

Thadomal Shahani Engineering College

Bandra (W.), Mumbai- 400 050.

© CERTIFICATE ©

Certify that Mr./Miss OM ANIRUDHA SHETE
of COMPUTER Department, Semester ✓ with
Roll No.2103163 has completed a course of the necessary
experiments in the subject Data Warehouse Mining under my
supervision in the **Thadomal Shahani Engineering College**
Laboratory in the year 2023 - 2024


Teacher In-Charge

Head of the Department

Date 23/10/2023

Principal

CONTENTS

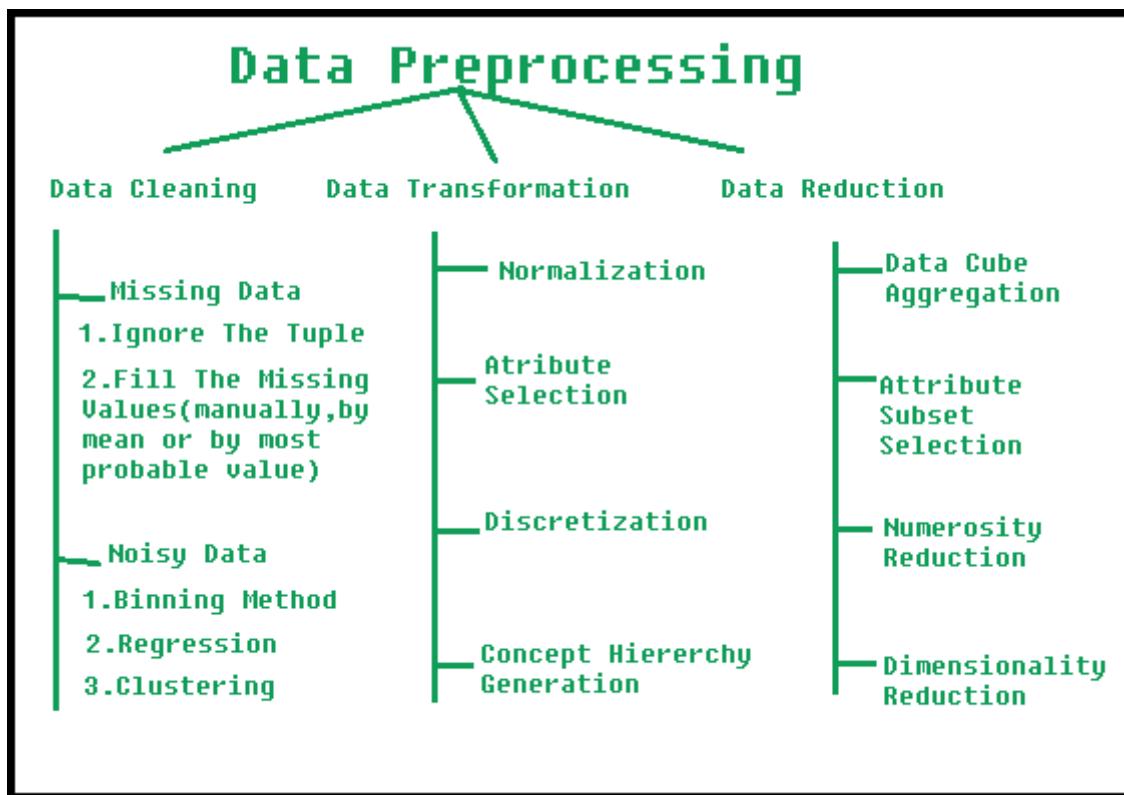
SR. NO.	EXPERIMENTS	PAGE NO.	DATE	TEACHERS SIGN.
1)	Select a dataset and perform explanatory data analysis using python.	1-6	21/8/23	
2)	One case study on building data warehouse / data marts	7-16	31/6/23	
3)	Implementation of all dimension table and fact table based on experiment 2 case study	17-43	10/8/23	
4)	Implementation of OLAP operational Slice, dice, rollup, drilldown, and pivot.	44-46	24/8/23	
5)	Implementation of Bayesian theorem	47-51	31/8/23	
6)	Perform data preprocessing and perform classification, clustering and association rules	52-57	5/9/23	(1)
7)	Implementation of clustering algorithm	58-62	14/9/23	
8)	Implementation of any one hierarchical clustering	63-68	21/9/23	
9)	Implementation of association rule mining algorithm	69-72	19/10/23	
10)	Implementation of Page Rank/Hits Algorithm.	73-76	19/10/23	
11)	Assignment 1	77-78	21/10/23	
12)	Assignment 2	79-85	21/10/23	

EXPERIMENT - 01

Aim: Select a dataset and perform exploratory data analysis using Python

Theory:

Data Preprocessing:



Data preprocessing is a crucial step in the data mining process. It involves cleaning, transforming, and preparing raw data to make it suitable for analysis. The steps involved include:

- Data cleaning: This involves identifying and correcting errors or inconsistencies in the data, such as missing values, outliers, and duplicates. Various techniques can be used for data cleaning, such as imputation, removal, and transformation.
- Data Integration: This involves combining data from multiple sources to create a unified dataset. Data integration can be challenging as it requires handling data with different formats, structures, and semantics. Techniques such as record linkage and data fusion can be used for data integration.

- Data Transformation: This involves converting the data into a suitable format for analysis. Common techniques used in data transformation include normalization, standardization, and discretization. Normalization is used to scale the data to a common range, while standardization is used to transform the data to have zero mean and unit variance. Discretization is used to convert continuous data into discrete categories.
- Data Reduction: This involves reducing the size of the dataset while preserving the important information. Data reduction can be achieved through techniques such as feature selection and feature extraction. Feature selection involves selecting a subset of relevant features from the dataset, while feature extraction involves transforming the data into a lower-dimensional space while preserving the important information.

Data Transformation:

Data transformation in data mining refers to the process of converting raw data into a format that is suitable for analysis and modelling. The goal of data transformation is to prepare the data for data mining so that it can be used to extract useful insights and knowledge. Data transformation typically involves several steps, including:

- Smoothing: It is a process that is used to remove noise from the dataset using some algorithms. It allows for highlighting important features present in the dataset. It helps in predicting the patterns. When collecting data, it can be manipulated to eliminate or reduce any variance or any other noise form. The concept behind data smoothing is that it will be able to identify simple changes to help predict different trends and patterns. This serves as a help to analysts or traders who need to look at a lot of data which can often be difficult to digest for finding patterns that they wouldn't see otherwise.
- Aggregation: Data collection or aggregation is the method of storing and presenting data in a summary format. The data may be obtained from multiple data sources to integrate these data sources into a data analysis description. This is a crucial step since the accuracy of data analysis insights is highly dependent on the quantity and quality of the data used. Gathering accurate data of high quality and a large enough quantity is necessary to produce relevant results. The collection of data is useful for everything from decisions concerning financing or business

strategy of the product, to pricing, operations, and marketing strategies. For example, Sales, data may be aggregated to compute monthly & annual total amounts.

- Normalization: Data normalization involves converting all data variables into a given range.
- Generalization: It converts low-level data attributes to high-level data attributes using concept hierarchy. For Example, Age initially in Numerical form (22, 25) is converted into categorical value (young, old). For example, Categorical attributes, such as house addresses, may be generalized to higher-level definitions, such as town or country.
- Attribute Construction: Where new attributes are created & applied to assist the mining process from the given set of attributes. This simplifies the original data & makes the mining more efficient.

Data Discretization:

Data discretization refers to a method of converting a huge number of data values into smaller ones so that the evaluation and management of data become easy. In other words, data discretization is a method of converting attribute values of continuous data into a finite set of intervals with minimum data loss.

- Histogram analysis: Histogram refers to a plot used to represent the underlying frequency distribution of a continuous data set. Histogram assists the data inspection for data distribution. For example, Outliers, skewness representation, normal distribution representation, etc.
- Binning: Binning refers to a data smoothing technique that helps to group a huge number of continuous values into smaller values. For data discretization and the development of idea hierarchy, this technique can also be used.

Data Visualization:

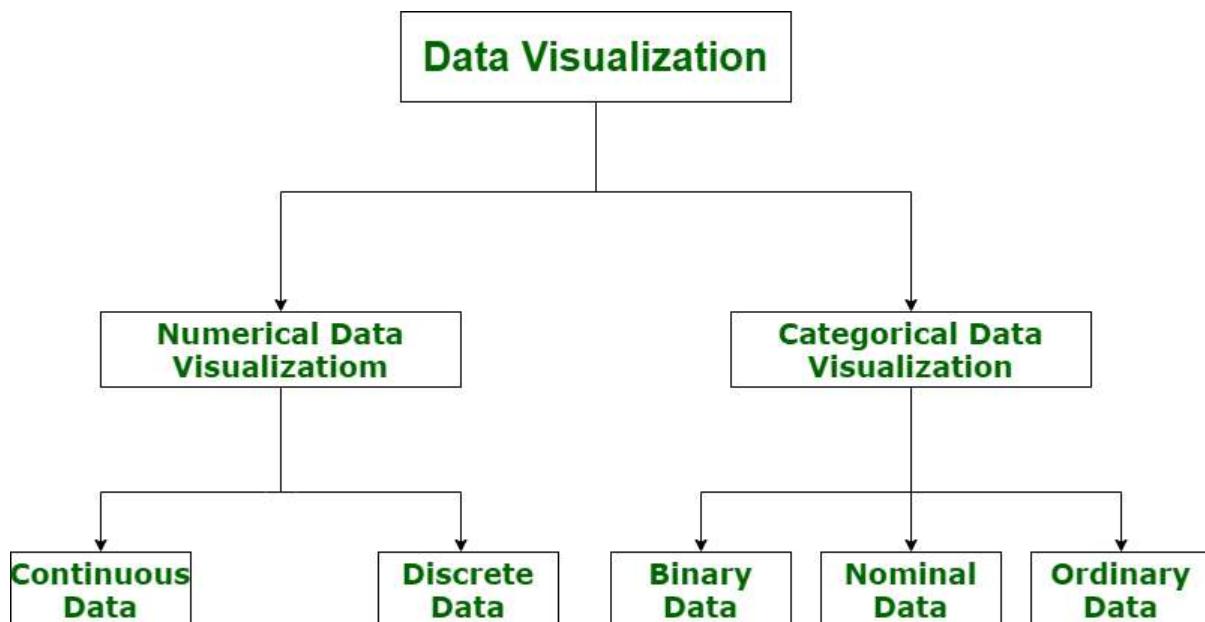
Data visualization is actually a set of data points and information that are represented graphically to make it easy and quick for user to understand. Data visualization is good if it has a clear meaning, purpose, and is very easy to interpret, without requiring context. Tools of data visualization provide an accessible way to see and understand trends, outliers, and patterns in data by using visual effects or elements such as a chart, graphs, and maps.

Categories of Data Visualization :

- Numerical Data : Numerical data is also known as Quantitative data. Numerical data is any data where data generally represents amount such as height, weight, age of a person, etc. Numerical data visualization is easiest way to visualize data. It is generally used for helping others to digest large data sets and raw numbers in a way that makes it easier to interpret into action. Numerical data is categorized into two categories :

1. Continuous Data- It can be narrowed or categorized (Example: Height measurements).
2. Discrete Data- This type of data is not “continuous” (Example: Number of cars or children’s a household has).

The type of visualization techniques that are used to represent numerical data visualization is Charts and Numerical Values. Examples are Pie Charts, Bar Charts, Averages, Scorecards, etc.



- Categorical Data : Categorical data is also known as Qualitative data. Categorical data is any data where data generally represents groups. It simply consists of categorical variables that are used to represent characteristics such as a person's ranking, a person's gender, etc. Categorical data visualization is all about depicting key themes, establishing connections, and lending context. Categorical data is classified into three categories :
1. Binary Data- In this, classification is based on positioning (Example: Agrees or Disagrees).
 2. Nominal Data- In this, classification is based on attributes (Example: Male or Female).
 3. Ordinal Data- In this, classification is based on ordering of information (Example: Timeline or processes).

The type of visualization techniques that are used to represent categorical data is Graphics, Diagrams, and Flowcharts. Examples are Word clouds, Sentiment Mapping, Venn Diagram, etc.

Code:

```
import csv
import random
from faker import Faker
from datetime import datetime, timedelta

# Set the number of entries for each
dimension num_products = 1000
num_warehouses = 100
num_suppliers = 500
num_contacts = 1000
num_entries =

1000 fake = Faker()

# Function to generate a unique ID
def generate_unique_id(existing_ids):
    new_id = fake.random_int(min=1000,
    max=9999) while new_id in existing_ids:
        new_id = fake.random_int(min=1000,
        max=9999) existing_ids.add(new_id)
    return new_id

# Function to generate random date within a
range def generate_random_date(start_date,
end_date):
    time_delta = end_date - start_date
    random_days = random.randint(0,
    time_delta.days) return start_date +
    timedelta(days=random_days)

# Generate data and write to a single CSV file
with open('inventory_data.csv', mode='w', newline='') as
    file: writer = csv.writer(file)
        writer.writerow(['transaction_id', 'product_id', 'product_name',
"product_category", "product_brand", "warehouse_id", "warehouse_name",
"warehouse_location", "city_id", "supplier_id", "supplier_name", "contact_id", "phone",
"email", "arrival_date",
```

```

"dispatch_date", "number_of_boxes"])

product_ids = set()
warehouse_ids =
set() supplier_ids =
set()

contact_ids = set()

for i in
range(num_entries):
transaction_id = i + 1

# Generate product data
product_id =
generate_unique_id(product_ids)
product_name = fake.word()
product_category = fake.word()
product_brand =
fake.company()

# Generate warehouse data
warehouse_id =
generate_unique_id(warehouse_ids)
warehouse_name = fake.company()
warehouse_location = fake.address()
city_id =
generate_unique_id(contact_ids)

# Generate supplier data
supplier_id =
generate_unique_id(supplier_ids)
supplier_name = fake.company()

# Generate contact data
contact_id = generate_unique_id(contact_ids)
phone = fake.phone_number()
email = fake.email()

# Generate transaction date data
arrival_date = generate_random_date(datetime(2020, 1, 1), datetime(2023, 8, 1))
dispatch_date = generate_random_date(arrival_date, datetime(2023, 8,

```

1)) number_of_boxes = random.randint(1, 100)

```
writer.writerow([transaction_id, f"P{product_id}", product_name, product_category,
product_brand, f"W{warehouse_id}", warehouse_name, warehouse_location,
f"C{city_id}", f"S{supplier_id}", supplier_name, f"C{contact_id}", phone, email,
arrival_date.strftime('%Y-%m-%d'), dispatch_date.strftime('%Y-%m-%d'),
number_of_boxes])
```

Output:

Original Dataset:

	A	B	C	D	E	F	G	H	I	J	K	L	M	
1	transaction_id	product_id	product_name	product_category	product_brand	warehouse_id	warehouse_name	warehouse_address	city_id	supplier_id	supplier_name	contact_id	phone	email
2	1P2899	Titan T27	watches	David LLC	W7276	BonneyX	Mumbai	FSC 4244, Box 5124	C7597	53535	Stewart, Clayton and Martinez	C2483	001-954-497-8463	xbamri@example.com
3	2P5937	Samsung QLED 21	tv	Reynolds and Sons	W7001	Haasy and Sons	AP0 AP 43657	4002 Julia Park	C7099	53793	Carroll, Brady and Hancock	C3454	601-443-1307x13981	vherando@exam
4	3P1553	Philips Iron H	Electrical App	Greene Hernandez	W1296	Lucas Inc	8000 8th Street, SU 888-0	47251 Thomas Avenue	C6006	52142	Singh, Johnson	C1265	(645)485-0326x71214	cospelan@exam
5	4P7296	Nike Neo	Wear	Weeks Group	WB451	Norris LLC	North Main Street, NH 32249	4003 10th Avenue	C8093	53115	Boyer Ltd	C3118	(869)959-1199x547	monacod6@exam
6	5P9813	US Polo Black Denim	Wear	Doyal Ltd	WB571	Brewer McDonald	371 1st Street, Suite 1650	30282 Bartaria Trail	C6338	58822	Davis Ltd	C7263	6267958495x1940@exam	renanabed@exam
7	6P5209	Samsung Washing Machine	Electrical App	Section Unit	WA498	Graeber Thomas	9007 8th Street, NJ 08229	20282 Bartaria Trail	C5676	54540	Senior, Brown and Peters	C7667	959-354-0180x12000	leopoldo@exam
8	7P1489	Balance Hair Dryer	Electrical App	Percussion, Thomas and Bryant	WB100	Marquez Ross	2nd Floor, 10th Avenue	3001 2nd Avenue, Apt. 362	C0001	53887	Madden Ltd	C1342	001-518-337-3662x8615	jeffreyt@exam
9	8P9353	LG Oven P32	tv	West, Howard and Huff	WB195	Willis Shaw	700 7th Avenue	10th Floor, 10th Avenue	C3348	53681	Pearson, Wolf and Shelton	C3949	904-831-7333x9687	heavymix@exam
10	9P4413	Sony Bravia S76	tv	Palmer, Hanna and Smith	WB184	Fields Ellis	100 7th Street, Box 4902	100 7th Street, Box 4902	C6039	58009	Ellis Ltd	C1699	409-236-0115	thompsonach@exam
11	10P3825	Fossil F31	watches	Johnson, Rutz and Todd	WB206	Ramirez, Cooper and Coogan	6452 Megan Underpass	West Haven, NY 53595	C2366	53143	Johnson Group	C5700	6193156094xrobert@exam	robert@exam
12	11P7873	program	occult	Quinn, Buster and Gregory	WB578	Gilmore-Woods	800 8th Street, Box 2001	400 8th Avenue, Apt. 2001	C4486	58420	Mccoy LLC	C9238	+1-410-300-4668x5728	thomasmartin@exam
13	12P1596	next	challenge	Brown Inc	WB169	Jones Inc	00074 Bennett Circles	South Portland, ME 04105	C0004	53161	Smit Ltd.	C7601	+1-85-763-3704x2089	stevies@exam
14	13P3993	visit	stop	Coleman LLC	WB639	Powell LLC	78349 Print Port	78349 Print Port	C0003	58841	Small PLC	C5995	+1-389-224-5753x751	kmsn32@exam
15	14P4110	always	next	Cook Ltd	WA420	Weaver PLC	700 7th Street, Apt. 1044	700 7th Street, Apt. 1044	C5176	53895	Carney, Edwards and Chen	C9699	(869)951-5017	natalie47@exam
16	15P235	data	right	Johnson, Ross and Phillips	WB572	Cotley Flores	002 Mid Mount	Port Matthew, HI 98471	C2046	58996	Ruiz Garcia	C1627	8268290411xlongstrey@exam	longstrey@exam
17	16P2301	education	top	Detwiler, Lane and Jacobson	WB373	Webb-Costa	East 10th Street, NM 58734	7739 Brandon Pass Suite 102	C2088	53534	Pater and Sons	C5607	(550)487-4208x82519	coreywilson@exam
18	17P5883	memory	however	King and Sons	WB067	Garcia-Morales	591 Clarkie Radnor	591 Clarkie Radnor	C3778	53719	Rose-Banter	C9300	(646)541-0134	irbyanderson@exam
19	18P7062	in	Mr	Copeland-Evans	WB176	Gonzalez Ltd	590 5th Street, UT 80021	590 5th Street, UT 80021	C4822	56126	Witteman, Wray and Martin	C8900	866-525-1911	bucksharr@exam
20	19P8087	collector	certain	Lawrence PLC	WA563	Gilespie and Sons	PSC 2647, Box 6951	PSC 2647, Box 6951	C4710	53204	Pice, Sparks and Frey	C2979	001-712-557-3078	akbar06@exam
21	20P4582	surface	number	Fowler, Todd and Williams	WB152	Ritter-Bishop	2241 The Mission Mission Suite 1	2241 The Mission Mission Suite 1	C5170	53996	Newton Group	C8997	952-647-3640x389	duspmptson@exam
22	21P2828	research	rock	Brennan, Johnson and Odonnell	WB480	Greer and Sons	3649 Smart Mills Suite 701	3649 Smart Mills Suite 701	C4203	57329	Walker-Espinoza	C9008	001-881-9077x81268	sheawbey@exam
23	22P4013	executive	behavior	Hernandez, McGee and Ray	WB157	Gozar Evans and Ramirez	310 Axon Station Suite 072	310 Axon Station Suite 072	C4485	53289	Taylor Ltd	C7302	331-511-0151x1354	adrian27@exam
24	23P8866	sense	light	Johnson-Sosa	WB268	Lucas, White and Smith	Timothy, UT 80023	Timothy, UT 80023	C9650	53005	Hernandez, Williams and Martinez	C8994	782-886-6325x408	mmyas@exam
25	24P1531	clear	dark	Rivera-Bowers	WB269	Janes Carter	7715 Greene Haven Suite 103	7715 Greene Haven Suite 103	C7781	52273	Lee PLC	C7742	+1-562-509-2362	wrightson@exam
26	25P2845	boy	candidate	Nelson-Cantillo	WB737	Wilson, Smith and Varon	Admiral, MA 42955	Admiral, MA 42955	C5620	53180	Baley LLC	C7185	+1-882-693-8171x943	laure50@exam

After data preprocessed:

