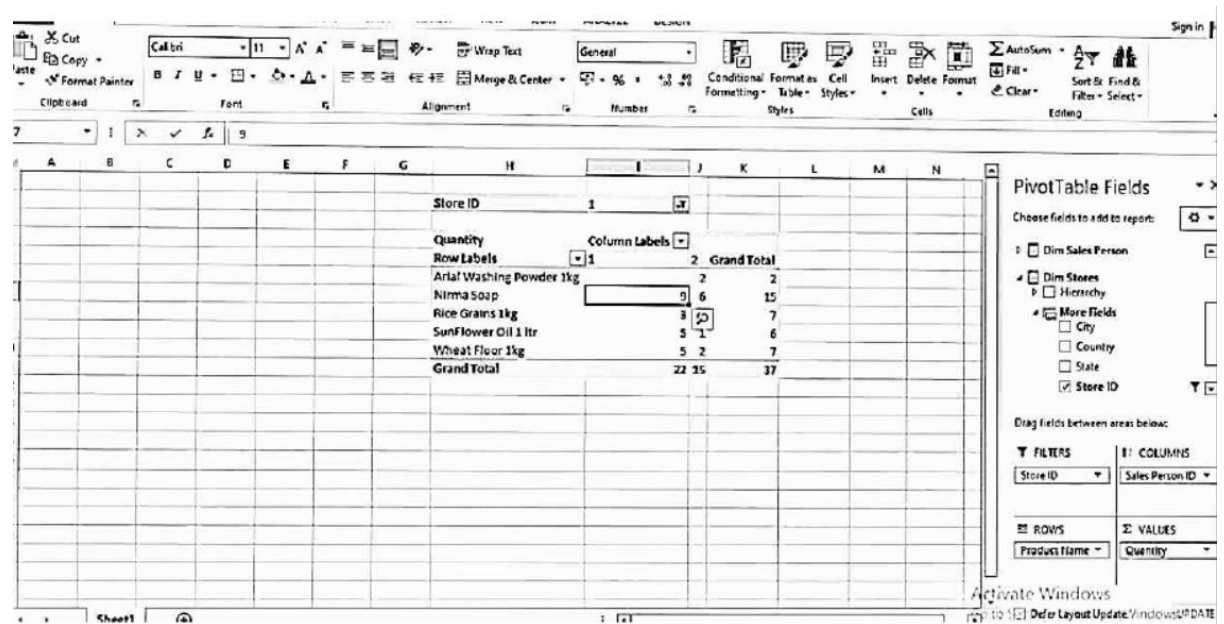
≡ ROWS Σ VALUES

BusinessType (Multiple Items)

Sum of InsuredValue

Row Labels	MI	MN	NJ	NY	VT	WI	Grand Total
Rural	3145700						3145700
Urban	3255300	1834200	6884700	35974222	10700000	23592848	82241270
Grand Total	3145700	3255300	1834200	6884700	35974222	10700000	23592848

4b)



Practical 5

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

Definition of What-if Analysis:

- 1) What-if analysis is often used to compare different scenarios and their potential outcomes based on changing conditions.
- 2) For using this feature, we can conditions. create What-it parameter and interact with the variable as a slicer.
- 3) Based on the input in the slicer, we can visualize and quantify changes in report data.

Definitions of Data Visualization-:

- 1) It is the representation of data through use of common graphics, such as charts, plots, infographic and even animation.
- 2) These visual displays of information communicate complex data relationships and data-driven to insight in a way easy to understand.
- 3) Visualizations are used to effectively present data and are the basic are the basic building your of any Business intelligence tool.

Output:-

Qty	Price	Total
25	20	500
35	20	700
40	30	1200

Goal Seek Status

Goal Seeking with Cell #11 found a solution.

Target value: 500

Current value: 500

OK Cancel

Sheet1 Sheet2 Sheet3

Qty	Price	Total
25	20	500
35	20	700
40	30	1200

Goal Seek

Set cell: F13

To value: 2000

By changing cell: \$D\$13

OK Cancel

Sheet1 Sheet2 Sheet3

Qty	Price	Total
25	20	500
35	20	700
66.66667	30	2000

Goal Seek Status

Goal Seeking with Cell F13 found a solution.

Target value: 2000

Current value: 2000

OK Cancel

	A	B	C	D	E
1	Book Store				
2					
3		total number of books	% sold for the highest price		
4		100	60%		
5					
6			number of books	unit profit	
7		highest price	60	\$50	
8		lower price	40	\$20	
9					
10			total profit	\$3,800	
11					

Practical 6

Aim-: Perform the data classification using classification algorithm.

Definition of data classification-:

- 1) Data classification is a method available in Power B.I. which allows users to tag dashboards. that alert consumers to sensitivity in their data.
- 2) Data classifications are enabled and configured at the tenant level, Once established, a visible tag will be present on the dashboards.
- 3) Data classification is not a data security implementation.
- 4) It is a tag for dashboards and can only be applied on the service, not on Power BI Desktop.

Definition of classification algorithm-:

- 1) The Classification algorithm is a Supervised Learning technique that is used to identify the category of new observations on the basis of training.
- 2) In Classification, a program learns from the given dataset or observations and then Classifies new observation into a no. of classes or groups.

Output-:

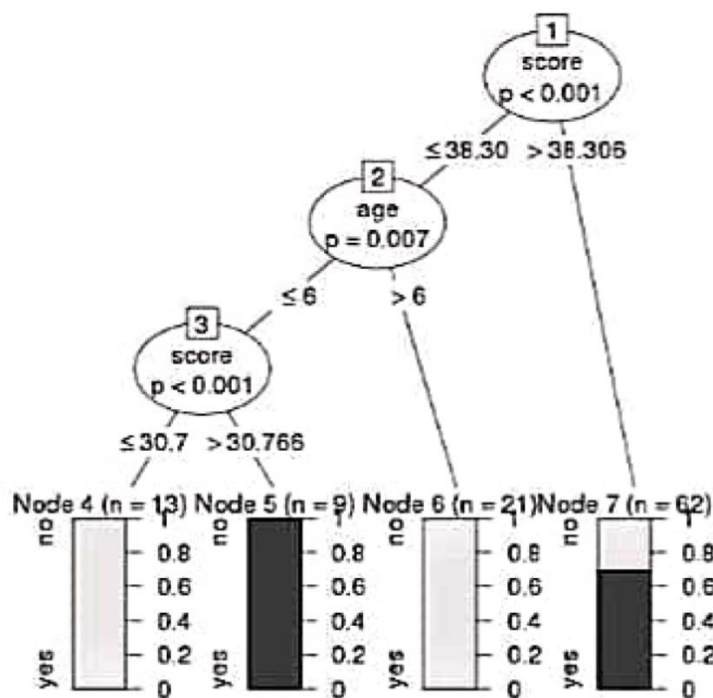
```
Loading required package: methods
Loading required package: grid
Loading required package: mvtnorm
Loading required package: modeltools
Loading required package: stats4
Loading required package: strucchange
Loading required package: zoo

Attaching package: 'zoo'

The following objects are masked from 'package:base':

  as.Date, as.Date.numeric

Loading required package: sandwich
```



```
# Load the party package. It will automatically load other
# dependent packages.
library(party)

# Create the input data frame.
input.dat <- readingSkills[c(1:105),]

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

# Create the tree.
output.tree <- ctree(
  nativeSpeaker ~ age + shoeSize + score,
```

Practical 7

Aim-: Perform Data Analysis using Time series Analysis

Definition of Time Series Analysis:-

- 1) Time series is a series of data each data point is associated with a timestamp.
- 2) Simple example is the price of a stock in the stock market at different points of time on a given data.
- 3) The time data for the time series stored in an R object called time-series object.
- 4) R language uses many plot functions to create the time series and manipulate data.
- 5) Following is the description of the parameters used-:
 - Data is a vector containing the value used in timeseries
 - Start specifies the start time for 1st observation in time series
 - End specifies the end time for last observation in time series
 - Frequency specifies the number of observation per unit.
- 6) The time series object is created by using ts() function.
The basic syntax for ts() function in time series analysis.

```
[timeseries.object.name<-ts(data,start,end,frequency)]
```