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R		
1	transaction_id	product_id	product_name	product_category	product_brand	warehouse_id	warehouse_name	warehouse_address	city_id	supplier_id	supplier_name	contact_id	phone	email	arrival_date	dispatch_date	data_number_of_boxes		
2	1P2899	Titan T27	watches	Team	WB7276	BonneyX	Mumbai	C7597	53535	Stewart, Clayton and Martinez	C2483	001-954-497-8463	xbamri@example.com	2021-11-29	2023-01-15	16			
3	2P5937	Samsung QLED 21	tv	Samsung	W7001	BonneyX	AP0 AP 43657	4002 Julia Park	C7099	53793	Carroll, Brady and Hancock	C3454	601-443-1307x13981	vherando@exam	2021-09-18	2022-05-26	77		
4	3P1553	Philips Iron H	Electrical App	Greene Hernandez	W1296	Bangalore	C9293	53115	Boyer Ltd	C3118	(985)459-1156x547	monica64@example.net	2021-06-21	2023-10-13	71				
5	4P7296	Nike Neo	Wear	Weeks Group	WB451	Bangalore	C6338	58822	Stewart, Clayton and Martinez	C2483	001-954-0180x36615	lucyanderson@exam	2022-10-29	2023-11-07	37				
6	5P9813	US Polo Black Denim	Wear	US Polo	WB571	KeralaS	CE366	59622	Davis Ltd	C1342	001-881-9077x81268	leopoldo@exam	2022-09-01	2023-09-11	73				
7	6P5209	Samsung Washing Machine	Electrical App	Shapiro	WB498	Gurgoen IDES	C5095	52603	Madson Inc	C1342	001-881-337-3662x615	lucyanderson@exam	2020-10-18	2023-05-16	5				
8	7P1489	Balance Hair Dryer	Electrical App	Reliance	WB100	Hyderabad	C5955	52003	Stewart, Clayton and Martinez	C5170	53895	Carney, Edwards and Chen	C9699	952-647-3640x389	duspmptson@exam	2020-09-18	2022-03-25	31	
9	8P9353	LG 40inch P32	tv	LG	WA4035	BonneyX	CD100	53566	Pearson, Wolf and Shelton	C5170	53996	Newton Group	C8997	952-647-3640x389	duspmptson@exam	2020-09-18	2022-03-25	31	
10	9P4413	Sony Bravia S76	tv	Palmer, Hanna and Smith	WB184	Fields Ellis	CD100	53566	Rosie-Banter	C2046	58996	Ruiz Garcia	C1627	8268290411xlongstrey@exam	longstrey@exam	2020-09-18	2022-03-25	31	
11	10P3825	Fossil F31	watches	Fossil	WB206	BonneyZ	Mumbai	C2366	53145	Johnson Group	C4203	57329	Walker-Espinoza	C9008	001-881-9077x81268	sheawbey@exam	2020-09-18	2022-03-25	31
12	11P7873	Gardien C21	watches	Gardien	WB578	Bangalore	C9270	53226	Deveraux	C4485	53273	Witteman, Wray and Martin	C8994	782-886-6325x408	mmyas@exam	2020-09-18	2022-03-25	31	
13	12P1596	Ompleta O20	tv	Ompleta	WB563	KeralaS	C4710	53284	Price, Sparks and Frey	C5238	5613150094	hoffmann7@exam.com	2020-09-18	2022-03-25	31				
14	13P3993	Omega C20	tv	Ompleta	WB169	Gurgoen IDES	C5099	56927	Goodman LLC	C3997	562-547-3640x3645	luisperalta@exam.com	2020-09-18	2022-03-25	31				
15	14P4110	NanCook Smart G21	Electrical App	Autocook	WA420	Hyderabad	C5270	53995	Levy, Williams and Mariano	C4710	53273	MacKenzie and Lewis	C5995	952-647-3640x389	duspmptson@exam	2020-09-18	2022-03-25	31	
16	15P235	CRMX Home	Wear	CRMX	WB195	Bangalore	C5270	53995	Carney, Edwards and Chen	C5170	53996	Newton Group	C8997	952-647-3640x389	duspmptson@exam	2020-09-18	2022-03-25	31	
17	16P2301	CRMX Home	Wear	CRMX	WB073	Dehraps	C2088	53304	Poter and Sons	C5667	050-487-8200x2519	coreywilson@exam.org	2020-04-03	2022-06-17	74				
18	17P5883	Gardeon C26	watches	Gardeon	WB067	Dehraps	C2090	53718	Rosie-Banter	C3000	562-547-3640x3614	benjamin7@exam.com	2020-10-29	2022-11-27	86				
19	18P7062	Gardeon C23	tv	Gardeon	WB508	Bangalore	C5270	53995	Witteman, Wray and Martin	C4485	53273	Walker-Espinoza	C9008	001-881-9077x81268	sheawbey@exam	2020-10-29	2022-11-27	86	
20	19P087	Ompleta O20	tv	Ompleta	WB563	KeralaS	C4710	53284	Price, Sparks and Frey	C5238	5613150094	hoffmann7@exam.com	2020-03-26	2023-06-24	18				
21	20P4582	Gardeon C25	watches	Gardeon	WB152	Gurgoen IDES	C5099	56927	Goodman LLC	C3997	562-547-3640x3645	luisperalta@exam.com	2020-03-26	2023-06-24	18				
22	21P2828	Gardeon C21	tv	Gardeon	WB578	Bangalore	C5270	53995	Deveraux	C4485	53273	Witteman, Wray and Martin	C8994	782-886-6325x408	mmyas@exam	2020-03-26	2023-06-24	18	
23	22P4013	Samsung L21	tv	Samsung	WB169	Bangalore	C7265	59110	Carney, Edwards and Chen	C3778	531-511-0151x1354	luisperalta@exam.com	2020-03-11	2023-07-06	85				
24	23P5666	Philips Iron E2	Electrical App	Philips	WB578	Bangalore	C7265	59110	Blaize, Jacobson and Long	C796	290-772-054-392	mackenzie7@exam.com	2020-08-20	2023-09-17	1				
25	24P2237	Agfa Photo Camera B	watches	Agfa	WB135	Dehraps	C2179	56932	Bryant, McCoy and Hayes	C8797	001-860-736-9552x7159	luisperalta@exam.com	2020-08-20	2023-08-27	96				
26	25P4664	LG 40inch P32	tv	LG	WA305	Dehraps	C5270	53995	Deveraux	C4485	53273	Witteman, Wray and Martin	C8994	782-886-6325x408	mmyas@exam	2020-08-20	2023-08-27	96	
27	26P4582	Ompleta Wash Machine	Electrical App	Ompleta	WB152	BonneyX	C7265	59110	Carney, Edwards and Chen	C3778	531-511-0151x1354	luisperalta@exam.com	2021-12-02	2023-06-27	66				
28	27P4413	Gardeon T22	tv	Gardeon	WB152	BonneyX	C7265	59110	Deveraux	C4485	53273	Witteman, Wray and Martin	C8994	782-886-6325x408	mmyas@exam	2021-12-02	2023-06-27	66	
29	28P4474	Sony Bravia A19	watches	Sony	WB174	BonneyX	C1886	53878	Witteman, Wray and Martin	C4484	53273	Deveraux	C8994	782-886-6325x408	mmyas@exam	2022-10-22	2023-03-09	83	
30	29P1063	Sonata Men M23	tv	Sonata	WB152	BonneyX	C1886	53878	Deveraux	C4484	53273	Witteman, Wray and Martin	C8994	782-886-6325x408	mmyas@exam	2022-10-22	2023-03-09	83	
31	30P2973	Mi Andromeda	tv	Andromeda	WB508	Bangalore	C5658	53813	Richard, Munoz and Kelley	C4658	338-297-9489	australia6@exam.org	2023-03-02	2023-07-17	32				
32	31P2747	Haier Dryer D504E	Electrical App	Philips	WB152	Bangalore	C5658	53813	Richard, Munoz and Kelley	C4658	338-2								

IMPLEMENTATION

DWM_Expt_Remote.ipynb

```

File Edit View Insert Runtime Tools Help All changes saved
Files + Code + Text
Loading the libraries and Download the dataset from Kaggle or other sources
[1]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np

Read the file and select appropriate file read function according to the datatype of file

[4]: df= pd.read_csv('content/inventory.csv')

Describe the attribute names

[5]: df.columns
Index(['transaction_id', 'product_id', 'product_name', 'product_category',
       'product_brand', 'warehouse_id', 'warehouse_name', 'warehouse_location',
       'city_id', 'supplier_id', 'supplier_name', 'contact_id', 'phone',
       'email', 'arrival_date', 'dispatch_date', 'number_of_boxes'],
      dtype='object')

[6]: df.sample(10)
transaction_id product_id product_name product_category product_brand warehouse_id warehouse_name warehouse_location city_id supplier_id supplier_name contact_id phone
129 129 P0248 Fossil F31 watches Fossil W9155 DelhiCaps Delhi C7099 S7412 Whitehead, Williams and Lewis C0982 732256-6491xx677 william72
114 115 P0369 DNMX Hoodie Wear DNMX W2275 DelhiCaps Delhi C7099 S4025 Campbell, Mann and Briggs C6315 704-531-7361 zachary41
80 61 P0814 Sony 55inch P79 tv Sony W8001 KeralaSwap Kerala C6336 S5648 Lucas-Pratt C6089 -3467 catherine35

```

DWM_Expt_Remote.ipynb

```

File Edit View Insert Runtime Tools Help All changes saved
Files + Code + Text
[6]: df.sample(10)
transaction_id product_id product_name product_category product_brand warehouse_id warehouse_name warehouse_location city_id supplier_id supplier_name contact_id phone
129 129 P0248 Fossil F31 watches Fossil W9155 DelhiCaps Delhi C7099 S7412 Whitehead, Williams and Lewis C0982 732256-6491xx677 william72
114 115 P0369 DNMX Hoodie Wear DNMX W2275 DelhiCaps Delhi C7099 S4025 Campbell, Mann and Briggs C6315 704-531-7361 zachary41
80 61 P0814 Sony 55inch P79 tv Sony W8001 KeralaSwap Kerala C6336 S5648 Lucas-Pratt C6089 -3467 catherine35
125 126 P1313 Hair-Hair Dryer T2 Electrical App Hale Jack & Jones Denim Pants wear Jack & Jones W2573 Gurgaon-DE2 Gurgaon C5676 52262 Dean-Figueras C1252 (666)955-3280 karen98
160 161 P1052 Jack & Jones Denim Pants wear Jack & Jones W2573 Gurgaon-DE2 Gurgaon C5676 98547 Perez-Martin C2843 959-369-3791 scottrogen
14 15 P1235 DNMX Shell S21 Wear DNMX W5472 BombayX Mumbai C7997 S9566 Ruiz-Garcia C1627 8288290416 longastyle
17 18 P7062 Sonata Kit K2 watches Sonata W1796 BangaloreS Bangalore C5293 96126 Williamson, Ware and Martin C8100 866-525-1911 barkerstand
195 196 P0145 LG hair dryer Electrical App LG W9621 Gurgaon-DE2 Gurgaon C5676 85845 Dixon and Sons C3036 001-946-351-1630 danielwashington
59 60 P4621 Rolex R-52 watches Rolex W7062 BangaloreS Bangalore C5293 88931 Sandow Group C4931 (552)286-7551xx1701 gnat
132 130 P4124 Titan WD watches Titan W3569 Gurgaon-DE2 Gurgaon C5676 89627 Dovit PLC C6994 #CRW90 brandaberry

```

DWM_Expt_Remote.ipynb

```

File Edit View Insert Runtime Tools Help All changes saved
Files + Code + Text
[7]: df.head()
transaction_id product_id product_name product_category product_brand warehouse_id warehouse_name warehouse_location city_id supplier_id supplier_name contact_id phone
0 1 P2559 Titan T27 watches Titan W7276 BombayX Mumbai C2597 S3395 Stewart, Clayton and Martinez C2483 001-946-497-8463 xbarrettg
1 2 P5937 Samsung QLED 21 tv Samsung W7001 DelhiCaps Delhi C7099 S7973 Carol, Brady and Hancock C0454 601-443-1007xx3981 vitemanuel0
2 3 P1553 Philips Iron 14 Electrical App Philips W1296 BombayZ Mumbai C2597 S2142 High-Johnson C0265 (645)485-0256xx1014 cooperdavid0
3 4 P7296 Nike Nix Wear Nike W6451 BangaloreS Bangalore C3293 S3115 Boyer L0 C3118 (969)559-1196xx547 monica666
4 5 P9813 US Polo Black Denim Wear US Polo W6571 KeralaSwap Kerala C6336 S8922 Davis Ltd C7283 5267955495 kiranraj

```

Count number of Values

```

[8]: unique_value_counts = df.unique()
print(unique_value_counts)

transaction_id 366
product_id 286
product_name 84
product_category 5
product_brand 33
warehouse_id 286
warehouse_name 7
warehouse_location 6
city_id 5
supplier_id 286
supplier_name 286
contact_id 286
phone 177
email 200

```

C32_Om_2103163_DWM

DWM_Expt1_Remote.ipynb

```

File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
Count number of Values
[1] unique_value_counts = df.unique()
print(unique_value_counts)

transaction_id    200
product_id       200
product_name      88
product_category   5
product_brand      32
warehouse_id      200
warehouse_name     7
warehouse_location  6
city_id           5
supplier_id       200
supplier_name      196
contact_id        200
phone             217
email            200
arrival_date      185
dispatch_date     180
number_of_boxes     98
dtype: int64

Checking for Null Entries
[2] df.isnull().sum()

transaction_id    a
product_id       a
product_name      a
product_category   a
product_brand      a
warehouse_id      a
warehouse_name     a
warehouse_location  a
city_id           a
supplier_id       a
supplier_name      a
contact_id        a
phone             a
email            a
arrival_date      a
dispatch_date     a
number_of_boxes     a
dtype: int64
0s completed at 15:57

```

DWM_Expt1_Remote.ipynb

```

File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
Eliminating rows having value of any column as zero or null
[3] df.dropna(inplace=True)

Counting duplicated entries
[4] df.duplicated().sum()
0

Command to remove duplicate entries
[5] df.drop_duplicates(keep='First', inplace=True)

Retrieving the no of columns and no of rows using df.shape
[6] df.shape
(200, 17)

Finding MIN and MAX for 'number_of_boxes'
[7] maximum_value=df['number_of_Boxes'].max()
minimum_value=df['number_of_Boxes'].min()
print("Minimum value for Number of Boxes is ", minimum_value)
print("Maximum value for Number of Boxes is ", maximum_value)

Minimum value for Number of Boxes is 1
Maximum value for Number of Boxes is 99

Printing the Datatypes of fields/columns
[8] df.dtypes
0s completed at 15:57

```

DWM_Expt1_Remote.ipynb

```

File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
Printing the Datatypes of fields/columns
[9] df.dtypes

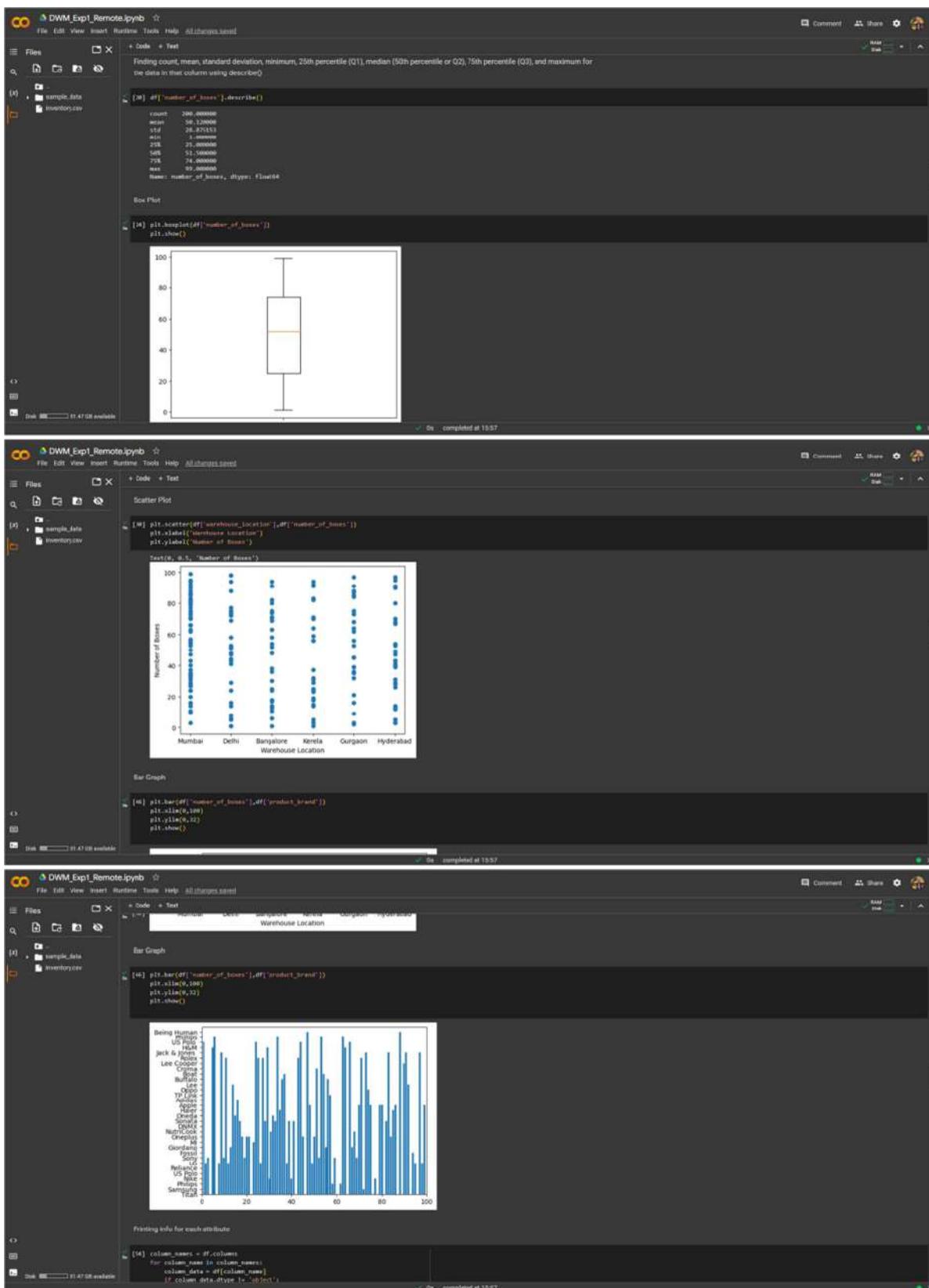
transaction_id    int64
product_id       object
product_name      object
product_category   object
product_brand      object
warehouse_id      object
warehouse_name     object
warehouse_location  object
city_id           object
supplier_id       object
supplier_name      object
contact_id        object
phone             object
email            object
arrival_date      object
dispatch_date     object
number_of_boxes     int64
dtype: object

Finding Mean, Median, Mode for 'number_of_boxes'
[10] meanCalc=df['number_of_Boxes'].mean()
medianCalc=df['number_of_Boxes'].median()
modeCalc=df['number_of_Boxes'].mode()
print("Mean = ", meanCalc)
print("Median = ", medianCalc)
print("Mode = ", modeCalc)
print("Mode = ", modeCalc[0])

Mean = 50.13
Median = 51.5
Mode =
0: 24
1: 2
2: 91
Mode: number_of_Boxes, dtype: int64
0s completed at 15:57

```

C32_Om_2103163_DWM



DWM_Expt1_Remote.ipynb

```

File Edit View Insert Runtime Tools Help
File File Editor + Code + Test
Print info for each numeric attribute
for column_name in df.columns:
    if column_name in column_names:
        column_data = df[column_name]
        if column_data.dtype != <class 'object'>:
            print("Attribute Name: " + column_name)
            count = column_data.count()
            print("Count of Values: " + str(count))
            min_val = column_data.min()
            print("Minimum Value: " + str(min_val))
            max_val = column_data.max()
            print("Maximum Value: " + str(max_val))
            quartiles = column_data.quantile([0.25, 0.5, 0.75])
            data_type = column_data.dtype
            print("Data Type: " + str(data_type))
            range_val = max_val - min_val
            print("Range: " + str(range_val))
            outliers = column_data[(column_data <= min_val) | (column_data >= max_val)]
            print("Outliers: " + str(outliers))
            print("Number of Outliers: " + str(len(outliers)))
            print("Outlier Range: " + str(outliers.min()) + " - " + str(outliers.max()))
            print("Box Plot for " + column_name)
            plt.figure(figsize=(8, 6))
            plt.boxplot(column_data, vert=True)
            plt.title(column_name)
            plt.show()
    else:
        print("Attribute Name: " + column_name)
        count = column_data.count()
        min_val = column_data.min()
        max_val = column_data.max()
        data_type = column_data.dtype
        print("Count of Values: " + str(count))
        print("Minimum Value: " + str(min_val))
        print("Maximum Value: " + str(max_val))
        print("Data Type: " + str(data_type))
        print("Range: " + str(max_val - min_val))
        print("Outliers: " + str(column_data[(column_data <= min_val) | (column_data >= max_val)]))
        print("Number of Outliers: " + str(len(column_data[(column_data <= min_val) | (column_data >= max_val)])))
        print("Outlier Range: " + str(column_data[(column_data <= min_val) | (column_data >= max_val)].min()) + " - " + str(column_data[(column_data <= min_val) | (column_data >= max_val)].max())))
        print("Box Plot for " + column_name)
        plt.figure(figsize=(8, 6))
        plt.boxplot(column_data, vert=True)
        plt.title(column_name)
        plt.show()

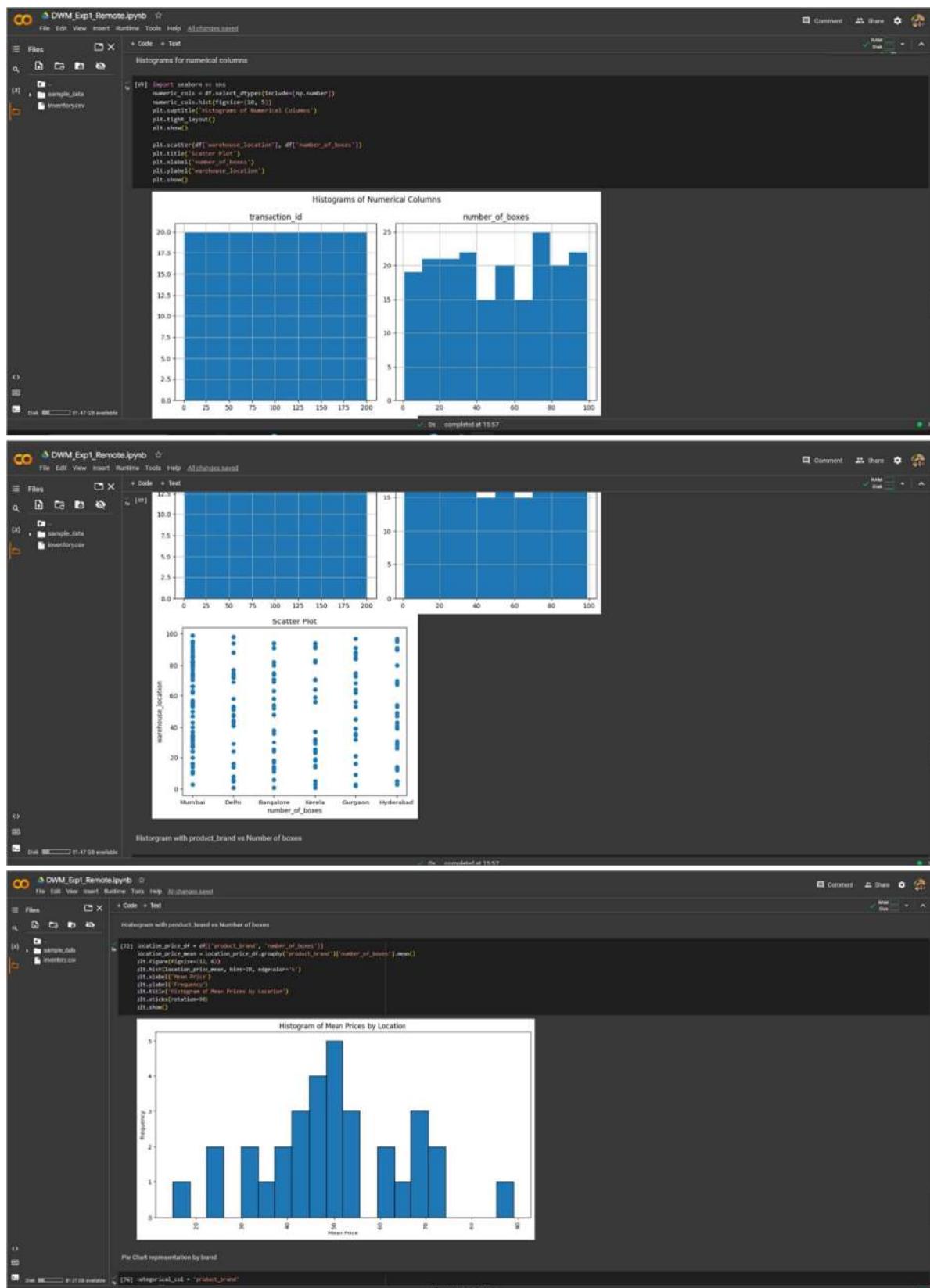
```

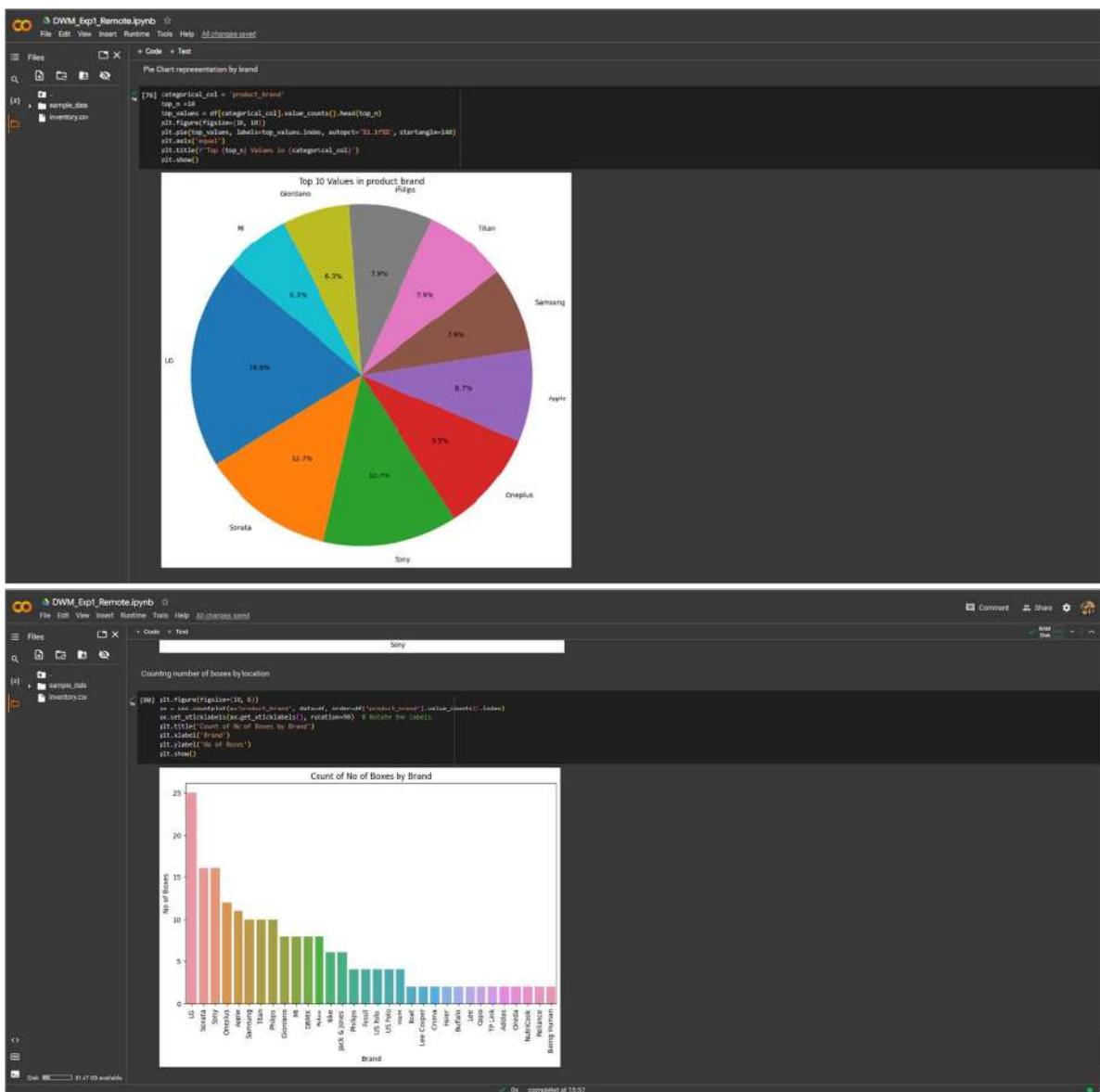
Attribute Name: transaction_id
Count of Values: 200
Minimum Value: 3
Maximum Value: 199
Data Type: int64
Range: 196
Quantiles:
0.25 50.75
0.50 100.50
0.75 150.25
Name: transaction_id, dtype: float64

Box Plot for transaction_id

Attribute Name: number_of_boxes
Count of Values: 200
Minimum Value: 3
Maximum Value: 99
Data Type: int64
Range: 96
Quantiles:
0.25 25.0
0.50 32.5
0.75 74.0
Name: number_of_boxes, dtype: float64

Box Plot for number_of_boxes





EXPERIMENT - 02

Aim: One case study on building Data warehouse/Data Mart.

Theory:

Difference between the two terms – DBMS & Data Warehouse

A Database Management System (DBMS) stores data in the form of tables, uses ER model and the goal is ACID properties. For example, a DBMS of college has tables for students, faculty, etc.

A Data Warehouse is separate from DBMS, it stores a huge amount of data, which is typically collected from multiple heterogeneous sources like files, DBMS, etc. The goal is to produce statistical results that may help in decision makings. For example, a college might want to see quick different results, like how the placement of CS students has improved over the last 10 years, in terms of salaries, counts, etc.

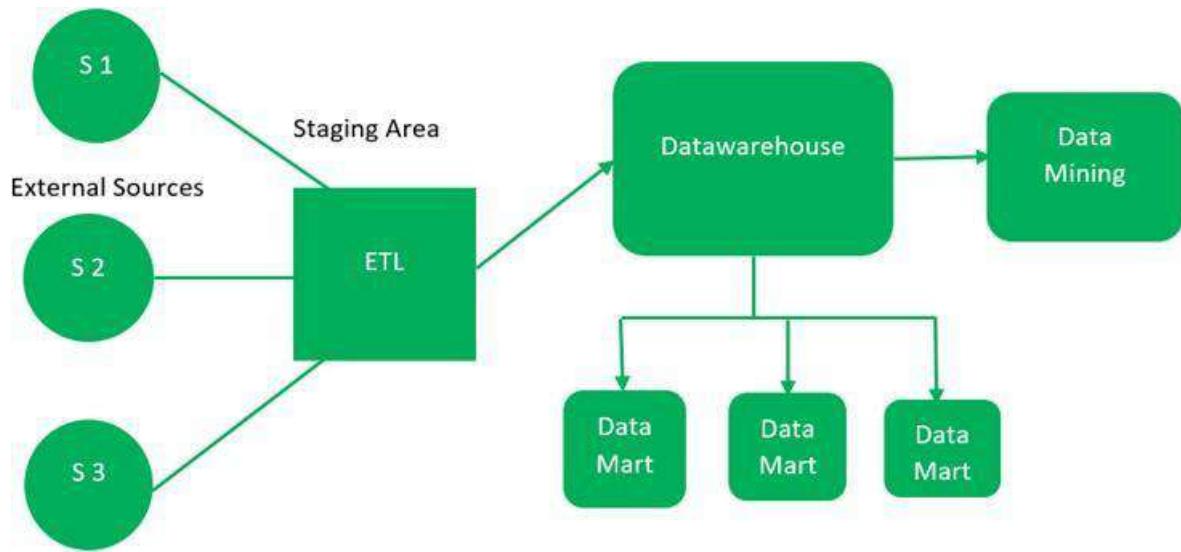
Benefits of Data Warehouse

1. Better business analytics: Data warehouse plays an important role in every business to store and analysis of all the past data and records of the company. which can further increase the understanding or analysis of data to the company.
2. Faster Queries: Data warehouse is designed to handle large queries that's why it runs queries faster than the database.
3. Improved data Quality: In the data warehouse the data you gathered from different sources is being stored and analyzed it does not interfere with or add data by itself so your quality of data is maintained and if you get any issue regarding data quality then the data warehouse team will solve this.
4. Historical Insight: The warehouse stores all your historical data which contains details about the business so that one can analyze it at any time and extract insights from it

Construction of a Data Warehouse

There are 2 approaches for constructing data-warehouse: Top-down approach and Bottom-up approach are explained as below.

A. Top-Down Approach:



The essential components are discussed below:

1. External Sources –

External source is a source from where data is collected irrespective of the type of data. Data can be structured, semi structured and unstructured as well.

2. Stage Area –

Since the data, extracted from the external sources does not follow a particular format, so there is a need to validate this data to load into data warehouse. For this purpose, it is recommended to use ETL tool.

3. E(Extracted): Data is extracted from External data source.

4. T(Transform): Data is transformed into the standard format.

5. L(Load): Data is loaded into data warehouse after transforming it into the standard format.

6. Data-warehouse –

After cleansing of data, it is stored in the data warehouse as central repository. It actually stores the meta data and the actual data gets stored in the data marts. Note that data warehouse stores the data in its purest form in this top-down approach.

7. Data Marts –

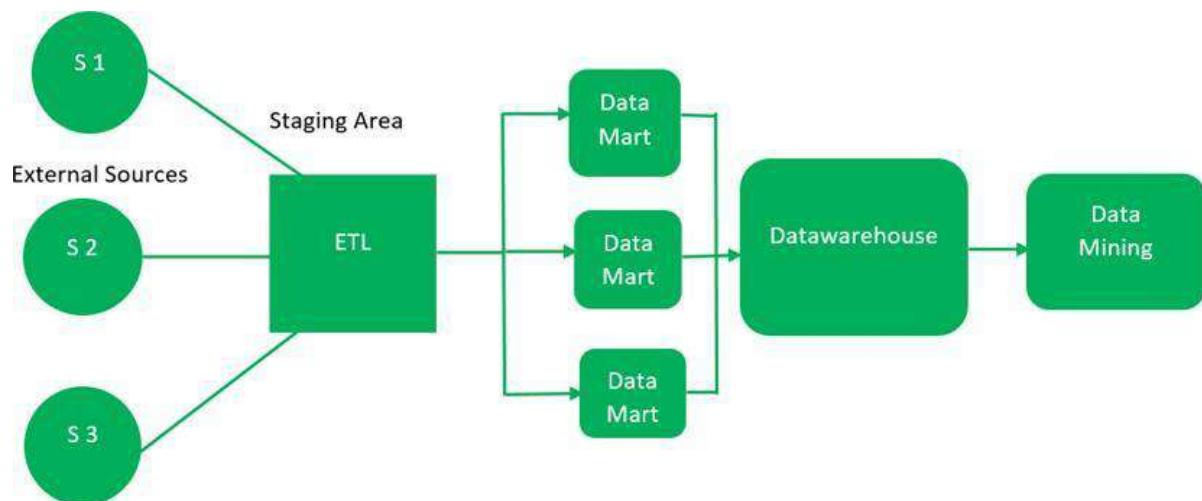
Data mart is also a part of storage component. It stores the information of a particular function of an organisation which is handled by single authority. There can be as many number of data marts in an organisation depending upon the functions. We can also say that data mart contains subset of the data stored in data warehouse.

8. Data Mining –

The practice of analysing the big data present in data warehouse is data mining. It is used to find the hidden patterns that are present in the database or in data warehouse with the help of algorithm of data mining.

This approach is defined by Inmon as – data warehouse as a central repository for the complete organisation and data marts are created from it after the complete data warehouse has been created.

B. Bottom-Up Approach:



1. First, the data is extracted from external sources (same as happens in top-down approach).
2. Then, the data go through the staging area (as explained above) and loaded into data marts instead of data warehouse. The data marts are created first and provide reporting capability. It addresses a single business area.
3. These data marts are then integrated into data warehouse.

This approach is given by Kinball as – data marts are created first and provides a thin view for analyses and data warehouse is created after complete data marts have been created.

Fact Table

- In a data warehouse, a fact table is a table that stores the measurements, metrics, or facts related to a business operation.
- It is located at the centre of a star or snowflake schema and is surrounded by dimension tables.
- When multiple fact tables are used, they can be organized using a "fact constellation schema."
- A fact table has two types of columns: those that contain the facts and those that serve as foreign keys linking to dimension tables.
- The primary key of a fact table is often a composite key made up of all of the foreign keys in the table.
- Fact tables can hold various types of measurements, such as additive, non-additive, and partly additive measures, and store important information in the data warehouse.
- They are useful for evaluating dimensional attributes because they provide additive values that can act as independent variables.

Dimension Table

- Dimension tables contain descriptions of the objects in a fact table and provide information about dimensions such as values, characteristics, and keys.
- These tables are usually small, with a number of rows ranging from a few hundred to a few thousand.
- The term "dimension table" refers to a set of data related to any quantifiable event and is the foundation for dimensional modelling.

- Dimension tables have a column that serves as a primary key, allowing each dimension row or record to be uniquely identified. This key is used to link the dimension table to the fact tables. A surrogate key, which is a system-generated key, is often used to uniquely identify the rows in the dimension table.

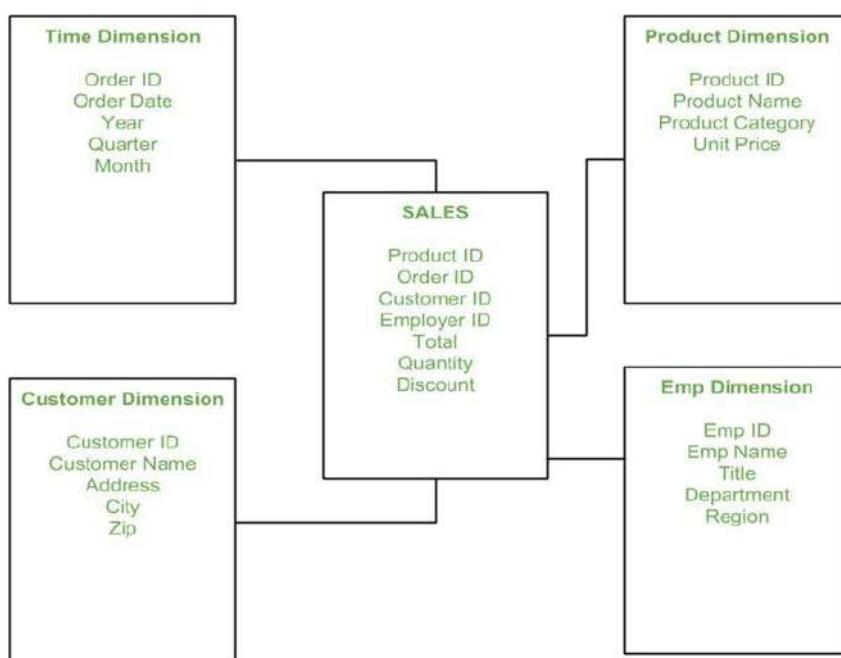
Star Schema

A star schema is a type of data modelling technique used in data warehousing to represent data in a structured and intuitive way. In a star schema, data is organized into a central fact table that contains the measures of interest, surrounded by dimension tables that describe the attributes of the measures.

The fact table in a star schema contains the measures or metrics that are of interest to the user or organization. For example, in a sales data warehouse, the fact table might contain sales revenue, units sold, and profit margins. Each record in the fact table represents a specific event or transaction, such as a sale or order.

The dimension tables in a star schema contain the descriptive attributes of the measures in the fact table. These attributes are used to slice and dice the data in the fact table, allowing users to analyse the data from different perspectives. For example, in a sales data warehouse, the dimension tables might include product, customer, time, and location.

In a star schema, each dimension table is joined to the fact table through a foreign key relationship. This allows users to query the data in the fact table using attributes from the dimension tables.



Advantages of Star Schema

1. Simpler Queries – Join logic of star schema is quite cinch in comparison to other join logic which are needed to fetch data from a transactional schema that is highly normalized.
2. Simplified Business Reporting Logic – In comparison to a transactional schema that is highly normalized, the star schema makes simpler common business reporting logic, such as of reporting and period-over-period.
3. Feeding Cubes – Star schema is widely used by all OLAP systems to design OLAP cubes efficiently. In fact, major OLAP systems deliver a ROLAP mode of operation which can use a star schema as a source without designing a cube structure.

Disadvantages of Star Schema

1. Data integrity is not enforced well since in a highly de-normalized schema state.
2. Not flexible in terms of analytical needs as a normalized data model.
3. Star schemas don't reinforce many-to-many relationships within business entities – at least not frequently.

Snowflake Schema

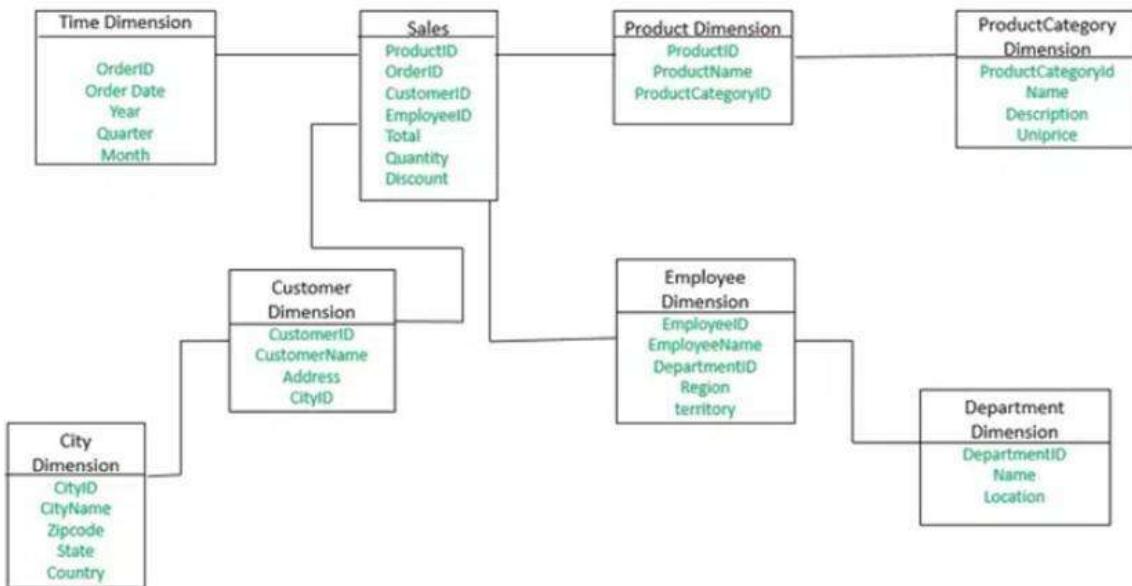
The snowflake schema is a variant of the star schema. Here, the centralized fact table is connected to multiple dimensions. In the snowflake schema, dimensions are present in a normalized form in multiple related tables. The snowflake structure materializes when the dimensions of a star schema are detailed and highly structured, having several levels of relationship, and the child tables have multiple parent tables. The snowflake effect affects only the dimension tables and does not affect the fact tables.

In a snowflake schema, the dimension tables are normalized into multiple related tables, creating a hierarchical or “snowflake” structure.

In a snowflake schema, the fact table is still located at the centre of the schema, surrounded by the dimension tables. However, each dimension table is further broken down into multiple related tables, creating a hierarchical structure that resembles a snowflake.

The Employee dimension table now contains the attributes: EmployeeID, EmployeeName, DepartmentID, Region, and Territory. The DepartmentID attribute links with the Employee table with the Department dimension table. The Department dimension is used to provide detail about each department, such as the Name and Location of the department.

The Customer dimension table now contains the attributes: CustomerID, CustomerName, Address, and CityID. The CityID attributes link the Customer dimension table with the City dimension table. The City dimension table has details about each city such as city name, Zipcode, State, and Country.



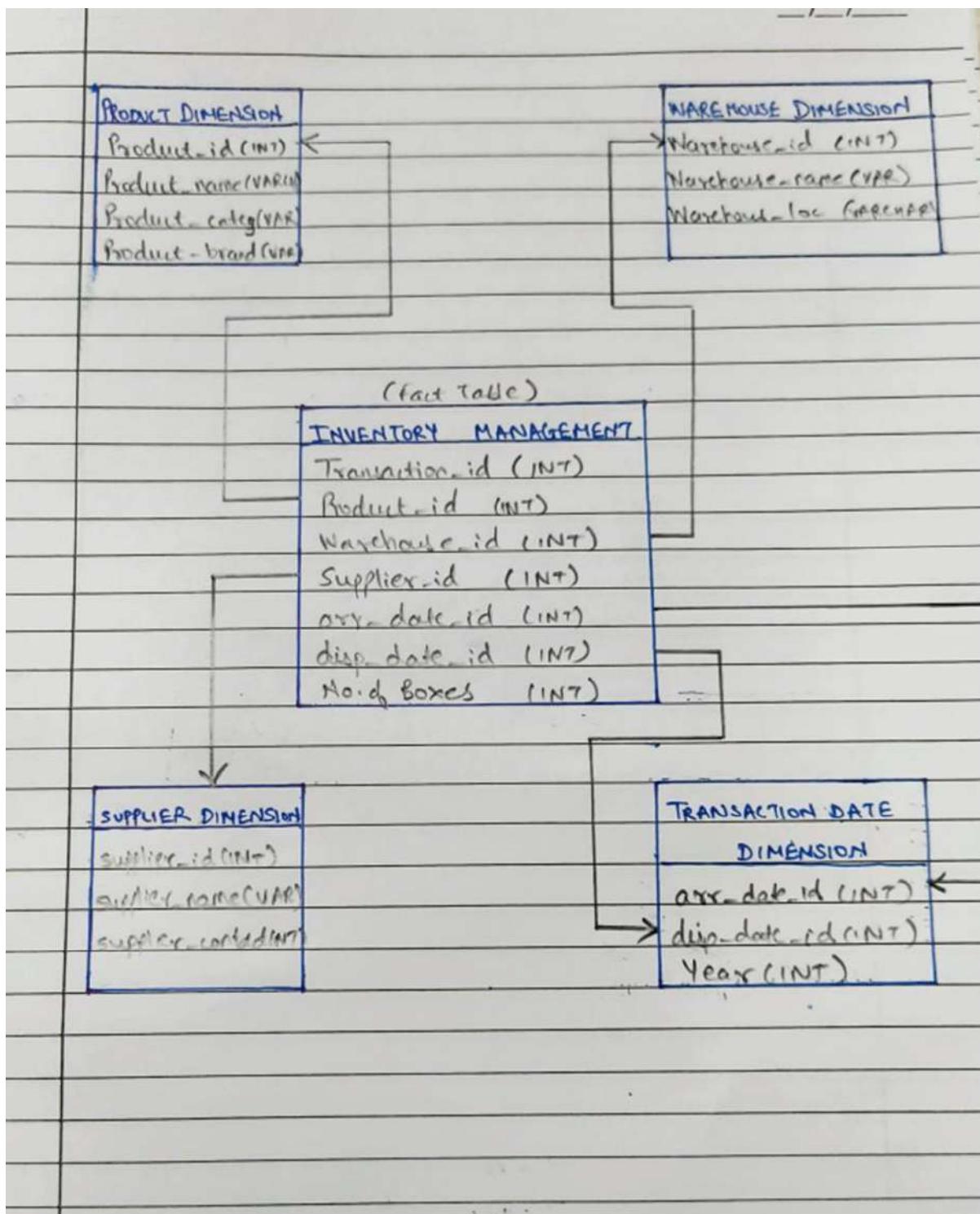
Advantages of Snowflake Schema

1. It provides structured data which reduces the problem of data integrity.
2. It uses small disk space because data are highly structured.