Output-:

```
Get the data points in form of a R vector.
rainfall <-
c(799,1174.8,865.1,1334.6,635.4,918.5,685.5,998.6,784.2,985,882.8,1071)

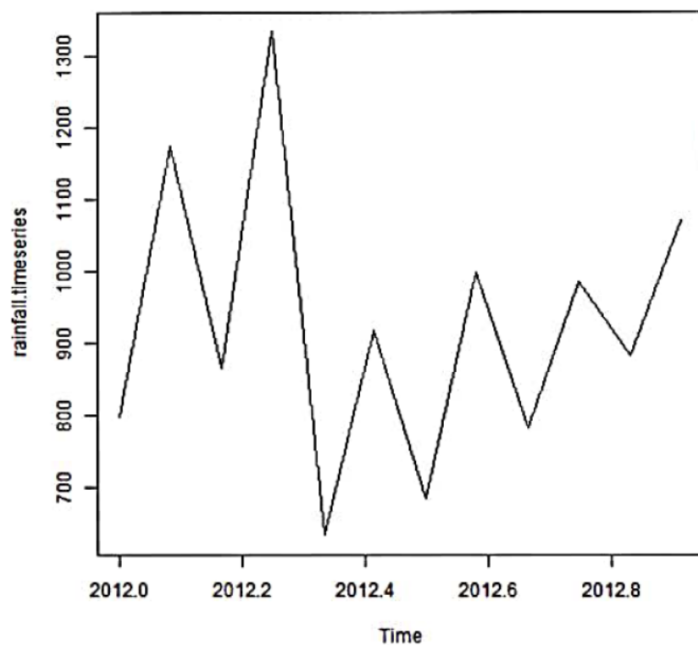
Convert it to a time series object.
rainfall.timeseries <- ts(rainfall,start = c(2012,1),frequency = 12)

Print the timeseries data.
print(rainfall.timeseries)

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

Plot a graph of the time series.
plot(rainfall.timeseries)

Save the file.
dev.off()
```



Jan	Feb	Mar	Apr	May	Jun	Jul
	Aug	Sep				
2012	799.0	1174.8	865.1	1334.6		
	Oct	635.4	918.5	685.5	998.6	
		784.2				
		Nov	Dec			
2012	985.0	882.8	1071.0			

Practical 8

Aim-: Perform the data clustering using clustering algorithm.

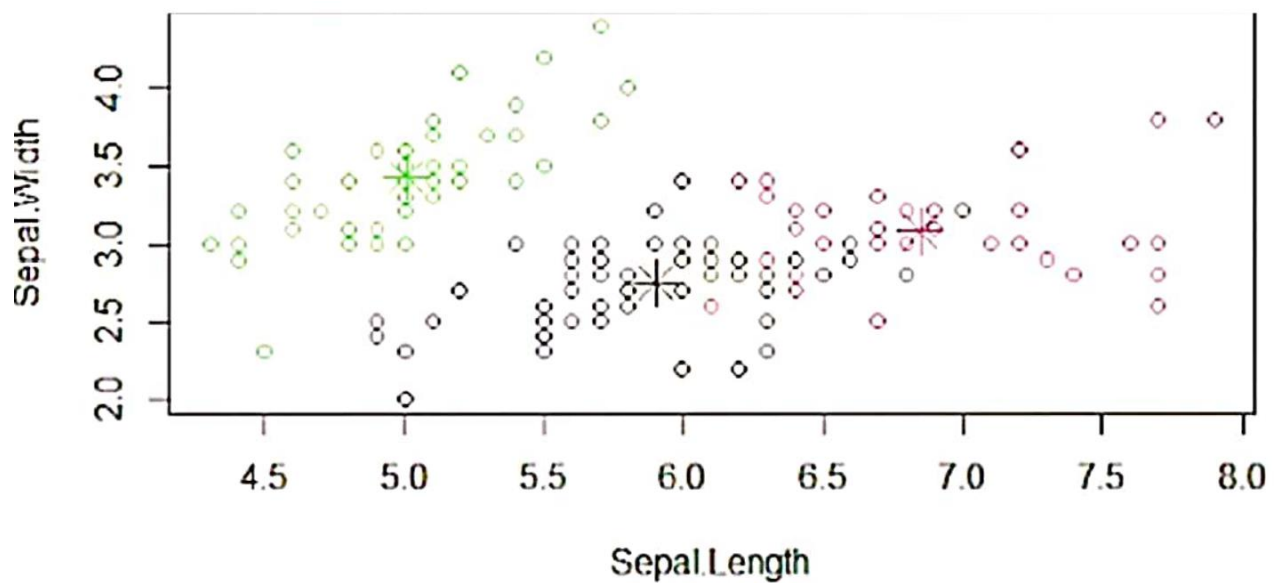
Definition of Data Clustering-:

- 1) Data clustering is an unsupervised machine learning algorithm that looks for patterns in data by dividing it into clusters.
- 2) These clusters are created such that the points are homogenous within the clusters and heterogenous across clusters.
- 3) It is commonly used in market segmentation, customer segmentation, image analysis, bioinformatics, data compression & computer graphics.

Definition of clustering algorithm-:

- 1) The Clustering algorithm provides two methods assigning for creating clusters and data points to the clusters.
- 2) The first, the K-means algorithm, is a hard clustering method. This means that data point can belong to only one clusters, and that single probability is calculated for the membership of each data point in that a cluster.

Output-:



```
#removing the species from data to cluster
newiris<-iris
newiris$Species<-NULL
#apply the k-means functions to newiris and store the result in kmeans.results
(kc<-kmeans(newiris,3))
print(kc)
#compare the clusters label
table(iris$Species,kc$cluster)
# plot the cluster and their centers
plot(newiris[,c("Sepal.Length","Sepal.Width")],col=kc$cluster)
# plot the cluster center
points(kc$centers[,c("Sepal.Length","Sepal.Width")],col=1:3,pch=8,cex=2)
```

Practical 9

Aim-: Perform the linear regression on the given data warehouse.

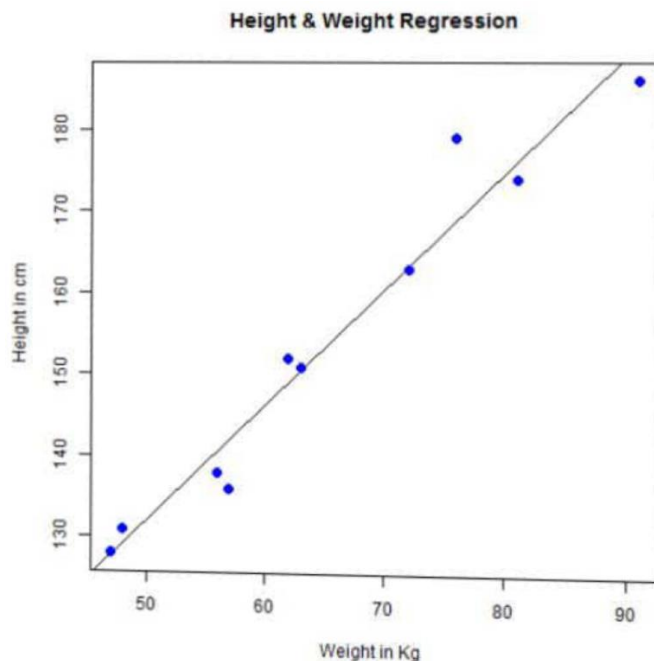
Definition of linear regression-:

- 1) linear regression is a statistical model apply to businesses to help forecast events based historical trend analysis.
- 2) Simple linear regression uses one variable, called the independent variable and other variable called as dependent variable. coordinate in a cartesian coordinate system and finds a linear function.