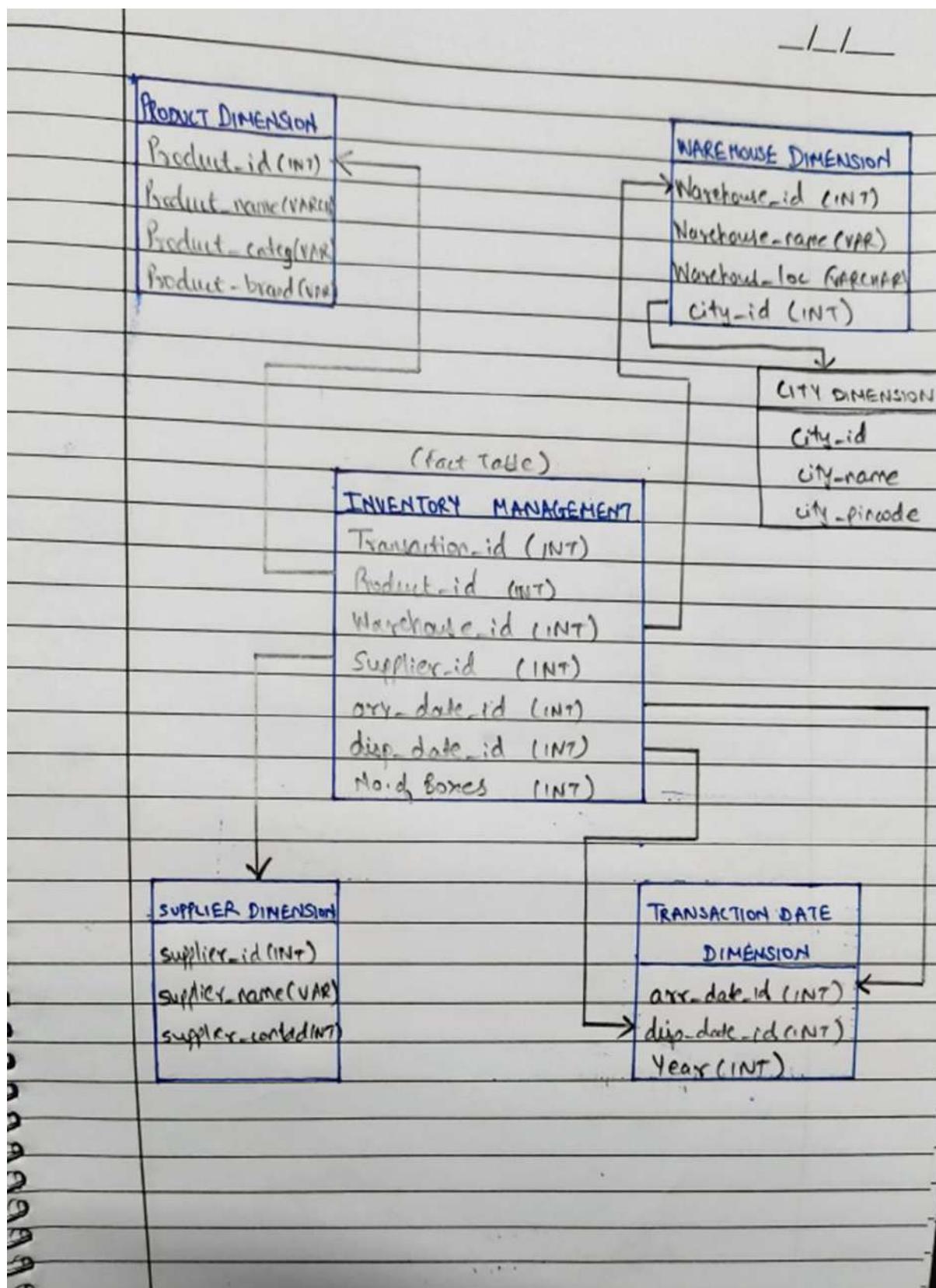
Disadvantages of Snowflake Schema

1. Snowflaking reduces space consumed by dimension tables but compared with the entire data warehouse the saving is usually insignificant.
2. Avoid snowflaking or normalization of a dimension table, unless required and appropriate.
3. Do not snowflake hierarchies of dimension table into separate tables. Hierarchies should belong to the dimension table only and should never be snowflakes.
4. Multiple hierarchies that can belong to the same dimension have been designed at the lowest possible detail.

Star Schema



Snowflake Schema



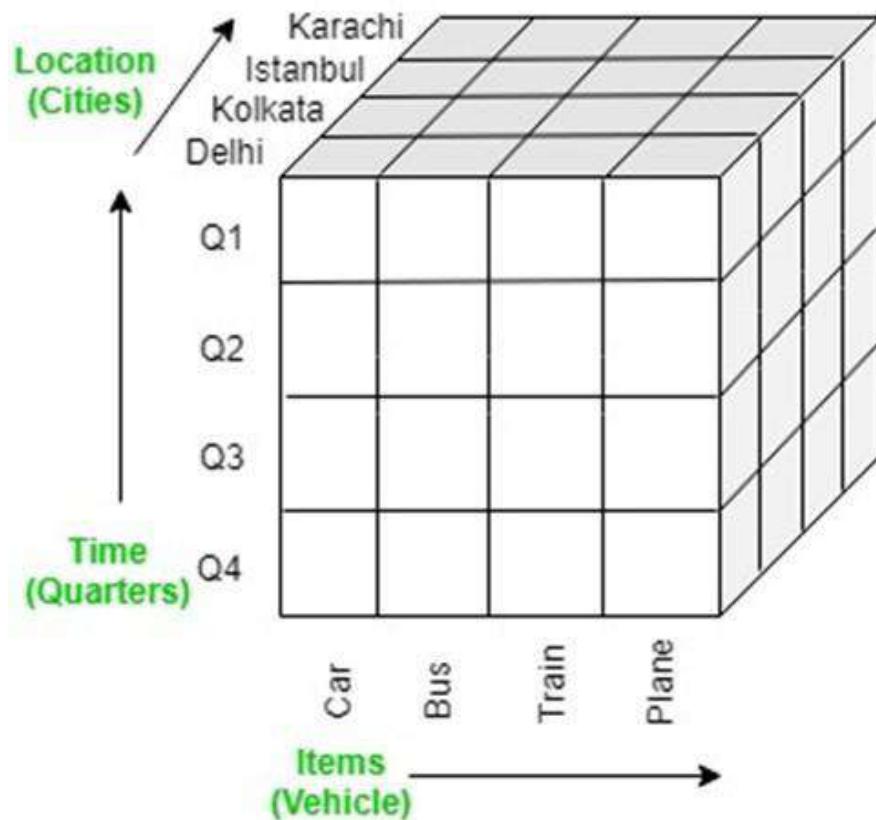
EXPERIMENT - 03

Aim: Implementation of all dimension tables and fact tables and OLAP (in Oracle)

Theory:

In the multidimensional model, the records are organized into various dimensions, and each dimension includes multiple levels of abstraction described by concept hierarchies. This organization supports users with the flexibility to view data from various perspectives. Several OLAP data cube operations exist to demonstrate these different views, allowing interactive queries and a search of the record at hand. Hence, OLAP supports a user-friendly environment for interactive data analysis.

Consider the OLAP operations which are to be performed on multidimensional data. The figure shows data cubes for sales of a shop. The cube contains the dimensions, location, time and item, where the location is aggregated about city values, time is aggregated concerning quarters, and an item is aggregated concerning item types.



OLAP operations:

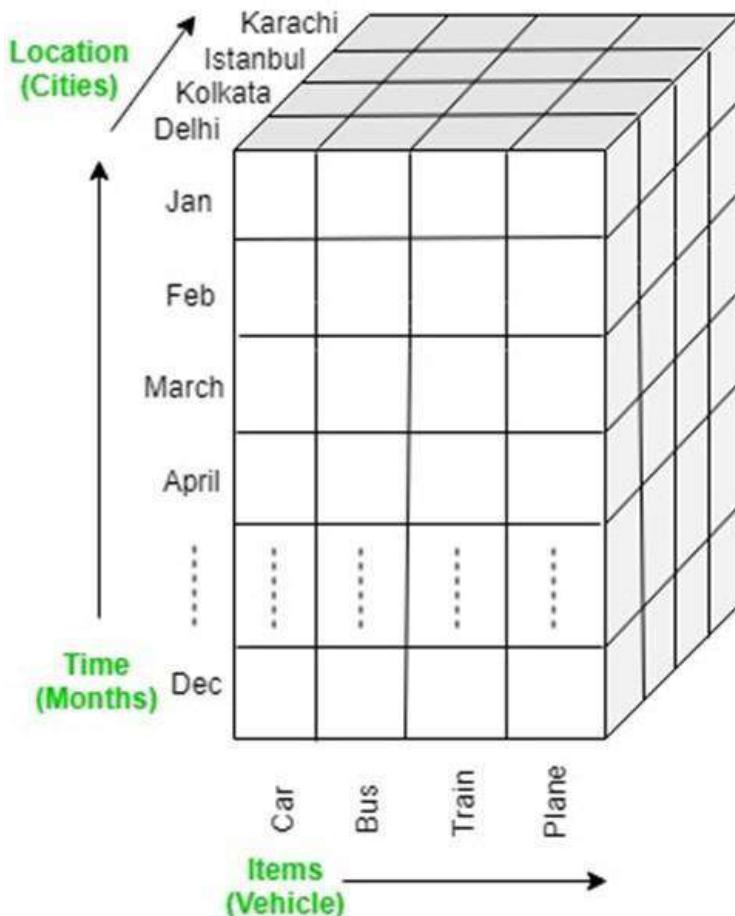
Five basic analytical operations can be performed on an OLAP cube:

1. Drill down:

In drill-down operation, the less detailed data is converted into highly detailed data. It can be done by:

- Moving down in the concept hierarchy
- Adding a new dimension

In the cube given in the overview section, the drill-down operation is performed by moving down in the concept hierarchy of the Time dimension (Quarter -> Month).

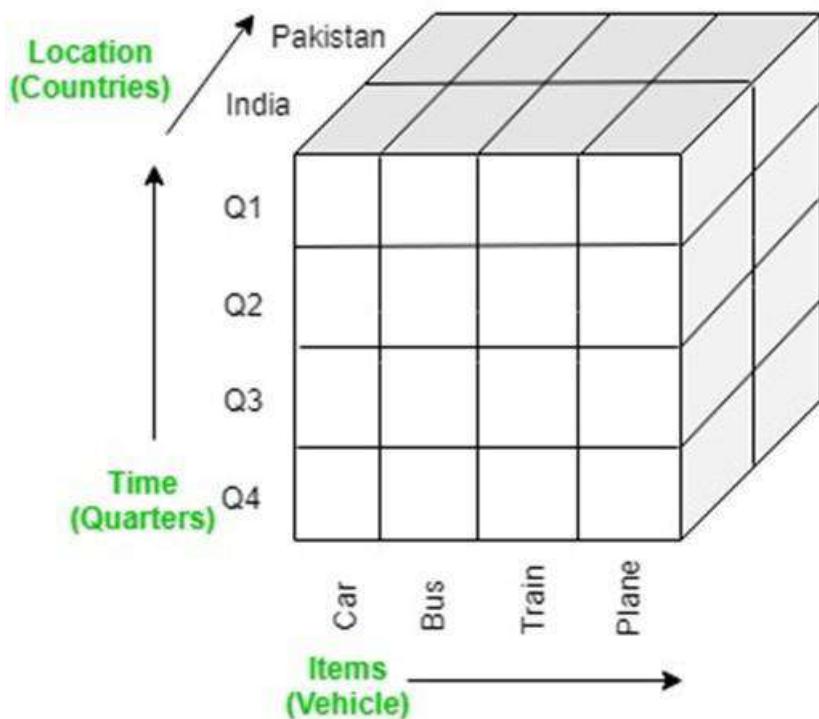


2. Roll up:

It is just the opposite of the drill-down operation. It performs aggregation on the OLAP cube. It can be done by:

- Climbing up in the concept hierarchy
- Reducing the dimensions

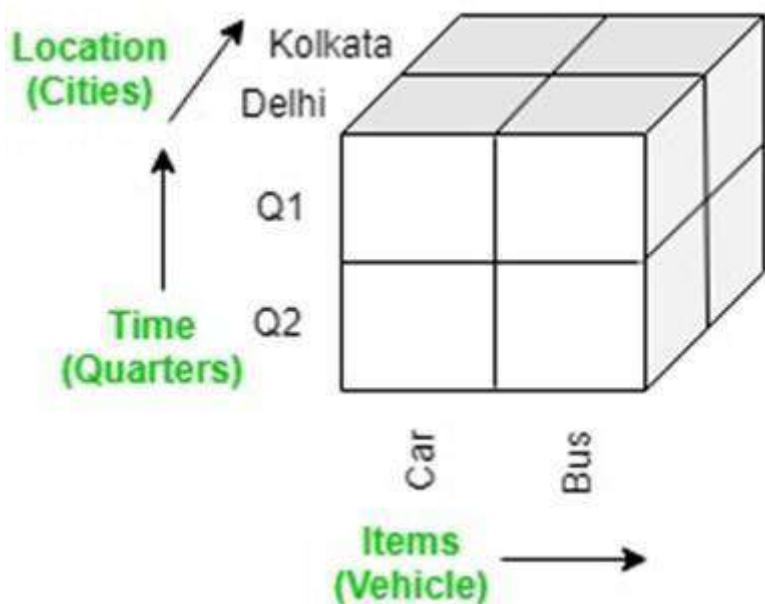
In the cube given in the overview section, the roll-up operation is performed by climbing up in the concept hierarchy of Location dimension (City -> Country).



3. Dice:

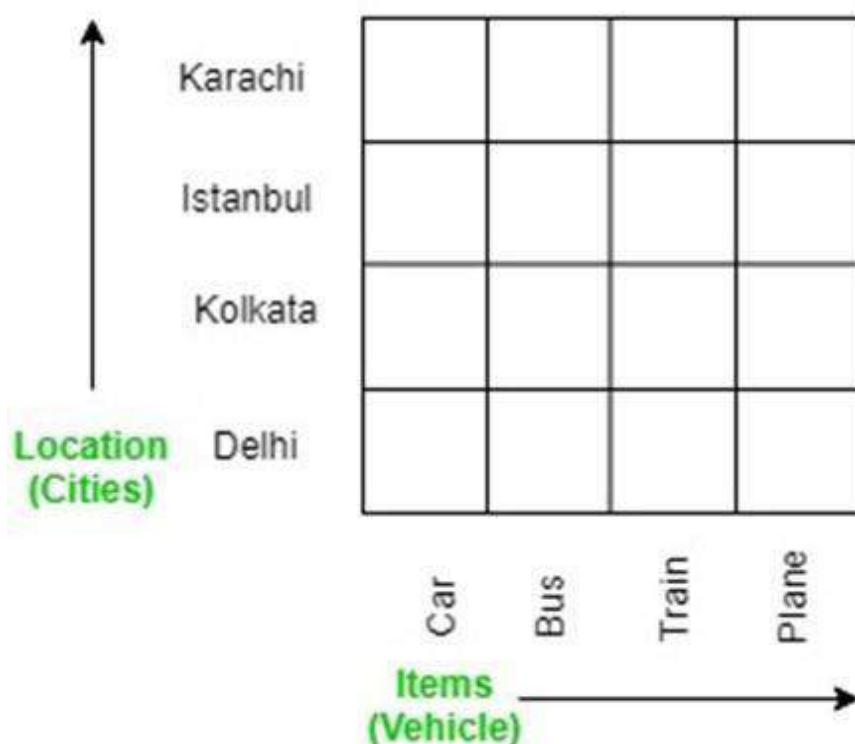
It selects a sub-cube from the OLAP cube by selecting two or more dimensions. In the cube given in the overview section, a sub-cube is selected by selecting the following dimensions with criteria:

- Location = “Delhi” or “Kolkata”
- Time = “Q1” or “Q2”
- Item = “Car” or “Bus”



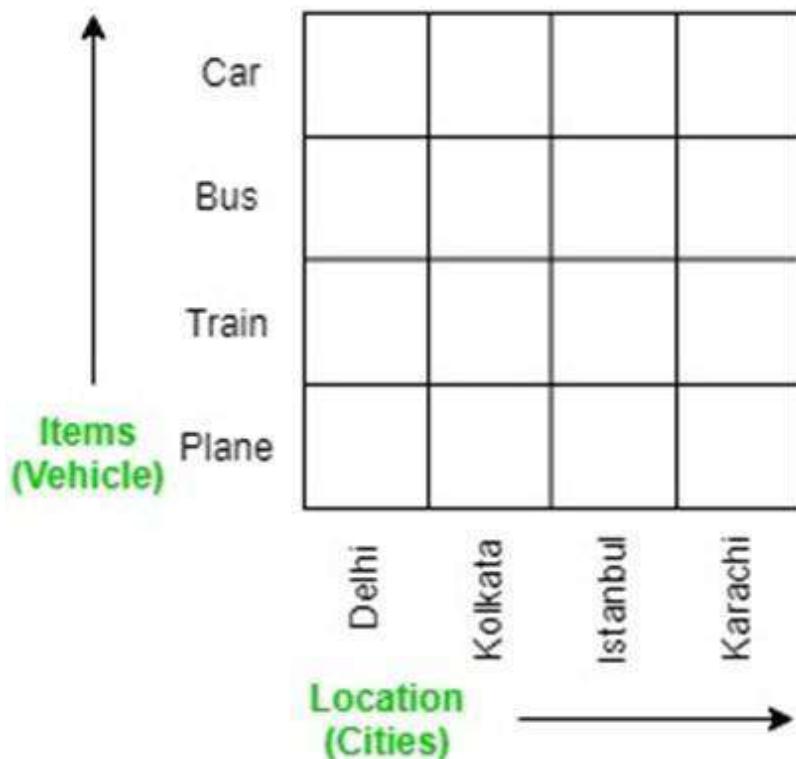
4. Slice:

It selects a single dimension from the OLAP cube which results in a new sub-cube creation. In the cube given in the overview section, Slice is performed on the dimension Time = "Q1".



5. Pivot:

It is also known as rotation operation as it rotates the current view to get a new view of the representation. In the sub-cube obtained after the slice operation, performing the pivot operation gives a new view of it.



Code:

Creating Fact Table and Dimension tables.