Definition of lm() function:-

- 1) The lm() function is used to fit linear models. to data frame in & R language.
- 2) It can be used to carry out regression, single stratum analysis of variance, and analysis of covariance to predict the value that is not in dataframe.

Output:-



```
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)
```

```
# The predictor vector.
x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)
```

```
# The response vector.
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)
```

```
# Create the predictor and response variable.
x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)
y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)
relation <- lm(y~x)
```

```
# 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 10

Aim:- Perform the logistic regression on the given data warehouse.

Definition of logistic regression:-

Logistic regression is a supervised machine learning algorithm that accomplishes binary classification task by predicting the probability of an outcome, even binary for observation. The model delivers a or dichotomous outcome limited to two possible outcomes, yes/no, 0/1, or true/false.

Dataset used:- mtcars

- 1) mtcars (motor trend car road sets)
- 2) There is a popular built-in data set in R called "mtcars" which is retrieved from the 1974 motor Trend Us Magazine.
- 3) The mtcars dataset is a built-in dataset in R. that contains measurements on 11 different attributes for 32 different cars.

GLM(Generalized linear Model):-

- 1) The glm (function in R can be used to generalized linear models.
- 2) This function is particularly useful for things for fitting logistic regression models, Poisson regression models and other complex models.

Output-:

```
> input=mtcars[,c("am","cyl","hp","wt")]
> print(head(input))
      am cyl  hp   wt
Mazda RX4      1   6 110 2.620
Mazda RX4 Wag  1   6 110 2.875
Datsun 710     1   4  93 2.320
Hornet 4 Drive  0   6 110 3.215
Hornet Sportabout 0   8 175 3.440
Valiant        0   6 105 3.460
> input=mtcars[,c("am","cyl","hp","wt")]
> am.data=glm(formula=am~cyl|hp|wt,data=input,family=binomial)
> summary(am.data)

Call:
glm(formula = am ~ cyl | hp | wt, family = binomial, data = input)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.021  -1.021  -1.021   1.342   1.342

Coefficients: (1 not defined because of singularities)
              Estimate Std. Error z value Pr(>|z|)
(Intercept)    -0.3795     0.3599  -1.054   0.292
cyl | hp | wtTRUE         NA         NA      NA      NA

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 43.23  on 31  degrees of freedom
Residual deviance: 43.23  on 31  degrees of freedom
AIC: 45.23

Number of Fisher Scoring iterations: 4
```

```
> mtcars
      mpg  cyl  disp  hp drat   wt  qsec vs am gear carb
Mazda RX4    21.0   6 160.0 110 3.90 2.620 16.46 0  1   4    4
Mazda RX4 Wag 21.0   6 160.0 110 3.90 2.875 17.02 0  1   4    4
Datsun 710   22.8   4 108.0  93 3.85 2.320 18.61 1  1   4    1
Hornet 4 Drive 21.4   6 258.0 110 3.08 3.215 19.44 1  0   3    1
Hornet Sportabout 18.7   8 360.0 175 3.15 3.440 17.02 0  0   3    2
Valiant      18.1   6 225.0 105 2.76 3.460 20.22 1  0   3    1
Duster 360   14.3   8 360.0 245 3.21 3.570 15.84 0  0   3    4
Merc 240D    24.4   4 146.7  62 3.69 3.190 20.00 1  0   4    2
Merc 230     22.8   4 140.8  95 3.92 3.150 22.90 1  0   4    2
Merc 280     19.2   6 167.6 123 3.92 3.440 18.30 1  0   4    4
Merc 280C    17.8   6 167.6 123 3.92 3.440 18.90 1  0   4    4
Merc 450SE   16.4   8 275.8 180 3.07 4.070 17.40 0  0   3    3
Merc 450SL   17.3   8 275.8 180 3.07 3.730 17.60 0  0   3    3
Merc 450SLC  15.2   8 275.8 180 3.07 3.780 18.00 0  0   3    3
Cadillac Fleetwood 10.4   8 472.0 205 2.93 5.250 17.98 0  0   3    4
```