DIMENSIONS TABLE – CITY DIMENSIONS

```
CREATE TABLE city_dim (
    city_id INT PRIMARY KEY,
    city_name VARCHAR(255),
    c_state VARCHAR(255),
    pincode VARCHAR(10)
);
```

DIMENSIONS TABLE – CONTACT DIMENSIONS

```
CREATE TABLE contact_dim (
    contact_id VARCHAR(255) PRIMARY
    KEY,
    phone VARCHAR(15),
    email VARCHAR(255)
);
```

DIMENSIONS TABLE – SUPPLIER DIMENSIONS

```
CREATE TABLE supplier_dim (
    supplier_id VARCHAR(255) PRIMARY
    KEY,
    supplier_name VARCHAR(255),
    contact_id INT
);
```

DIMENSIONS TABLE – PRODUCT DIMENSIONS

```
CREATE TABLE product_dim (
    product_id VARCHAR(255) PRIMARY KEY,
    product_name VARCHAR(255),
    product_category VARCHAR(255),
    product_brand VARCHAR(255)
);
```

DIMENSIONS TABLE – TRANSACTION_DATE DIMENSIONS

```
CREATE TABLE transaction_date_dim (
    arrival_date DATE,
    dispatch_date DATE,
    duration_in_warehouse INT
);
```

DIMENSIONS TABLE – WAREHOUSE DIMENSIONS

```
CREATE TABLE warehouse_dim (
    warehouse_id VARCHAR(255) PRIMARY KEY,
    warehouse_name VARCHAR(255),
    warehouse_location VARCHAR(255),
    city_id INT
);
```

FACT TABLE - INVENTORY

```
CREATE TABLE inventory_management_fact (
    transaction_id INT PRIMARY KEY,
    product_id VARCHAR(255),
    warehouse_id VARCHAR(255),
    quantity INT
);
```

```
supplier_id  
VARCHAR(255),  
arrival_date_id DATE,  
dispatch_date_id DATE,  
number_of_boxes INT,  
FOREIGN KEY (product_id) REFERENCES product_dim(product_id),  
FOREIGN KEY (warehouse_id) REFERENCES  
warehouse_dim(warehouse_id), FOREIGN KEY (supplier_id) REFERENCES  
supplier_dim(supplier_id)  
);
```

Inserting Values in the above tables:

city_dimension:

```
INSERT INTO city_dim (city_id, city_name, c_state,  
pincode) VALUES (7597, 'Mumbai', 'Maharashtra',  
'400001');  
  
INSERT INTO city_dim (city_id, city_name, c_state,  
pincode) VALUES (7099, 'Bangalore', 'Karnataka', '560001');  
  
INSERT INTO city_dim (city_id, city_name, c_state,  
pincode) VALUES (9293, 'Pune', 'Maharashtra', '411001');  
  
INSERT INTO city_dim (city_id, city_name, c_state,  
pincode) VALUES (6336, 'Hyderabad', 'Telangana',  
'500001');  
  
INSERT INTO city_dim (city_id, city_name, c_state,
```

```
pincode) VALUES (5676, 'Gurgaon', 'Haryana', '122001');

INSERT INTO city_dim (city_id, city_name, c_state, pincode)
VALUES (5055, 'Delhi', 'Delhi', '110001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (7298, 'Kota', 'Rajasthan', '324001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (6244, 'Jaipur', 'Rajasthan', '302001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (4098, 'Chennai', 'Tamil Nadu', '600001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (7544, 'Kolkata', 'West Bengal', '700001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (6003, 'Ahmedabad', 'Gujarat', '380001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (8081, 'Lucknow', 'Uttar Pradesh',
'226001'); INSERT INTO city_dim (city_id, city_name,
c_state, pincode) VALUES (5210, 'Chandigarh',
'Chandigarh', '160001'); INSERT INTO city_dim (city_id,
city_name, c_state, pincode) VALUES (7142, 'Bhopal',
'Madhya Pradesh', '462001'); INSERT INTO city_dim
(city_id, city_name, c_state, pincode) VALUES (6001,
'Patna', 'Bihar', '800001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (9012, 'Guwahati', 'Assam', '781001');

INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (7659, 'Trivandrum', 'Kerala', '695001');
```

```
INSERT INTO city_dim (city_id, city_name, c_state, pincode)
VALUES (7245, 'Bhubaneswar', 'Odisha', '751001');
```

```
INSERT INTO city_dim (city_id, city_name, c_state,
pincode) VALUES (5222, 'Dehradun', 'Uttarakhand',
'248001'); INSERT INTO city_dim (city_id, city_name,
c_state, pincode) VALUES (6211, 'Shimla', 'Himachal
```

The screenshot shows a SQL worksheet interface with the following details:

- SQL Worksheet:** The main area displays the following SQL code:


```
40 VALUES (6211, 'Shimla', 'Himachal Pradesh', '171001');
41
42 SELECT * FROM city_dim;
```
- My Session:** A sidebar on the left lists various sections: Home, SQL Worksheet, My Session, Schema, Quick SQL, My Scripts, My Tutorials, and Code Library.
- Schema:** A table view showing the data inserted into the city_dim table:

CITY_ID	CITY_NAME	C_STATE	PINCODE
7597	Mumbai	Maharashtra	400001
7099	Bangalore	Karnataka	560001
9293	Pune	Maharashtra	411001
6336	Hyderabad	Telangana	500001
5676	Gurgaon	Haryana	122001
5055	Delhi	Delhi	110001
7298	Kota	Rajasthan	324001
6244	Jaipur	Rajasthan	302001
4898	Chennai	Tamil Nadu	600001
7544	Kolkata	West Bengal	700001

Pradesh', '171001');

contact_dim:

```
INSERT INTO contact_dim (contact_id, phone, email)
VALUES ('C2483', '9988776655',
'xbarrett@example.com'); INSERT INTO contact_dim
(contact_id, phone, email)
VALUES ('C8454', '9876543210',
'vernandez@example.com'); INSERT INTO contact_dim
(contact_id, phone, email)
VALUES ('C8265', '8899776644',
'cooperdavid@example.org'); INSERT INTO contact_dim
(contact_id, phone, email)
VALUES ('C3118', '7766554433',
```

```
'monica66@example.net'); INSERT INTO contact_dim  
(contact_id, phone, email) VALUES ('C7283', '9870123456',  
'krivera@example.org'); INSERT INTO contact_dim  
(contact_id, phone, email)  
VALUES ('C7667', '9876012345',  
'rosarioalbert@example.net'); INSERT INTO contact_dim  
(contact_id, phone, email)  
VALUES ('C1342', '8765432109', 'leejennifer@example.net');  
INSERT INTO contact_dim (contact_id, phone, email)  
VALUES ('C3989', '8877665544',  
'breynolds@example.com'); INSERT INTO contact_dim  
(contact_id, phone, email)  
VALUES ('C1699', '9988001122',  
'thompsonzachary@example.org'); INSERT INTO contact_dim  
(contact_id, phone, email)  
VALUES ('C5700', '9012345678',  
'robert84@example.org'); INSERT INTO contact_dim  
(contact_id, phone, email)  
VALUES ('C9238', '9878998899',  
'hoffmanbeth@example.net'); INSERT INTO contact_dim  
(contact_id, phone, email)  
VALUES ('C7601', '9988778899', 'sdavis@example.net');  
INSERT INTO contact_dim (contact_id, phone, email)  
VALUES ('C5895', '9123456789', 'kristina32@example.net');  
INSERT INTO contact_dim (contact_id, phone, email)  
VALUES ('C9589', '8899889988', 'natalie47@example.org');
```

```

INSERT INTO contact_dim (contact_id, phone, email)

VALUES ('C1627', '9900990099',

'longashley@example.org'); INSERT INTO contact_dim

(contact_id, phone, email)

VALUES ('C5607', '8765432101',

'coreywilson@example.org'); INSERT INTO contact_dim

(contact_id, phone, email)

VALUES ('C9300', '8901234567',

'ebryan@example.com'); INSERT INTO contact_dim

(contact_id, phone, email)

VALUES ('C8900', '9876543211',

'burkesharon@example.net'); INSERT INTO contact_dim

(contact_id, phone, email)

VALUES ('C2870', '8877665599',

'deborah09@example.org'); INSERT INTO contact_dim

(contact_id, phone, email)

VALUES ('C8997', '9876543232', 'dustinpeterson@example.org');

```

The screenshot shows the Oracle Live SQL interface. On the left is a sidebar with navigation links: Home, SQL Worksheet (selected), My Session, Schema, Quick SQL, My Scripts, My Tutorials, and Code Library. The main area is titled "SQL Worksheet". It contains the following SQL code:

```

37 ✓ INSERT INTO contact_dim (contact_id, phone, email)
38 VALUES ('C2870', '8877665599', 'deborah09@example.org');
39 ✓ INSERT INTO contact_dim (contact_id, phone, email)
40 VALUES ('C8997', '9876543232', 'dustinpeterson@example.org');
41
42 SELECT * FROM contact_dim
43

```

Below the code, a table displays the data inserted into the contact_dim table:

CONTACT_ID	PHONE	EMAIL
C2870	8877665599	deborah09@example.org
C8997	9876543232	dustinpeterson@example.org
C3483	9988776655	xbarrett@example.com
C0454	9876543210	vhernandez@example.com
C8265	8899776644	cooperdavid@example.org
C3118	7766554433	monica66@example.net
C7283	9870123456	krivera@example.org
C6667	9876012345	rosarioalbert@example.net
C1342	8765432109	leejennifer@example.net
C3989	8877665544	breyynolds@example.com

At the bottom of the interface, there is footer text: "2023 Oracle - Live SQL 23.4.1, running Oracle Database 19c EE Extreme Perf - 19.17.0.0.0 - Database Documentation - Ask Tom - Dev Gym" and "Built with ❤ using Oracle APEX - Privacy · Terms of Use".

supplier_dimension:

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S3535', 'Stewart, Clayton and Martinez',  
'C2483');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S7973', 'Carroll, Brady and Hancock', 'C8454');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S2142', 'Singh-Johnson', 'C8265');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S3115', 'Boyer Ltd', 'C3118');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name, contact_id)  
VALUES ('S8922', 'Davis Ltd', 'C7283');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S4540', 'Sandoval, Brown and Peters',  
'C7667');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S2603', 'Madden Inc', 'C1342');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S3568', 'Pearson, Wall and Shelton', 'C3989');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S8035', 'Ellis Ltd', 'C1699');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S3141', 'Johnson Group', 'C5700');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,  
contact_id) VALUES ('S8420', 'Mccoy LLC', 'C9238');
```

```
INSERT INTO supplier_dim (supplier_id, supplier_name,
```

```

contact_id) VALUES ('S1561', 'Smith Ltd', 'C7601');

INSERT INTO supplier_dim (supplier_id, supplier_name,
                           contact_id) VALUES ('S8841', 'Small PLC', 'C5895');

INSERT INTO supplier_dim (supplier_id, supplier_name,
                           contact_id) VALUES ('S3895', 'Carney, Edwards and Chen', 'C9589');

INSERT INTO supplier_dim (supplier_id, supplier_name,
                           contact_id) VALUES ('S9366', 'Ruiz-Garcia', 'C1627');

INSERT INTO supplier_dim (supplier_id, supplier_name,
                           contact_id) VALUES ('S3534', 'Potter and Sons', 'C5607');

INSERT INTO supplier_dim (supplier_id, supplier_name, contact_id)
VALUES ('S7178', 'Rowe-Benitez', 'C9300');

INSERT INTO supplier_dim (supplier_id, supplier_name,
                           contact_id) VALUES ('S6126', 'Williamson, Ware and Martin',
                           'C8900');

INSERT INTO supplier_dim (supplier_id, supplier_name,
                           contact_id) VALUES ('S3284', 'Price, Sparks and Frey', 'C2870');

INSERT INTO supplier_dim (supplier_id, supplier_name,
                           contact_id) VALUES ('S9796', 'Newton Group', 'C8997');

```

The screenshot shows the Oracle SQL Developer interface with a "Live SQL" tab selected. The left sidebar includes sections for Home, SQL Worksheet (selected), My Session, Schema, Quick SQL, My Scripts, My Tutorials, and Code Library. The main area displays a SQL Worksheet with the following code and results:

```

SQL Worksheet
56    SELECT * FROM supplier_dim
57
58    INSERT INTO supplier_dim (supplier_id, supplier_name, contact_id)
59    VALUES ('S9796', 'Newton Group', 'C8997');
60
61    SELECT * FROM supplier_dim
62

```

The results of the last query show the following data:

SUPPLIER_ID	SUPPLIER_NAME	CONTACT_ID
S3535	Stewart, Clayton and Martinez	C2483
S7973	Carroll, Brady and Hancock	C8454
S2142	Singh-Johnson	C8265
S3115	Boyer Ltd	C3118
S8922	Davis Ltd	C7283
S4548	Sandoval, Brown and Peters	C7667
S2603	Madden Inc	C1342
S3568	Pearson, Wall and Shelton	C3989
S8035	Ellis Ltd	C1699

At the bottom of the interface, there is a footer with the text: "2023 Oracle - Use SQL 23.4.1, running Oracle Database 19c EE Extreme Perf - 19.17.0.0.0 - Database Documentation - Ask Tom - Dev Gym".

product_dimension:

```
INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P2859', 'Titan T27', 'watches', 'Titan');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P5937', 'Samsung QLED 21', 'tv', 'Samsung');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P1553', 'Philips Iron I4', 'Electrical App', 'Philips');

INSERT INTO product_dim (product_id, product_name, product_category, product_brand)
VALUES ('P7296', 'Nike Nex', 'Wear', 'Nike');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P9813', 'US Polo Black Denim', 'Wear', 'US Polo');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P5206', 'Samsung Washing Machine', 'Electrical App',
'Samsung');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P1489', 'Reliance Hair Dryer', 'Electrical App', 'Reliance');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P9351', 'LG 42inch P32', 'tv', 'LG');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P4413', 'Sony Bravia S76', 'tv', 'Sony');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P3825', 'Fossil F31', 'watches', 'Fossil');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P7873', 'Giordano G21', 'watches', 'Giordano');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P1596', 'MI MicroLED M31', 'tv', 'MI');
```

```
INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P2693', 'Oneplus O8', 'tv', 'Oneplus');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P4110', 'NutriCook Mixer G21', 'Electrical App', 'NutriCook');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P1235', 'DNMX Shirt S21', 'Wear', 'DNMX');

INSERT INTO product_dim (product_id, product_name, product_category, product_brand)
VALUES ('P2301', 'DNMX Hoodie', 'Wear', 'DNMX');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P5882', 'Giordano G26', 'watches', 'Giordano');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P7062', 'Sonata Kid K2', 'watches', 'Sonata');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P8087', 'Oneplus O20', 'tv', 'Oneplus');

INSERT INTO product_dim (product_id, product_name, product_category,
product_brand) VALUES ('P4582', 'Oneplus O85', 'tv', 'Oneplus');
```

The screenshot shows the Oracle Live SQL interface. On the left, there's a sidebar with navigation links: Home, SQL Worksheet (selected), My Session, Schema, Quick SQL, My Scripts, My Tutorials, and Code Library. The main area is titled "SQL Worksheet". It contains the following SQL code:

```
40 VALUES ('P4582', 'Oneplus O85', 'tv', 'Oneplus');
41
42
43 SELECT * FROM product_dim
44
```

Below the code, the results of the query are displayed in a table:

PRODUCT_ID	PRODUCT_NAME	PRODUCT_CATEGORY	PRODUCT_BRAND
P2859	Titan T27	watches	Titan
P5937	Samsung QLED 21	tv	Samsung
P1553	Phillips Iron 14	Electrical App	Phillips
P7296	Nike Nex	Wear	Nike
P9813	US Polo black Denim	Wear	US Polo
P5206	Samsung Washing Machine	Electrical App	Samsung
P1489	Reliance Hair Dryer	Electrical App	Reliance
P9351	LG 42inch P32	tv	LG
P4413	Sony Bravia S26	tv	Sony

At the bottom of the interface, there's a footer with the text: "2023 Oracle - Live SQL 23.4.1, running Oracle Database 19c EE Extreme Perf - 19.17.0.0.0 - Database Documentation - Ask Tom - Dev Gym" and "Built with ❤ using Oracle APEX - Privacy · Terms of Use".

transaction_date_dimension:

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('29-11-2021', 'DD-MM-YYYY'), TO_DATE('15-01-2023',
'DD-MM-YYYY'), 413);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('18-09-2021', 'DD-MM-YYYY'), TO_DATE('26-05-2022',
'DD-MM-YYYY'), 250);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('14-05-2021', 'DD-MM-YYYY'), TO_DATE('16-05-2022',
'DD-MM-YYYY'), 367);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('21-06-2021', 'DD-MM-YYYY'), TO_DATE('13-10-2021',
'DD-MM-YYYY'), 115);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('29-10-2022', 'DD-MM-YYYY'), TO_DATE('07-11-2022', 'DD-MM-YYYY'),
9);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('04-11-2022', 'DD-MM-YYYY'), TO_DATE('17-11-2022',
'DD-MM-YYYY'), 13);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('18-10-2020', 'DD-MM-YYYY'), TO_DATE('16-05-2023',
'DD-MM-YYYY'), 926);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('04-08-2022', 'DD-MM-YYYY'), TO_DATE('23-05-2023',
'DD-MM-YYYY'), 292);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('01-11-2020', 'DD-MM-YYYY'), TO_DATE('09-10-2022',
'DD-MM-YYYY'), 712);
```

```
INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
```

```
VALUES (TO_DATE('11-09-2022', 'DD-MM-YYYY'), TO_DATE('20-01-2023',
'DD-MM-YYYY'), 131);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)

VALUES (TO_DATE('15-03-2020', 'DD-MM-YYYY'), TO_DATE('08-05-2023',
'DD-MM-YYYY'), 1158);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('29-10-2022', 'DD-MM-YYYY'), TO_DATE('10-06-2023',
'DD-MM-YYYY'), 225);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)

VALUES (TO_DATE('23-09-2022', 'DD-MM-YYYY'), TO_DATE('27-10-2022',
'DD-MM-YYYY'), 34);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('10-09-2020', 'DD-MM-YYYY'), TO_DATE('25-10-2022',
'DD-MM-YYYY'), 775);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('03-01-2020', 'DD-MM-YYYY'), TO_DATE('14-01-2023',
'DD-MM-YYYY'), 1093);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)

VALUES (TO_DATE('03-04-2020', 'DD-MM-YYYY'), TO_DATE('17-05-2022',
'DD-MM-YYYY'), 765);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
VALUES (TO_DATE('29-10-2020', 'DD-MM-YYYY'), TO_DATE('27-11-2022',
'DD-MM-YYYY'), 753);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)

VALUES (TO_DATE('20-06-2022', 'DD-MM-YYYY'), TO_DATE('24-06-2023',
'DD-MM-YYYY'), 369);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)

VALUES (TO_DATE('23-06-2023', 'DD-MM-YYYY'), TO_DATE('28-06-2023', 'DD-MM-YYYY'),
5);

INSERT INTO transaction_date_dim (arrival_date, dispatch_date, duration_in_warehouse)
```

```
VALUES (TO_DATE('23-01-2023', 'DD-MM-YYYY'), TO_DATE('23-05-2023',
'DD-MM-YYYY'), 122);
```

The screenshot shows the Oracle Live SQL interface. On the left, there's a sidebar with 'My Session' expanded, showing 'Quick SQL', 'My Scripts', 'My Tutorials', and 'Code Library'. The main area is titled 'SQL Worksheet' and contains the following code:

```
40 VALUES (TO_DATE('23-01-2023', 'DD-MM-YYYY'), TO_DATE('23-05-2023', 'DD-MM-YYYY'), 122);
41
42 SELECT * FROM transaction_date_dim;
43
```

Below the code, the results of the query are displayed in a table:

	ARRIVAL_DATE	DISPATCH_DATE	DURATION_IN_WAREHOUSE
1	29-NOV-21	15-JAN-22	413
2	18-SEP-21	26-MAY-22	258
3	14-MAY-21	16-MAY-22	367
4	21-JUN-21	13-OCT-21	115
5	29-OCT-22	07-NOV-22	9
6	04-NOV-22	17-NOV-22	13
7	18-OCT-20	16-MAY-23	926
8	04-AUG-22	23-MAY-23	292
9	01-NOV-20	09-OCT-22	712

At the bottom of the interface, it says '2023 Oracle - Live SQL 23.4.1, running Oracle Database 19c EE Extreme Perf - 19.17.0.0.0 - Database Documentation - Ask Tom - Dev Gym'.

warehouse_dimension:

```
INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)
```

```
VALUES ('W7276', 'BombayX', 'Mumbai', 7597);
```

```
INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)
```

```
VALUES ('W7001', 'DelhiCaps', 'Delhi', 7099);
```

```
INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)
```

```
VALUES ('W1296', 'BombayZ', 'Mumbai', 7597);
```

```
INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)
```

```
VALUES ('W8451', 'BangalorIS', 'Bangalore', 9293);
```

```
INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)
```

```
VALUES ('W8571', 'KerelaSweP', 'Kerela', 6336);
```

```
INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)
```

```
VALUES ('W4968', 'GurgaonHDES', 'Gurgaon', 5676);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W1700', 'HyderabadHub', 'Hyderabad', 5055);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W4035', 'BombayX', 'Mumbai', 7597);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W1934', 'DelhiCaps', 'Delhi', 7099);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W2006', 'BombayZ', 'Mumbai', 7597);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W5878', 'BangalorIS', 'Bangalore', 9293);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W1169', 'KerelaSweP', 'Kerela', 6336);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W6639', 'GurgaonHDES', 'Gurgaon', 5676);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W4150', 'HyderabadHub', 'Hyderabad', 5055);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name,
warehouse_location, city_id)

VALUES ('W5472', 'BombayX', 'Mumbai', 7597);
```

INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)

VALUES ('W8373', 'DelhiCaps', 'Delhi', 7099);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)

VALUES ('W3067', 'BombayZ', 'Mumbai', 7597);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)

VALUES ('W1796', 'BangalorIS', 'Bangalore', 9293);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)

VALUES ('W4563', 'KerelaSweP', 'Kerela', 6336);

INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)

VALUES ('W1452', 'GurgaonHDES', 'Gurgaon', 5676);

The screenshot shows the Oracle Live SQL interface. On the left, there's a sidebar with navigation links like Home, SQL Worksheet, My Session, Schema, Quick SQL, My Scripts, and My Tutorials. The main area is titled 'SQL Worksheet' and contains the following SQL code:

```

28: VALUES ('W4150', 'HyderabadHub', 'Hyderabad', 5055);
29: INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)
30: VALUES ('W5472', 'BombayX', 'Mumbai', 7597);
31: INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)
32: VALUES ('W8373', 'DelhiCaps', 'Delhi', 7099);
33: INSERT INTO warehouse_dim (warehouse_id, warehouse_name, warehouse_location, city_id)

```

Below the code, a table grid displays the data inserted into the warehouse_dim table:

WAREHOUSE_ID	WAREHOUSE_NAME	WAREHOUSE_LOCATION	CITY_ID
W7276	BombayX	Mumbai	7597
W7001	DelhiCaps	Delhi	7099
W1296	BombayZ	Mumbai	7597
W8451	BangalorIS	Bangalore	9293
W8571	KerelaSweP	Kerela	6336
W4968	GurgaonHDES	Gurgaon	5676
W1700	HyderabadHub	Hyderabad	5055
W4035	BombayX	Mumbai	7597

At the bottom of the interface, there's a footer with the text: "2023 Oracle - Live SQL 23.4.1, running Oracle Database 19c EE Extreme Perf - 19.17.0.0.0 - Database Documentation - Ask Tom - Dev Gym" and "Built with ❤️ using Oracle APEX - Privacy - Terms of Use".

inventory_fact_table:

CREATE TABLE inventory_management_fact

AS SELECT

1 AS

transaction_id,

p.product_id,

w.warehouse_id,

s.supplier_id,

ad.arrival_date,

ad.dispatch_date,

0 AS no_of_boxes

FROM

product_dim p

CROSS JOIN

warehouse_dim

w CROSS JOIN

supplier_dim s

CROSS JOIN

transaction_date_dim ad;

Inserting random values for transaction_id and

no_of_boxes BEGIN

FOR rec IN (SELECT rowid, product_id FROM inventory_management_fact) LOOP

UPDATE inventory_management_fact

SET

transaction_id = ROUND(DBMS_RANDOM.VALUE(50, 500)),

```

no_of_boxes = ROUND(DBMS_RANDOM.VALUE(1, 1000))

WHERE rowid =

rec.rowid; END LOOP;

COMMIT

; END;

/

```

SELECT * FROM inventory_management_fact

TRANSACTION_ID	PRODUCT_ID	WAREHOUSE_ID	SUPPLIER_ID	ARRIVAL_DATE	DISPATCH_DATE	NO_OF_BOXES
317	P1235	W1169	S1561	29-NOV-21	15-JAN-23	570
328	P1235	W1169	S2142	29-NOV-21	15-JAN-23	174
132	P1235	W1169	S2603	29-NOV-21	15-JAN-23	826
137	P1235	W1169	S3115	29-NOV-21	15-JAN-23	896
286	P1235	W1169	S3141	29-NOV-21	15-JAN-23	233
68	P1235	W1169	S3284	29-NOV-21	15-JAN-23	42
325	P1235	W1169	S3534	29-NOV-21	15-JAN-23	209
66	P1235	W1169	S3535	29-NOV-21	15-JAN-23	64
401	P1235	W1169	S3568	29-NOV-21	15-JAN-23	235
319	P1235	W1169	S3895	29-NOV-21	15-JAN-23	838

1. Slice:

CREATE MATERIALIZED VIEW mv_slice1

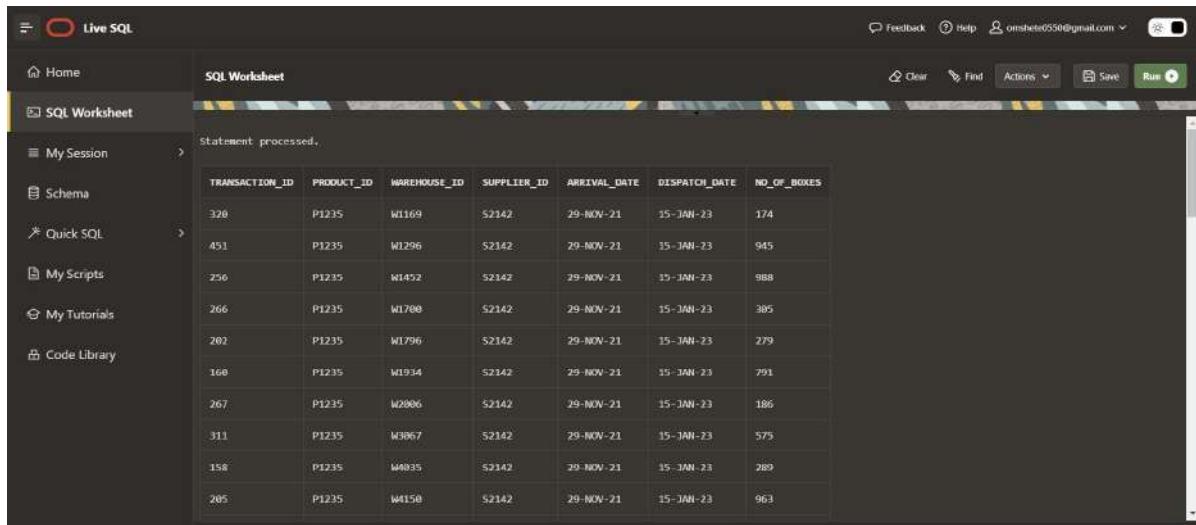
AS SELECT *

FROM

inventory_management_fact

WHERE SUPPLIER_ID = 'S2142';

SELECT * FROM mv_slice1;



The screenshot shows the DBeaver interface with the 'SQL Worksheet' tab selected. The left sidebar has 'My Session' expanded, showing 'Schema', 'Quick SQL', 'My Scripts', 'My Tutorials', and 'Code Library'. The main area displays a table titled 'Statement processed.' with the following data:

TRANSACTION_ID	PRODUCT_ID	WAREHOUSE_ID	SUPPLIER_ID	ARRIVAL_DATE	DISPATCH_DATE	NO_OF_BOXES
320	P1235	W1169	S2142	29-NOV-21	15-JAN-23	174
451	P1235	W1296	S2142	29-NOV-21	15-JAN-23	945
256	P1235	W1452	S2142	29-NOV-21	15-JAN-23	988
266	P1235	W1700	S2142	29-NOV-21	15-JAN-23	395
282	P1235	W1796	S2142	29-NOV-21	15-JAN-23	279
168	P1235	W1934	S2142	29-NOV-21	15-JAN-23	791
267	P1235	W2006	S2142	29-NOV-21	15-JAN-23	186
311	P1235	W3067	S2142	29-NOV-21	15-JAN-23	575
158	P1235	W4035	S2142	29-NOV-21	15-JAN-23	289
285	P1235	W4150	S2142	29-NOV-21	15-JAN-23	963

2. Dice

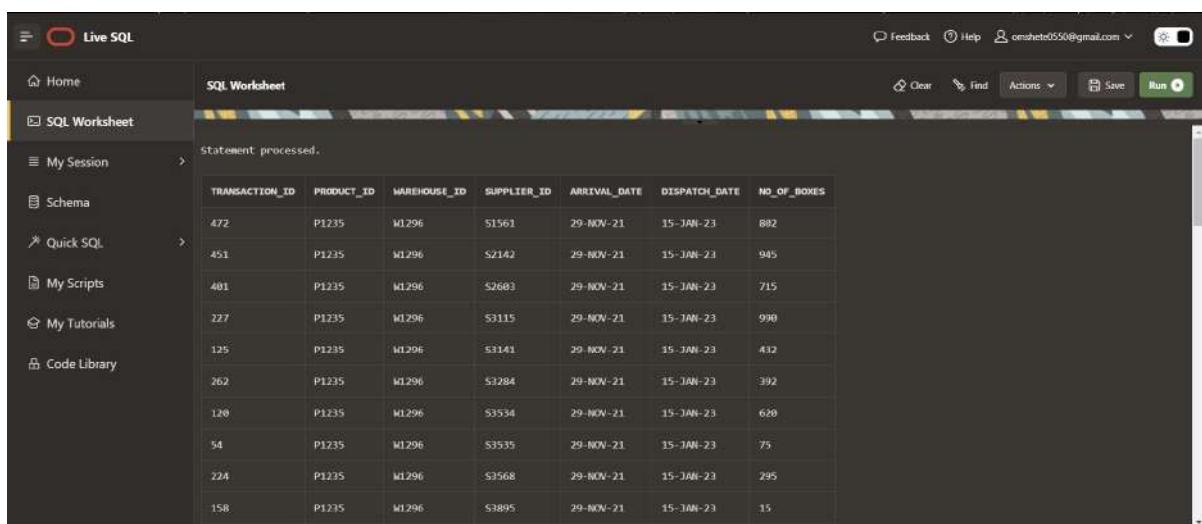
CREATE MATERIALIZED VIEW mv_dice

AS SELECT *

FROM inventory_management_fact

WHERE WAREHOUSE_ID = 'W1296' AND PRODUCT_ID = 'P1235';

SELECT * FROM mv_dice;



The screenshot shows the DBeaver interface with the 'SQL Worksheet' tab selected. The left sidebar has 'My Session' expanded, showing 'Schema', 'Quick SQL', 'My Scripts', 'My Tutorials', and 'Code Library'. The main area displays a table titled 'Statement processed.' with the following data:

TRANSACTION_ID	PRODUCT_ID	WAREHOUSE_ID	SUPPLIER_ID	ARRIVAL_DATE	DISPATCH_DATE	NO_OF_BOXES
472	P1235	W1296	S1561	29-NOV-21	15-JAN-23	892
451	P1235	W1296	S2142	29-NOV-21	15-JAN-23	945
481	P1235	W1296	S2683	29-NOV-21	15-JAN-23	715
227	P1235	W1296	S3115	29-NOV-21	15-JAN-23	990
125	P1235	W1296	S3141	29-NOV-21	15-JAN-23	432
262	P1235	W1296	S3284	29-NOV-21	15-JAN-23	392
128	P1235	W1296	S3534	29-NOV-21	15-JAN-23	620
54	P1235	W1296	S3535	29-NOV-21	15-JAN-23	75
224	P1235	W1296	S3568	29-NOV-21	15-JAN-23	295
158	P1235	W1296	S3895	29-NOV-21	15-JAN-23	15

3. Rollup

```

CREATE OR REPLACE VIEW vw_rollup
AS SELECT
CASE
    WHEN GROUPING(TRANSACTION_ID) = 317 THEN 'Grand
        Total' ELSE TO_CHAR(TRANSACTION_ID)
END AS TRANSACTION_ID,
SUM(NO_OF_BOXES) AS "Total
Quantity" FROM
inventory_management_fact GROUP BY
ROLLUP(TRANSACTION_ID);

```

```

CREATE MATERIALIZED VIEW mv_rollup_1 AS
SELECT * FROM vw_rollup;

```

The screenshot shows the Databricks Live SQL interface. The sidebar on the left has sections for Home, Schema, Quick SQL, My Scripts, My Tutorials, and Code Library. The main area is titled 'SQL Worksheet' and contains the results of a query. The results are presented in a table with two columns: 'TRANSACTION_ID' and 'Total_Quantity'. The data is as follows:

TRANSACTION_ID	Total_Quantity
50	88329
51	169269
52	192914
53	175933
54	170068
55	170292
56	189967
57	196045
58	182652
59	176034
60	186782

4. Drilldown

```

CREATE OR REPLACE VIEW vw_drilldown
AS SELECT

```

```
TRANSACTION_ID,  
WAREHOUSE_ID,  
SUM(no_of_boxes) AS "Total  
Boxes" FROM  
inventory_management_fact  
GROUP BY TRANSACTION_ID, WAREHOUSE_ID;
```

```
CREATE MATERIALIZED VIEW mv_drilldown
```

```
AS SELECT * FROM vw_drilldown;
```

The screenshot shows the Databricks Live SQL interface. On the left, there is a sidebar with navigation links: Home, SQL Worksheet (which is selected), My Session, Schema, Quick SQL, My Scripts, My Tutorials, and Code Library. The main area is titled "SQL Worksheet" and displays a table of data. The table has three columns: TRANSACTION_ID, WAREHOUSE_ID, and Total Boxes. The data consists of 12 rows:

TRANSACTION_ID	WAREHOUSE_ID	Total Boxes
317	W1169	8719
424	W1169	11037
491	W1296	2564
126	W1296	6300
54	W1296	8180
482	W1296	6363
174	W1296	10231
412	W1452	11205
225	W1452	5051
369	W1452	7487
482	W1452	13760

EXPERIMENT - 04

Implementation

: Pivot Table:

SUM OF WEEKLY_SALES	DRUG_NAME	Category	AllerClear 10	CardioGuard 50	Dermacare 100	DiabetCare 500	ImmunoGuard 7	MediHealth 100	NeuroCalm XR	Resigen Cough	StomachEase PI	VitaBoost Multi	Grand Total
Allergy	64000												64000
Analgesics										120000			120000
Cardiovascular			40000										40000
Cough and Cold										90000			90000
Dermatology				48000									48000
Endocrinology					37500								37500
Gastrointestinal										102000			102000
Immunology						70000							70000
Neurology									36000				36000
Vitamins and Supplements											84000		84000
Grand Total	64000		40000	48000	37500	70000	120000	36000	90000	102000	84000	691500	

Slice:

The screenshot shows a Microsoft Excel spreadsheet with a pivot table editor open. The main table area displays sales data by month and category. The pivot table editor sidebar shows fields for Sales Date, Category Name, Drug Name, Weekly Sales, and Monthly Sales. The Sales Date field is currently selected.

SALES_DATE	CATEGORY_NAME	SUM of MONTH
01/08/2023	Analgesics	80000
	Cough and Cold	60000
	Gastrointestinal	68000
01/08/2023 Total		208000
02/08/2023	Allergy	32000
	Cough and Cold	60000
	Dermatology	48000
	Endocrinology	50000
	Immunology	56000
	Vitamins and Su	84000
02/08/2023 Total		330000
03/08/2023	Allergy	32000
	Analgesics	80000
	Cardiovascular	40000
	Neurology	36000
	Vitamins and Su	84000
03/08/2023 Total		272000
04/08/2023	Allergy	32000
	Cardiovascular	40000
	Gastrointestinal	68000
	Neurology	36000
	Vitamins and Su	84000

Dice:

SUM of MONTH	CATEGORY_NAME							
SALES_DATE	Allergy	Analgesics	Cardiovascular	Cough and Cold	Dermatology	Endocrinology	G	Grand Total
01/08/2023		80000		60000				
02/08/2023	32000			60000	48000	50000		
03/08/2023	32000	80000	40000					

Endocrinology	Gastrointestinal	Immunology	Neurology	Vitamins and Su	Grand Total
	68000				208000
00	50000		56000	84000	330000
			36000	84000	272000

Rollup:

SUM of MONTH	CATEGORY_NAME							
SALES_DATE	Allergy	Analgesics	Cardiovascular	Cough and Cold	Dermatology	Endocrinology	G	Grand Total
01/08/2023		80000		60000				
02/08/2023	32000			60000	48000	50000		
03/08/2023	32000	80000	40000					
04/08/2023	32000		40000					
05/08/2023	32000	80000		60000	48000			
06/08/2023	32000	80000		60000		50000		
07/08/2023	32000		40000			50000		
08/08/2023				60000				
09/08/2023	32000	80000			48000			
10/08/2023	32000	80000	40000	60000	48000			
Grand Total	256000	480000	160000	360000	192000	150000		

Endocrinology	Gastrointestinal	Immunology	Neurology	Vitamins and Su	Grand Total
	68000				208000
00	50000		56000	84000	330000
			36000	84000	272000
	68000		36000	84000	260000
00	68000	56000			344000
	50000	68000	56000		346000
	50000	68000	56000	84000	330000
		68000			128000
00		56000	36000		252000
00			36000		296000
00	150000	408000	280000	144000	336000
					2766000

Drilldown:

table (3)

File Edit View Insert Format Data Tools Extensions Help

M26

	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	
3	SUM of MOI	CATEGORY	DRUG NAME														
4		=Alergy	Alergy Total	=Antigens	Antigens Total	=Cardiovascular	Cardiovascular Total	=Cough and Cough and Cough	Cough and Cough and Cough Total	=Dermatology	Dermatology Total	=Endocrinology	Endocrinology Total	=Gastroenterology	Gastroenterology Total	=Immunology	Immunology Total
5	SALES_DATE	MonthYear	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	Medicare	
6	01/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
7	02/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
8	03/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
9	04/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
10	05/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
11	06/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
12	07/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
13	08/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
14	09/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
15	10/09/2023	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	32000	
16	Grand Total	296000	296000	480000	480000	160000	160000	390000	390000	360000	360000	192000	192000	192000	192000	192000	
17																	
18																	
19																	
20																	
21																	
22																	
23																	
24																	
25																	
26																	
27																	
28																	
29																	
30																	

table (3).csv

	R	S	T	U	V	W	X	Y	Z	AA	AB
Enteral Enteral Enteral Enteral Gastrointestinal Gastrointestinal Immunology Immunology Immunology Neurology Neurology Vitamins and Vitamins and Su Grand Total											
7	DiabeticCare 500		StomachEase Plus		ImmunoGuard 75		NeuroCalm XR		VitaBoost Multivitamin		
8		68000		68000						208000	
9	50000	50000			56000	56000			84000	84000	130000
10							38000	38000	81000	81000	127000
11		68000	68000			35000	35000	84000	84000	126000	
12		68000	68000	58000	58000					144000	
13	50000	50000	68000	68000	56000	56000				146000	
14	50000	50000	68000	68000	59000	59000			84000	84000	130000
15		68000	68000							128000	
16				56000	56000	35000	35000			252000	
17					35000	35000	35000			196000	
18	180000	160000	408000	408000	280000	280000	144000	144000	336000	336000	2766000
19											
20											

EXPERIMENT - 05

Aim: Implementation of Bayesian algorithm

Theory:

Data Mining Bayesian Classifiers

The connection between the attribute set and the class variable is non-deterministic in numerous applications. In other words, we can say the class label of a test record can't be assumed with certainty even though its attribute set is the same as some of the training examples. These circumstances may emerge due to the noisy data or the presence of certain confusing factors that influence classification, but it is not included in the analysis. For example, consider the task of predicting the occurrence of whether an individual is at risk for liver illness based on the individual's eating habits and working efficiency. Although most people who eat healthy and exercise consistently have less probability of occurrence of liver disease, they may still do so due to other factors. For example, due to the consumption of high-calorie street foods and alcohol abuse. Determining whether an individual's eating routine is healthy or the workout efficiency is sufficient is also subject to analysis, which may introduce vulnerabilities to the learning issue.

Bayesian classification uses Bayes theorem to predict the occurrence of any event. Bayesian classifiers are the statistical classifiers with the Bayesian probability understandings. The theory expresses how a level of belief, is expressed as a probability.

Bayes theorem came into existence after Thomas Bayes, who first utilized conditional probability to provide an algorithm that uses evidence to calculate limits on an unknown parameter.

Bayes's theorem is expressed mathematically by the following equation that is given below. For proposition X and evidence Y,

- $P(X)$, the prior, is the primary degree of belief in X
- $P(X/Y)$, the posterior is the degree of belief having accounted for Y.

- The quotient $\frac{P(Y/X)}{P(Y)}$ represents the supports Y provides for

X. Bayes theorem can be derived from the conditional probability:

$$P(X/Y) = \frac{P(X \cap Y)}{P(Y)}, \text{ if } P(Y) \neq 0$$

$$P(Y/X) = \frac{P(Y \cap X)}{P(X)}, \text{ if } P(X) \neq 0$$

Where $P(X \cap Y)$ is the **joint probability** of both X and Y being true, because

$$P(Y \cap X) = P(X \cap Y)$$

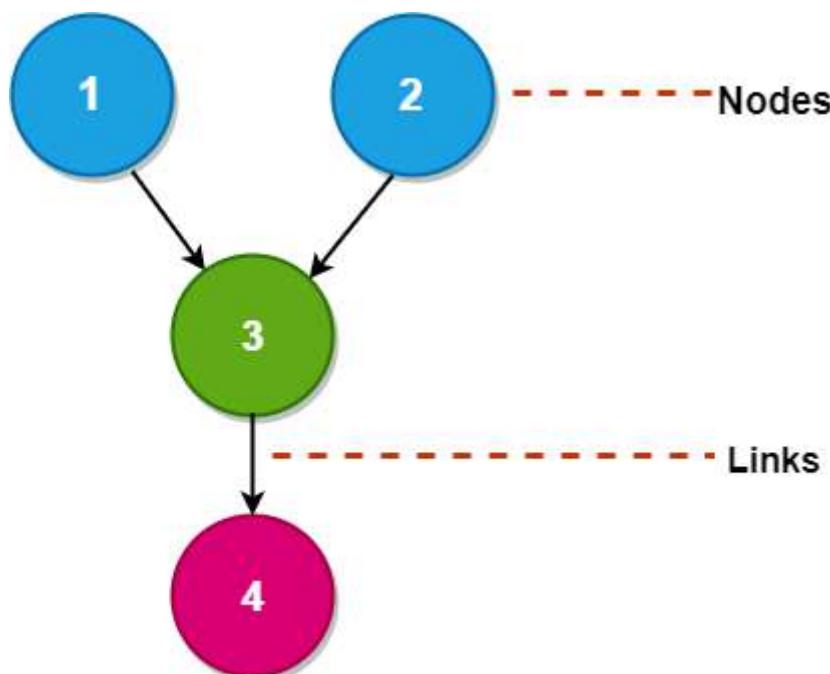
$$\text{or, } P(X \cap Y) = P(X/Y)P(Y) = P(Y/X)P(X)$$

$$\text{or, } P(X/Y) = \frac{P(Y/X)P(X)}{P(Y)}, \text{ if } P(Y) \neq 0$$

Bayesian network:

A Bayesian Network falls under the classification of Probabilistic Graphical Modelling (PGM) procedure that is utilized to compute uncertainties by utilizing the probability concept. Generally known as **Belief Networks, Bayesian Networks** are used to show uncertainties using **Directed Acyclic Graphs (DAG)**

A **Directed Acyclic Graph** is used to show a Bayesian Network, and like some other statistical graph, a DAG consists of a set of nodes and links, where the links signify the connection between the nodes.



The nodes here represent random variables, and the edges define the relationship between these variables.

A DAG models the uncertainty of an event taking place based on the Conditional Probability Distribution (CPD) of each random variable. A **Conditional Probability Table** (CPT) is used to represent the CPD of each variable in a network.

CODE:

```
import pandas as
pd import math
data =


pd.read_csv("/content/inventory_bayes.csv") price

= input("Enter Price (Low, Average, High): ")

number_of_boxes = int(input("Enter Number Of Boxes: (30, 40, 10): "))
product_brand = input("Enter Product Brand (Sony, LG, Nike):
") location = input("Enter location: (Suburbs, City, Downtown):
")

print("\n")
#prior probabilities
total_instances =
len(data)
package_left_yes = len(data[data['package_left'] ==
'Yes']) package_left_no = len(data[data['package_left'] ==
'No'])

p_yes = round(package_left_yes/total_instances,
4) p_no = round(package_left_no/total_instances,
4)

print(f"Total occurrences of Yes:
{package_left_yes}") print(f"Total occurrences of No:
{package_left_no}")

print(f"Prior probability for Yes
{p_yes}") print(f"Prior probability for
No: {p_no}")

#Likelihood Probabilities
```

```

def calculate_likelihood(attribute, value,
    package_left): subset = data[data['package_left'] ==
    package_left] count = len(subset[subset[attribute] ==
    value])

total = len(subset)

return count /
total

print('\n')
print('\n')
print("The Likelihood Probabilities are: ")

#
p_price_given_yes = round(calculate_likelihood('price', price, 'Yes'),
4) p_price_given_no = round(calculate_likelihood('price', price, 'No'),
4)

print(f"P(price = {price} / buy='Yes') =
{p_price_given_yes}") print(f"P(price = {price} / buy='No') =
{p_price_given_no}") print('\n')

#
p_number_of_boxes_given_yes = round(calculate_likelihood('number_of_boxes',
number_of_boxes, 'Yes'), 4)
p_number_of_boxes_given_no = round(calculate_likelihood('number_of_boxes',
number_of_boxes, 'No'), 4)

print(f"P(number_of_boxes = {number_of_boxes} / buy='Yes') =
{p_number_of_boxes_given_yes}")
print(f"P(number_of_boxes = {number_of_boxes} / buy='No') =
{p_number_of_boxes_given_no}")
) print('\n')

#
p_product_brand_given_yes = round(calculate_likelihood('product_brand',
product_brand, 'Yes'), 4)

```

```

p_product_brand_given_no = round(calculate_likelihood('product_brand',
product_brand, 'No'), 4)

print(f"P(product_brand = {product_brand} / buy='Yes') =
{p_product_brand_given_yes}") print(f"P(product_brand = {product_brand} / buy='No') =
{p_product_brand_given_no}") print('\n')

#
p_location_given_yes = round(calculate_likelihood('location', location, 'Yes'), 4)
p_location_given_no = round(calculate_likelihood('location', location, 'No'), 4)

print(f"P(location = {location} / buy='Yes') =
{p_location_given_yes}") print(f"P(location = {location} / buy='No') =
{p_location_given_no}") print('\n')

# Posterior probabilities
print("The posterior probabilities are: ")

p_yes_given_x = round(p_yes * p_price_given_yes * p_number_of_boxes_given_yes
* p_product_brand_given_yes * p_location_given_yes,4)
p_no_given_x = round(p_no * p_price_given_no * p_number_of_boxes_given_yes
* p_product_brand_given_no * p_location_given_no,4)

if p_yes_given_x > p_no_given_x:
    prediction =
    'Yes' else:
    prediction = 'No'

print(f'Probability          of      Purchased=Yes:
{p_yes_given_x}') print(f'Probability of Purchased=No:
{p_no_given_x}') print(f'Prediction: {prediction}')

```

Output:

```

DWM_Exp5.ipynb
File Edit View Insert Runtime Tools Help All changes saved
Comment Share RAM Disk
+ Code + Text
Enter Price (Low, Average, High): High
Enter Number Of Boxes: (30, 40, 10): 30
Enter Product Brand (Sony, LG, Nike): Sony
Enter location: (Suburbs, City, Downtown): Downtown

Total occurrences of Yes: 24
Total occurrences of No: 16
Prior probability for Yes 0.5854
Prior probability for No: 0.3902

The Likelihood Probabilities are:
P(price = High / buy='Yes') = 0.375
P(price = High / buy='No') = 0.4375

P(number_of_boxes = 30 / buy='Yes') = 0.2917
P(number_of_boxes = 30 / buy='No') = 0.0625

P(product_brand = Sony / buy='Yes') = 0.3333
P(product_brand = Sony / buy='No') = 0.25

P(location = Downtown / buy='Yes') = 0.4583
P(location = Downtown / buy='No') = 0.3125

The posterior probabilities are:

```

EXPERIMENT - 06

Aim: Perform data Pre-processing task and Demonstrate performing Classification, Clustering, and Association algorithms on data sets using a data mining tool (WEKA/R tool)

Theory:

Classification:

The term "classification" is usually used when there are exactly two target classes called binary classification. When more than two classes may be predicted, specifically in pattern recognition problems, this is often referred to as multinomial classification. However, multinomial classification is also used for categorical response data, where one wants to predict which category amongst several categories has the instances with the highest probability.

Classification is one of the most important tasks in data mining. It refers to a process of assigning pre-defined class labels to instances based on their attributes. There is a similarity between classification and clustering, it looks similar, but it is different. The major difference between classification and clustering is that classification includes the levelling of items according to their membership in pre-defined groups. Let's understand this concept with the help of an example; suppose you are using a self-organizing map neural network algorithm for image recognition where there are 10 different kinds of objects. If you label each image with one of these 10 classes, the classification task is solved.

Clustering:

Clustering refers to a technique of grouping objects so that objects with the same functionalities come together and objects with different functionalities go apart. In other words, we can say that clustering is a process of portioning a data set into a set of meaningful subclasses, known as clusters. Clustering is the same as classification in which data is grouped. However, unlike classification, the groups are not previously defined.

Instead, the grouping is achieved by determining similarities between data according to characteristics found in the real data. The groups are called Clusters.

Association Rules:

Association rule learning is a type of unsupervised learning technique that checks for the dependency of one data item on another data item and maps accordingly so that it can be more profitable. It tries to find some interesting relations or associations among the variables of the dataset. It is based on different rules to discover the interesting relations between variables in the database.

Association rule learning is one of the very important concepts of machine learning, and it is employed in **Market Basket analysis, Web usage mining, continuous production, etc.** Here market basket analysis is a technique used by the various big retailers to discover the associations between items. We can understand it by taking the example of a supermarket, as in a supermarket, all products that are purchased together are put together.



Association rule learning can be divided into three types of algorithms:

1. **Apriori**
2. **Eclat**
3. **F-P Growth Algorithm**

Support

Support is the frequency of A or how frequently an item appears in the dataset. It is defined as the fraction of the transaction T that contains the itemset X. If there are X datasets, then for transactions T, it can be written as:

$$\text{Supp}(X) = \frac{\text{Freq}(X)}{T}$$

Confidence

Confidence indicates how often the rule has been found to be true. Or how often the items X and Y occur together in the dataset when the occurrence of X is already given. It is the ratio of the transaction that contains X and Y to the number of records that contain X.

$$\text{Confidence} = \frac{\text{Freq}(X,Y)}{\text{Freq}(X)}$$

Lift

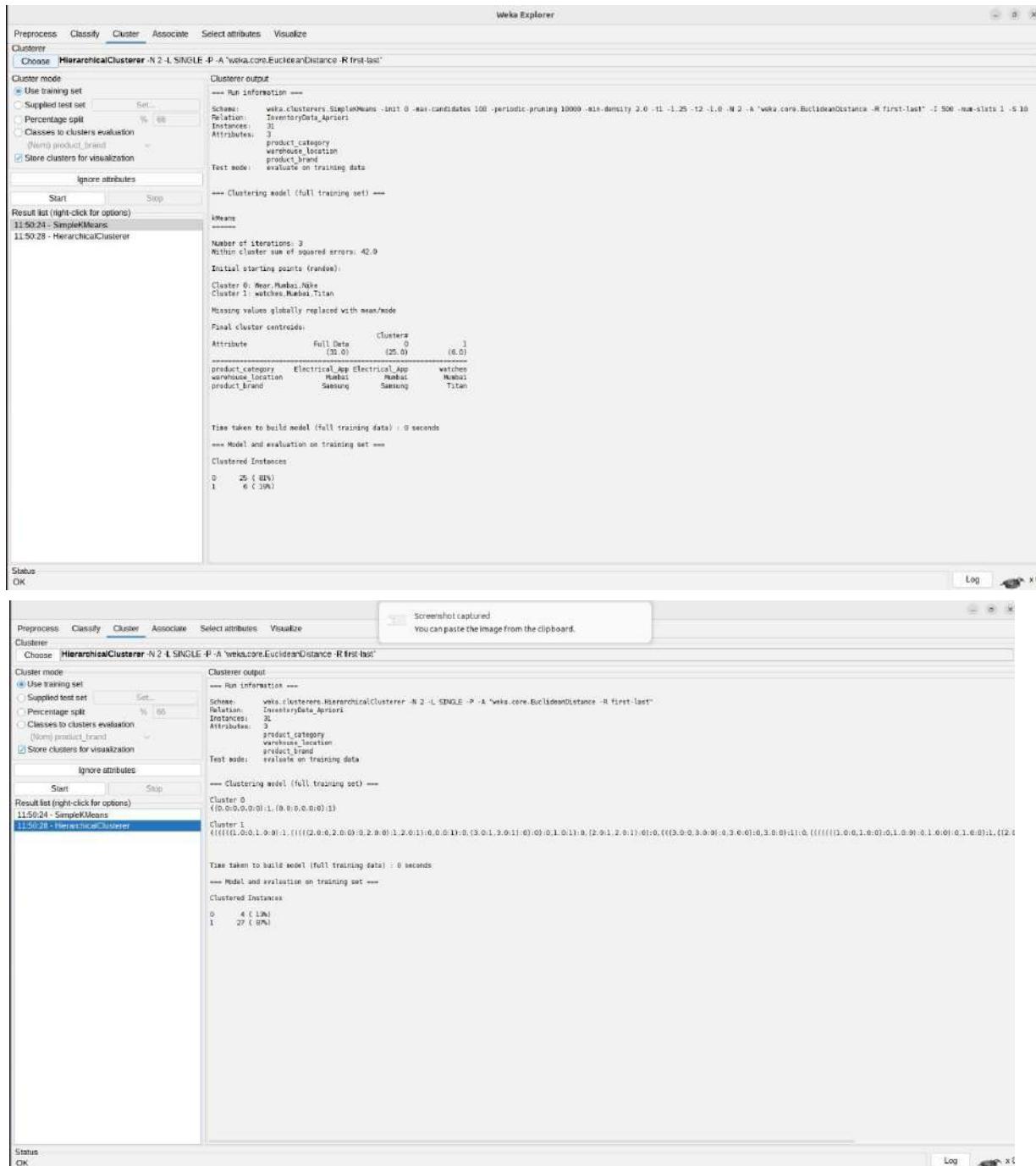
It is the strength of any rule, which can be defined as below formula:

$$\text{Lift} = \frac{\text{Supp}(X,Y)}{\text{Supp}(X) \times \text{Supp}(Y)}$$

It is the ratio of the observed support measure and expected support if X and Y are independent of each other. It has three possible values:

- If **Lift= 1**: The probability of occurrence of antecedent and consequent is independent of each other.
- **Lift>1**: It determines the degree to which the two itemsets are dependent to each other.
- **Lift<1**: It tells us that one item is a substitute for other items, which means one item has a negative effect on another.

Output:



Screenshot captured
You can paste the image from the clipboard.

Classifier output

Test note: 10-fold cross-validation

== Classifier model (full training set) ==

Naive Bayes Classifier

Attribute	Class	Titan	Seung	Phillips	Nike
	0.17	0.21	0.29	0.23	

(None) warehouse_location

Start Stop

Result list (right-click for options)

- 11:49:51 - bayes.NaiveBayes
- 11:50:06 - rules.DecisionTable
- 11:50:10 - rules.DecisionTable
- 11:53:07 - trees.J48
- 11:53:35 - trees.J48

Time taken to build model: 0 seconds

== Stratified cross-validation ==

== Summary ==

Correctly Classified Instances	13	41.855 %
Incorrectly Classified Instances	18	58.045 %
Kappa statistic	0.196	
Mean absolute error	0.3428	
Root mean squared error	0.418	
Relative absolute error	92.4127 %	
Root relative squared error	96.7671 %	
Total Number of Instances	36	

== Detailed Accuracy By Class ==

TP Rate	FP Rate	Precision	Recall	F-Measure	NCC	RCC Area	PRC Area	Class
0.400	0.938	0.667	0.400	0.500	0.459	0.898	0.456	Titan
0.700	0.020	0.438	0.700	0.538	0.254	0.595	0.395	Seung
0.800	0.020	0.800	0.800	0.800	0.800	0.960	0.960	Phillips
0.571	0.167	0.500	0.571	0.533	0.387	0.625	0.472	Nike
Weighted Avg.	0.419	0.205	0.362	0.416	0.375	0.179	0.576	0.388

== Confusion Matrix ==

a = 1	b = 2	c = 3	d = 4
21 0 2	0 1 0	0 7 3 0	0 0 0 1
0 7 3 0	0 0 0 1	0 0 0 1	0 0 0 1
1 6 0 2	0 0 0 1	0 0 0 1	0 0 0 1
0 2 1 4	0 0 0 1	0 0 0 1	0 0 0 1

Status OK

Weka Explorer

Classifier output

== Run information ==

Schema: weka.classifiers.trees.J48 - C 0.25 - N 2

Relation: InventoryData_Apriori

Instances: 31

Attributes: 3

(None) warehouse_location

Start Stop

Result list (right-click for options)

- 11:49:51 - bayes.NaiveBayes
- 11:50:06 - rules.DecisionTable
- 11:50:10 - rules.DecisionTable
- 11:53:07 - trees.J48
- 11:53:35 - trees.J48

Time taken to build model: 0 seconds

== Stratified cross-validation ==

== Summary ==

Correctly Classified Instances	17	54.839 %
Incorrectly Classified Instances	14	45.161 %
Kappa statistic	-0.0057	
Mean absolute error	0.5205	
Root mean squared error	0.5126	
Relative absolute error	162.2585 %	
Root relative squared error	103.2799 %	
Total Number of Instances	31	

== Detailed Accuracy By Class ==

TP Rate	FP Rate	Precision	Recall	F-Measure	NCC	RCC Area	PRC Area	Class
0.944	1.000	0.567	0.944	0.700	-0.155	0.556	0.456	Delhi
0.800	0.000	0.800	0.800	0.800	-0.155	0.556	0.356	Mumbai
0.548	0.000	0.500	0.548	0.411	-0.155	0.556	0.456	Delhi
Weighted Avg.	0.548	0.000	0.500	0.548	0.411	-0.155	0.556	0.456

== Confusion Matrix ==

a = 1	b = 2	c = 3	d = 4
17 1 1	1 0 1	0 0 0 1	0 0 0 1
1 0 1	0 0 0 1	0 0 0 1	0 0 0 1
0 0 0 1	0 0 0 1	0 0 0 1	0 0 0 1
0 0 0 1	0 0 0 1	0 0 0 1	0 0 0 1



EXPERIMENT - 07

Aim: Implementation of Clustering algorithm (K-means/Kmedoids)

Theory:

K means Clustering

K-Means Clustering is an Unsupervised Machine Learning algorithm, which groups the unlabeled dataset into different clusters.

Unsupervised Machine Learning is the process of teaching a computer to use unlabeled, unclassified data and enabling the algorithm to operate on that data without supervision. Without any previous data training, the machine's job in this case is to organize unsorted data according to parallels, patterns, and variations.

The goal of clustering is to divide the population or set of data points into a number of groups so that the data points within each group are more comparable to one another and different from the data points within the other groups. It is essentially a grouping of things based on how similar and different they are to one another.

We are given a data set of items, with certain features, and values for these features (like a vector). The task is to categorize those items into groups. To achieve this, we will use the K-means algorithm; an unsupervised learning algorithm. 'K' in the name of the algorithm represents the number of groups/clusters we want to classify our items into.

(It will help if you think of items as points in an n-dimensional space). The algorithm will categorize the items into k groups or clusters of similarity. To calculate that similarity, we will use the Euclidean distance as a measurement.

The algorithm works as follows:

1. First, we randomly initialize k points, called means or cluster centroids.
2. We categorize each item to its closest mean and we update the mean's coordinates, which are the averages of the items categorized in that cluster so far.
3. We repeat the process for a given number of iterations and at the end, we have our clusters.

The "points" mentioned above are called means because they are the mean values of the items categorized in them. To initialize these means, we have a lot of options. An intuitive

method is to initialize the means at random items in the data set. Another method is to initialize the means at random values between the boundaries of the data set (if for a feature x , the items have values in $[0,3]$, we will initialize the means with values for x at $[0,3]$).

The above algorithm in pseudocode is as follows:

Initialize k means with random values

--> For a given number of iterations:

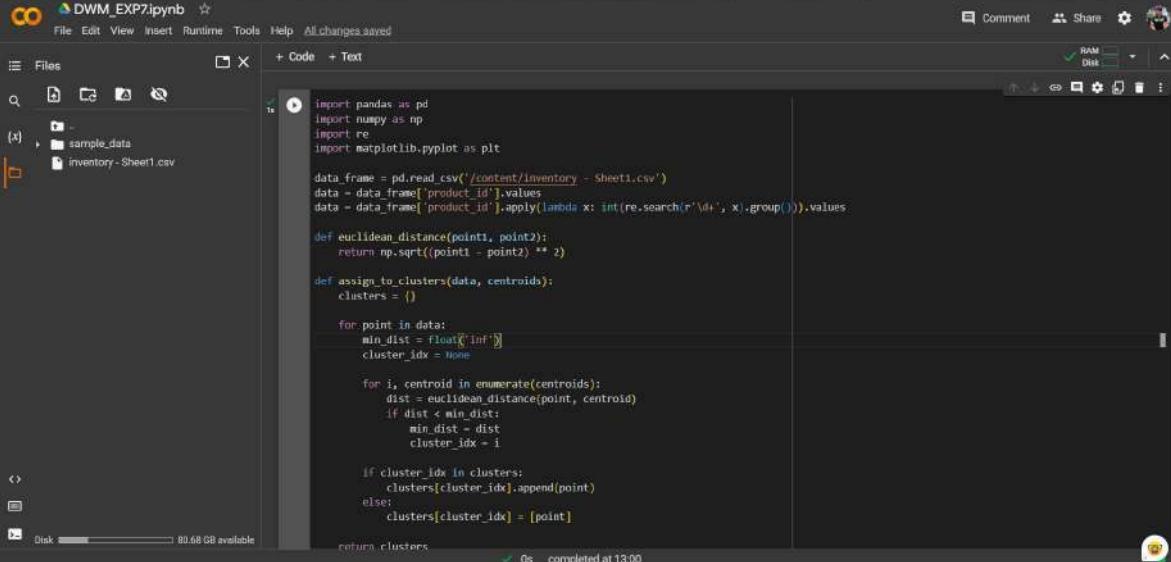
--> Iterate through items:

--> Find the mean closest to the item by calculating
the Euclidean distance of the item with each of the means

--> Assign item to mean

--> Update mean by shifting it to the average of the items in that cluster

Output:



```

File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
import pandas as pd
import numpy as np
import re
import matplotlib.pyplot as plt

data_frame = pd.read_csv('/content/inventory - Sheet1.csv')
data = data_frame['product_id'].values
data = data_frame['product_id'].apply(lambda x: int(re.search(r'\d+', x).group(0))).values

def euclidean_distance(point1, point2):
    return np.sqrt((point1 - point2) ** 2)

def assign_to_clusters(data, centroids):
    clusters = {}

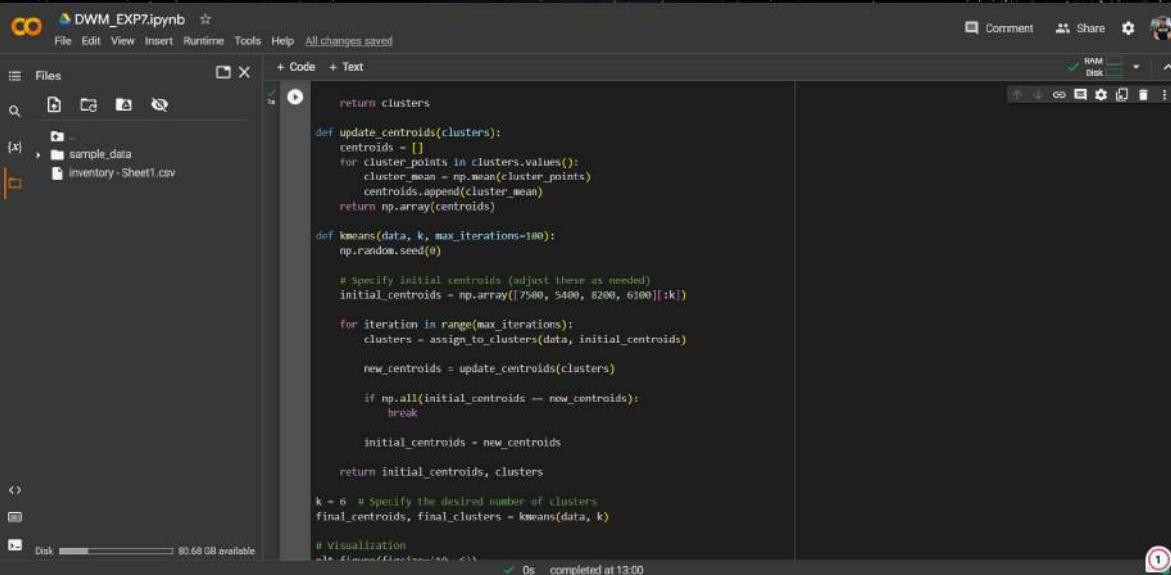
    for point in data:
        min_dist = float('inf')
        cluster_idx = None

        for i, centroid in enumerate(centroids):
            dist = euclidean_distance(point, centroid)
            if dist < min_dist:
                min_dist = dist
                cluster_idx = i

        if cluster_idx in clusters:
            clusters[cluster_idx].append(point)
        else:
            clusters[cluster_idx] = [point]

    return clusters

```



```

return clusters

def update_centroids(clusters):
    centroids = []
    for cluster_points in clusters.values():
        cluster_mean = np.mean(cluster_points)
        centroids.append(cluster_mean)
    return np.array(centroids)

def kmeans(data, k, max_iterations=100):
    np.random.seed(0)

    # Specify initial centroids (adjust these as needed)
    initial_centroids = np.array([7500, 5400, 8200, 6300])[:k]

    for iteration in range(max_iterations):
        clusters = assign_to_clusters(data, initial_centroids)

        new_centroids = update_centroids(clusters)

        if np.all(initial_centroids == new_centroids):
            break

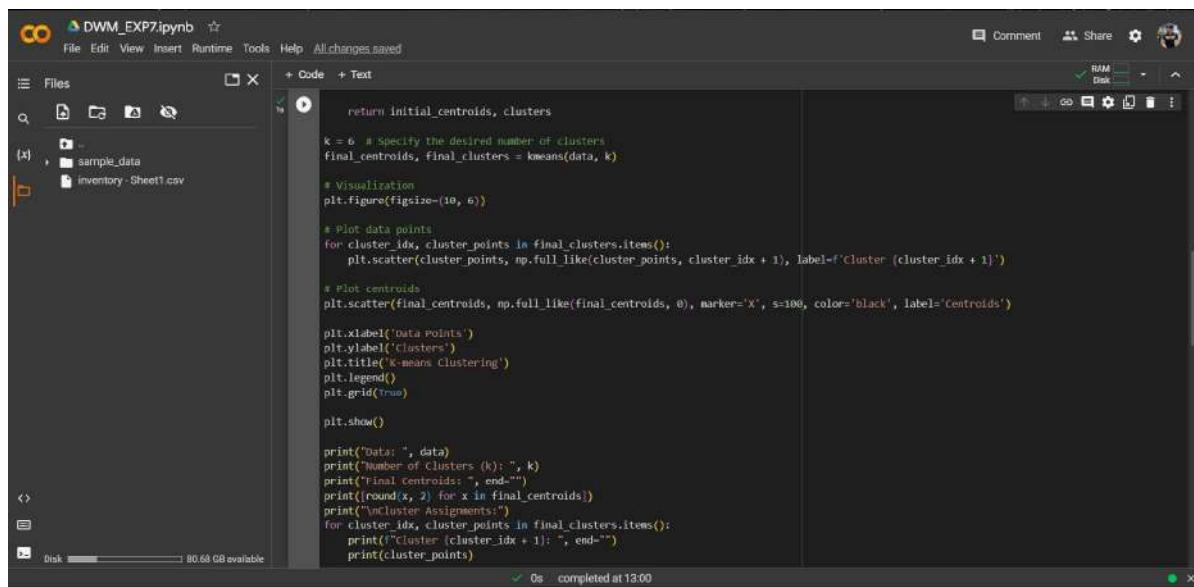
        initial_centroids = new_centroids

    return initial_centroids, clusters

k = 6 # Specify the desired number of clusters
final_centroids, final_clusters = kmeans(data, k)

# Visualization
plt.figure(figsize=(10, 6))

```



DWM_EXP7.ipynb

```

def kmeans_clustering(data, k):
    # Initialize centroids
    initial_centroids = data[0:k]
    final_centroids = initial_centroids

    # Visualization
    plt.figure(figsize=(10, 6))

    # Plot data points
    for cluster_idx, cluster_points in final_clusters.items():
        plt.scatter(cluster_points, np.full_like(cluster_points, cluster_idx + 1), label=f'Cluster {cluster_idx + 1}')

    # Plot centroids
    plt.scatter(final_centroids, np.full_like(final_centroids, 0), marker='x', s=100, color='black', label='Centroids')

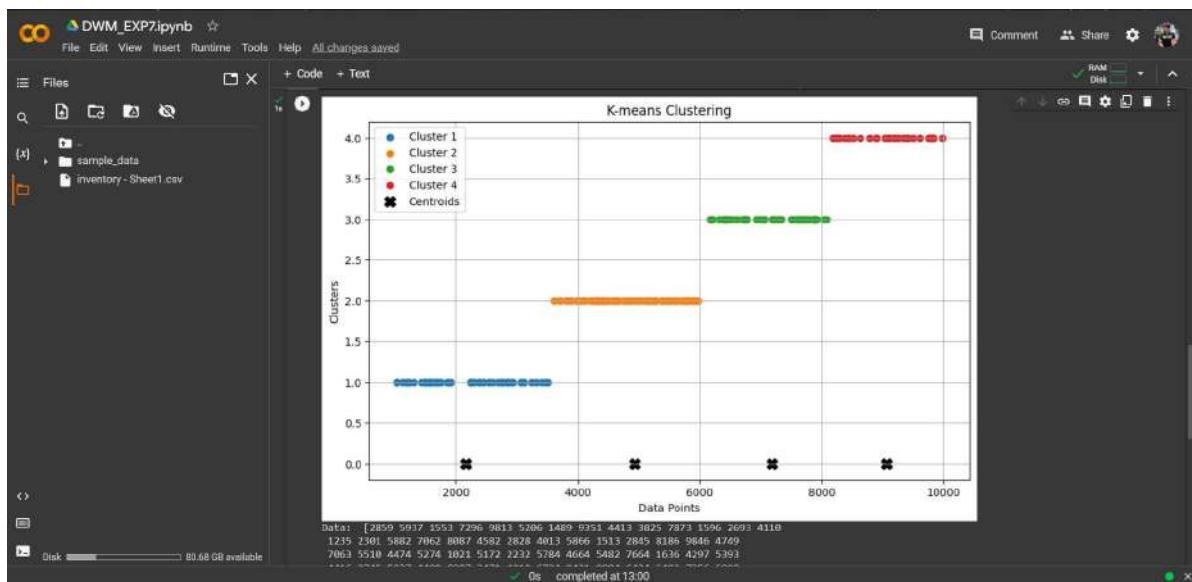
    plt.xlabel('Data Points')
    plt.ylabel('Clusters')
    plt.title('K-means Clustering')
    plt.legend()
    plt.grid(True)

    plt.show()

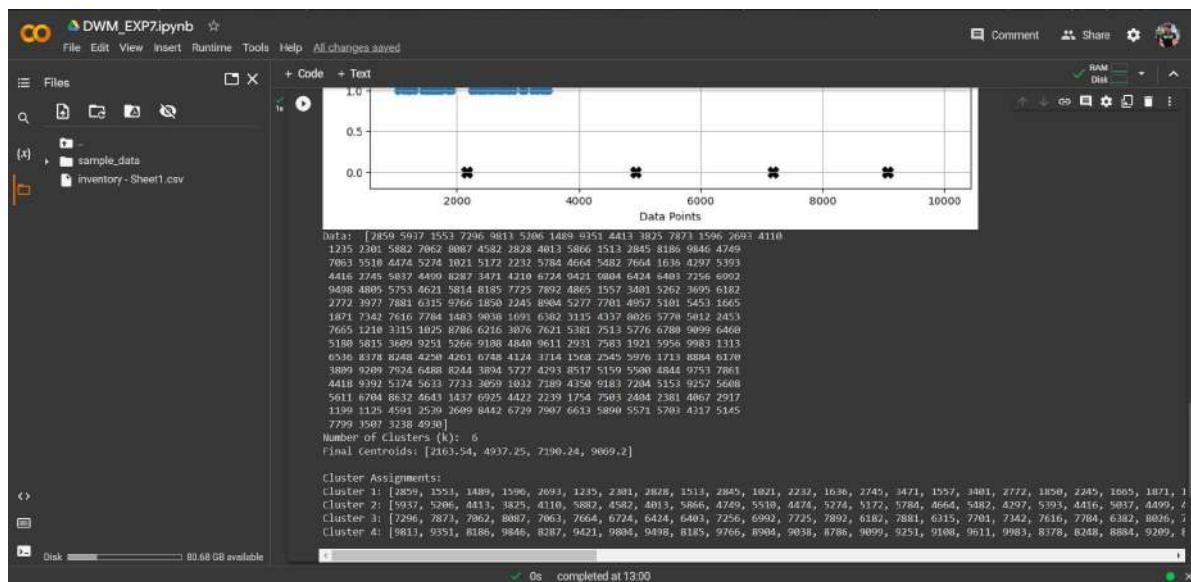
print("Data: ", data)
print("Number of Clusters (k): ", k)
print("Final Centroids: ", end="")
print([round(x, 2) for x in final_centroids])
print("\nCluster Assignments:")
for cluster_idx, cluster_points in final_clusters.items():
    print(f"Cluster {cluster_idx + 1}: ", end="")
    print(cluster_points)

```

0s completed at 13:00



C32_Om_2103163_DWM



EXPERIMENT - 08

Aim: Implementation of any one Hierarchical Clustering method

Theory:

A Hierarchical clustering method works via grouping data into a tree of clusters.

Hierarchical clustering begins by treating every data point as a separate cluster. Then, it repeatedly executes the subsequent steps:

Identify the 2 clusters which can be closest together, and

Merge the 2 maximum comparable clusters. We need to continue these steps until all the clusters are merged together.

In Hierarchical Clustering, the aim is to produce a hierarchical series of nested clusters. A diagram called a Dendrogram (A Dendrogram is a tree-like diagram that statistics the sequences of merges or splits) graphically represents this hierarchy and is an inverted tree that describes the order in which factors are merged (bottom-up view) or clusters are broken up (top-down view).

Hierarchical clustering is a method of cluster analysis in data mining that creates a hierarchical representation of the clusters in a dataset. The method starts by treating each data point as a separate cluster and then iteratively combines the closest clusters until a stopping criterion is reached. The result of hierarchical clustering is a tree-like structure, called a dendrogram, which illustrates the hierarchical relationships among the clusters.

Hierarchical clustering has a number of advantages over other clustering methods, including:

1. The ability to handle non-convex clusters and clusters of different sizes and densities.
2. The ability to handle missing data and noisy data.
3. The ability to reveal the hierarchical structure of the data, can be useful for understanding the relationships among the clusters. However, it also has some drawbacks, such as:
4. The a need for a criterion to stop the clustering process and determine the final number of clusters.
5. The computational cost and memory requirements of the method can be high, especially for large datasets.
6. The results can be sensitive to the initial conditions, linkage criterion, and distance metric used.

In summary, Hierarchical clustering is a method of data mining that groups similar data points into clusters by creating a hierarchical structure of the clusters.

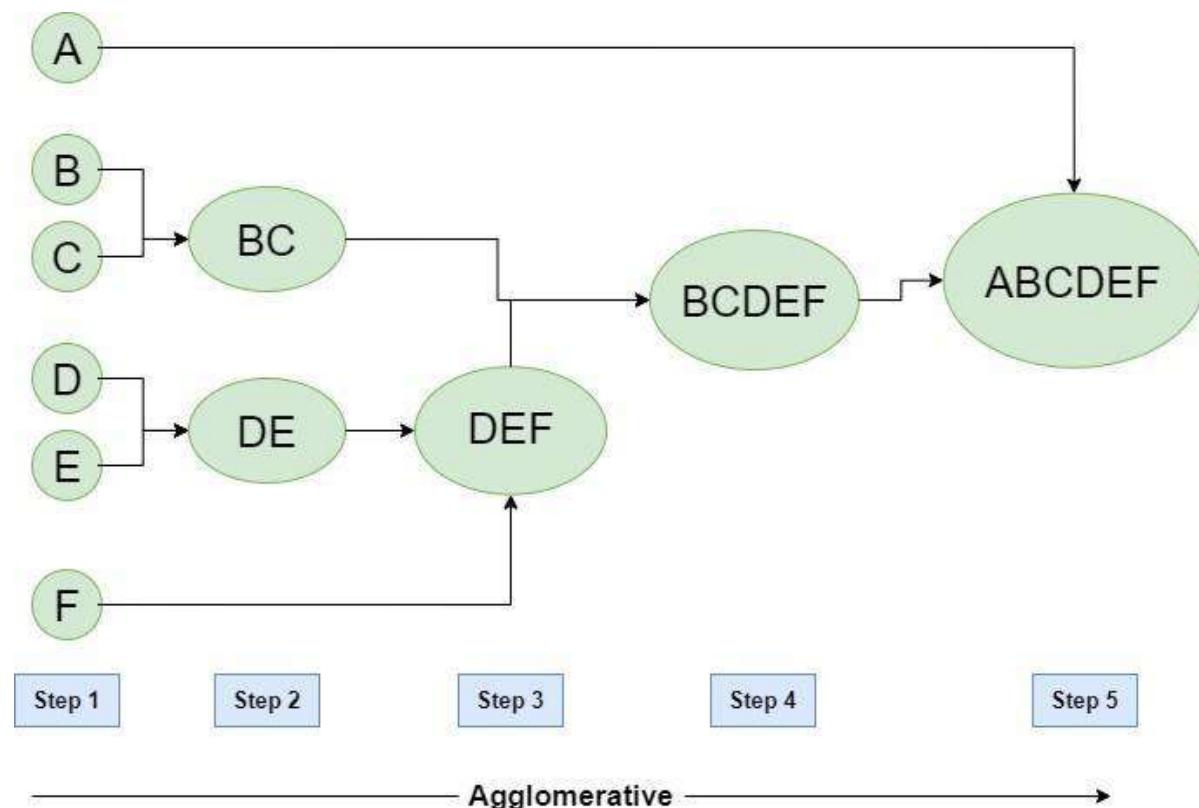
7. This method can handle different types of data and reveal the relationships among the clusters. However, it can have high computational cost and results can be sensitive to some conditions.

Agglomerative: Initially consider every data point as an individual Cluster and at every step, merge the nearest pairs of the cluster. (It is a bottom-up method). At first, every dataset is considered an individual entity or cluster. At every iteration, the clusters merge with different clusters until one cluster is formed.

The algorithm for Agglomerative Hierarchical Clustering is:

- Calculate the similarity of one cluster with all the other clusters (calculate proximity matrix)
- Consider every data point as an individual cluster
- Merge the clusters which are highly similar or close to each other.
- Recalculate the proximity matrix for each cluster
- Repeat Steps 3 and 4 until only a single cluster remains.

Let's say we have six data points A, B, C, D, E, and F.



Step 1: Consider each alphabet as a single cluster and calculate the distance of one cluster from all the other clusters.

Step 2: In the second step comparable clusters are merged together to form a single cluster. Let's say cluster (B) and cluster (C) are very similar to each other therefore we merge them in the second step similarly to cluster (D) and (E) and at last, we get the clusters [(A), (BC), (DE), (F)]

Step 3: We recalculate the proximity according to the algorithm and merge the two nearest clusters([(DE), (F)]) together to form new clusters as [(A), (BC), (DEF)]

Step 4: Repeating the same process; The clusters DEF and BC are comparable and merged together to form a new cluster. We're now left with clusters [(A), (BCDEF)].

Step 5: At last the two remaining clusters are merged together to form a single cluster [(ABCDEF)].

Output:

```

+ Code + Text
[ ] import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
# for better readability:
import matplotlib
matplotlib.rcParams['font.size'] = 16
matplotlib.rcParams['figure.figsize'] = (12, 6)
matplotlib.rcParams['figure.facecolor'] = '#e0e0e0'
df = pd.read_csv('/content/inventory.csv')
df

transaction_id product_id product_name product_category product_brand warehouse_id warehouse_name warehouse_location city_id supplier_id supplier_name contact_id
0 1 P2859 Titan T27 watches Titan W7276 BombayX Mumbai C7597 S3535 Stewart, Clayton and Martinez C2483 001-I
1 2 P5937 Samsung QLED 21 tv Samsung W7001 DelhiCaps Delhi C7099 S7973 Carroll, Brady and Hancock C8454 601.44-I
2 3 P1553 Philips Iron 14 Electrical App Philips W1296 BombayZ Mumbai C7597 S2142 Singh-Johnson C8265
3 4 P7296 Nike Nex Wear Nike W8451 BangaloreS Bangalore C9293 S3115 Boyer Ltd C3118 (989)Y
4 5 P9813 US Polo Black Denim Wear US Polo W8571 KeralaSwap Kerala C6336 S8922 Davis Ltd C7283
...
195 196 P5145 LG hair dryer Electrical App LG W9621 GurgaonHDES Gurgaon C5676 S5845 Dixon and Sons C9036 001-I
196 197 P7799 Apple watch series 8 watches Apple W4993 HyderabadHub Hyderabad C5055 S7793 Booth LLC C7568

```

```

+ Code + Text
[ ] dataset = df[['product_id', 'number_of_boxes']]
dataset

product_id number_of_boxes
0 P2859 16
1 P5937 77
2 P1553 62
3 P7296 71
4 P9813 37
...
195 P5145 45
196 P7799 80
197 P3507 38
198 P3238 54
199 P4930 75
200 rows × 2 columns

```

```

+ Code + Text
[ ] X = np.array(dataset)

class Distance_computation_grid(object):
    def __init__(self):
        pass

    def compute_distance(self, samples):
        Distance_mat = np.zeros((len(samples), len(samples)))

        for i in range(Distance_mat.shape[0]):
            for j in range(Distance_mat.shape[0]):
                if i != j:
                    Distance_mat[i, j] = float(self.distance_calculate(samples[i], samples[j]))
                else:
                    Distance_mat[i, j] = 10**4
        return Distance_mat

    # For two samples
    def distance_calculate(self, sample1, sample2):
        dist = []
        for i in range(len(sample1)):
            for j in range(len(sample2)):
                try:
                    # Try to convert elements to float and calculate distance
                    dist.append(np.linalg.norm(np.array(float(sample1[i])) - np.array(float(sample2[j]))))
                except (ValueError, TypeError):
                    # Handle non-numeric or missing values (e.g., NaN) here
                    pass

        if not dist:
            return float('Inf') # Return a large value for cases with no valid numeric values
        return np.array(dist)

```

```
+ Code + Text
    pass

    if not dist:
        return float('inf') # Return a large value for cases with no valid numeric values
    return min(dist)

    # For one sample and one cluster
    def intersampledist(self, s1, s2):
        if str(type(s2[0])) != '<class \'list\'>':
            s2 = [s2]
        if str(type(s1[0])) != '<class \'list\'>':
            s1 = [s1]
        n = len(s1)
        m = len(s2)
        dist = []
        if n >= m:
            for i in range(n):
                for j in range(m):
                    if (len(s2[i]) >= len(s1[j])) and str(type(s2[i][0])) != '<class \'list\'>':
                        dist.append(self.interclusterdist(s2[i], s1[j]))
                    else:
                        dist.append(np.linalg.norm(np.array(s2[i]) - np.array(s1[j])))
        else:
            for i in range(m):
                for j in range(n):
                    if (len(s1[i]) >= len(s2[j])) and str(type(s1[i][0])) != '<class \'list\'>':
                        dist.append(self.interclusterdist(s1[i], s2[j]))
                    else:
                        dist.append(np.linalg.norm(np.array(s1[i]) - np.array(s2[j])))
        return min(dist)

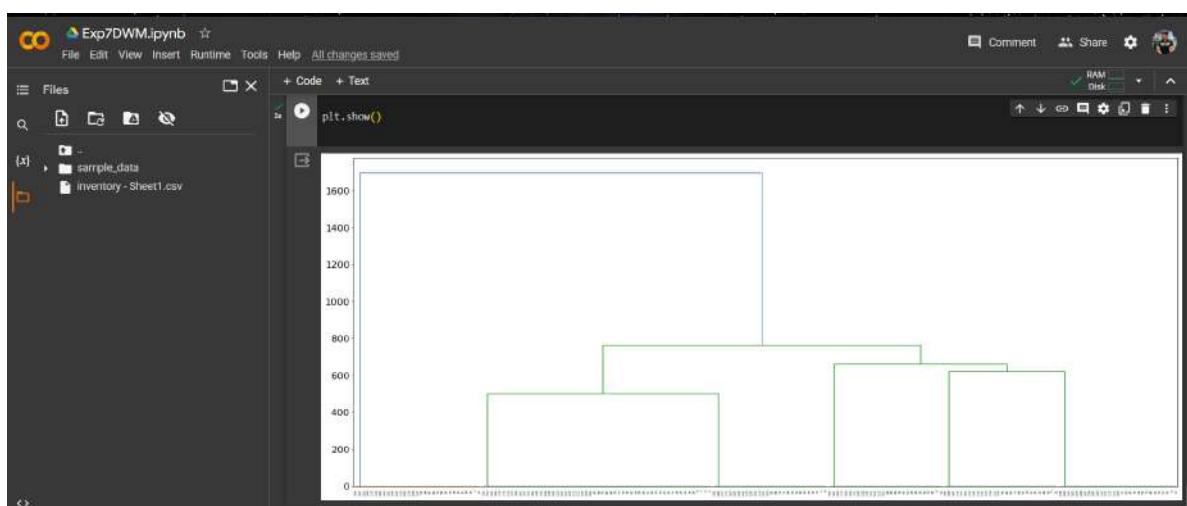
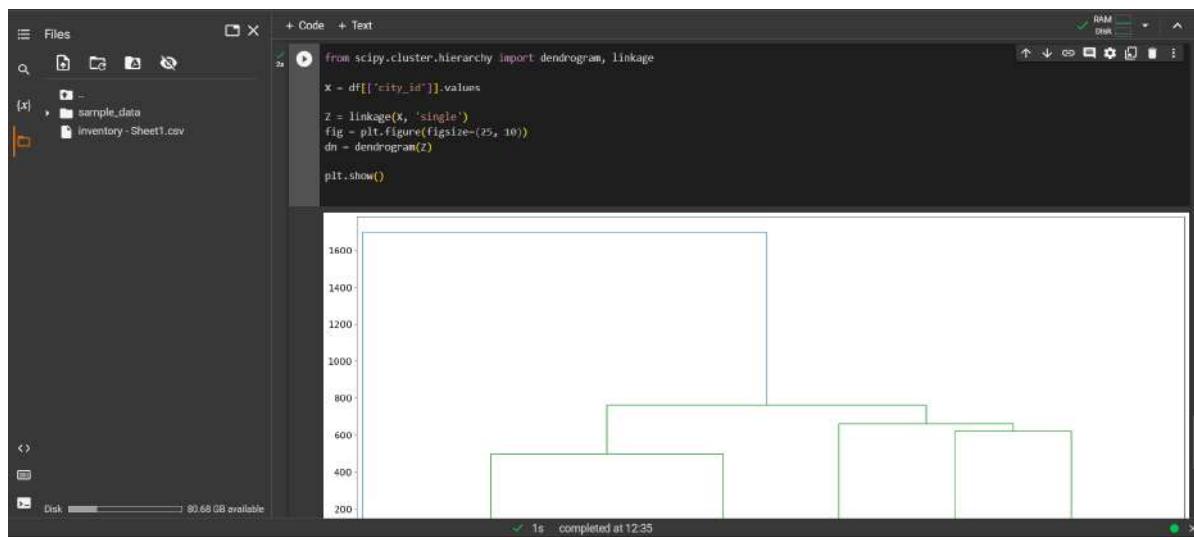
+ Code + Text
    # For two clusters
    def interclusterdist(self, cl, sample):
        dist = []
        for i in range(len(cl)):
            for j in range(len(sample)):
                try:
                    cl_value = float(cl[i])
                    sample_value = float(sample[j])
                    dist.append(np.linalg.norm(cl_value - sample_value))
                except (ValueError, TypeError):
                    # Handle non-numeric or missing values (e.g., NaN) here
                    pass
        if not dist:
            return float('inf') # Return a large value for non-numeric cases
        return min(dist)

progression = [{i} for i in range(x.shape[0])]
samples = [[list(x[i])] for i in range(x.shape[0])]
m = len(samples)
distcal = Distance_computation_grid()

while m > 1:
    print('Sample size before clustering:', m)
    Distance_mat = distcal.compute_distance(samples)
    sample_ind_needed = np.where(Distance_mat == Distance_mat.min())[0]
    value_to_add = samples.pop(sample_ind_needed[0])
    samples[sample_ind_needed[0]].append(value_to_add)
    print('Cluster Node 1:', progression[sample_ind_needed[0]])
    print('Cluster Node 2:', progression[sample_ind_needed[1]])
    progression[sample_ind_needed[0]].append(progression[sample_ind_needed[1]])
    progression[sample_ind_needed[0]] = [progression[sample_ind_needed[0]]]
    v = progression.pop(sample_ind_needed[1])
    m = len(samples)
    print('Progression (Current Sample):', progression)
    print('Cluster attained:', progression[sample_ind_needed[0]])
    print('Sample size after clustering:', m)
    print('\n')

    Sample size before clustering: 200
    Cluster Node 1: [0]
    Cluster Node 2: [1]
    Progression (Current Sample): [[[0, [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], [26]
    Cluster attained: [[0, [1]]]
    Sample size after clustering: 199

    Sample size before clustering: 199
    Cluster Node 1: [[0, [1]]]
    Cluster Node 2: [2]
    Progression (Current Sample): [[[0, [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], [26]
    Cluster attained: [[0, [1]], [2]]
    Sample size after clustering: 198
```



EXPERIMENT - 09

Aim: Implementation of Association Rule Mining algorithm (Apriori)

Theory:

Apriori algorithm refers to the algorithm which is used to calculate the association rules between objects. It means how two or more objects are related to one another. In other words, we can say that the apriori algorithm is an association rule leaning that analyzes that people who bought product A also bought product B.

The primary objective of the apriori algorithm is to create the association rule between different objects. The association rule describes how two or more objects are related to one another. Apriori algorithm is also called frequent pattern mining. Generally, you operate the Apriori algorithm on a database that consists of a huge number of transactions. Let's understand the apriori algorithm with the help of an example; suppose you go to Big Bazar and buy different products. It helps the customers buy their products with ease and increases the sales performance of the Big Bazar. In this tutorial, we will discuss the apriori algorithm with examples.

Apriori algorithm refers to an algorithm that is used in mining frequent products sets and relevant association rules. Generally, the apriori algorithm operates on a database containing a huge number of transactions. For example, the items customers buy at a Big Bazar.

Apriori algorithm helps the customers to buy their products with ease and increases the sales performance of the particular store.

Components of the Apriori algorithm

The given three components comprise the apriori algorithm.

1. Support
2. Confidence
3. Lift

Support

Support refers to the default popularity of any product. You find the support as a quotient of the division of the number of transactions comprising that product by the total number of transactions. Hence, we get

Support (Biscuits) = (Transactions relating to biscuits) / (Total transactions)

$$= 400/4000 = 10 \text{ percent.}$$

Confidence

Confidence refers to the possibility that the customers bought both biscuits and chocolates together. So, you need to divide the number of transactions that comprise both biscuits and chocolates by the total number of transactions to get confidence.

Hence,

Confidence = (Transactions relating to both biscuits and Chocolate) / (Total transactions involving Biscuits)

$$= 200/400$$

$$= 50 \text{ percent.}$$

It means that 50 per cent of customers who bought biscuits bought chocolates also.

Lift

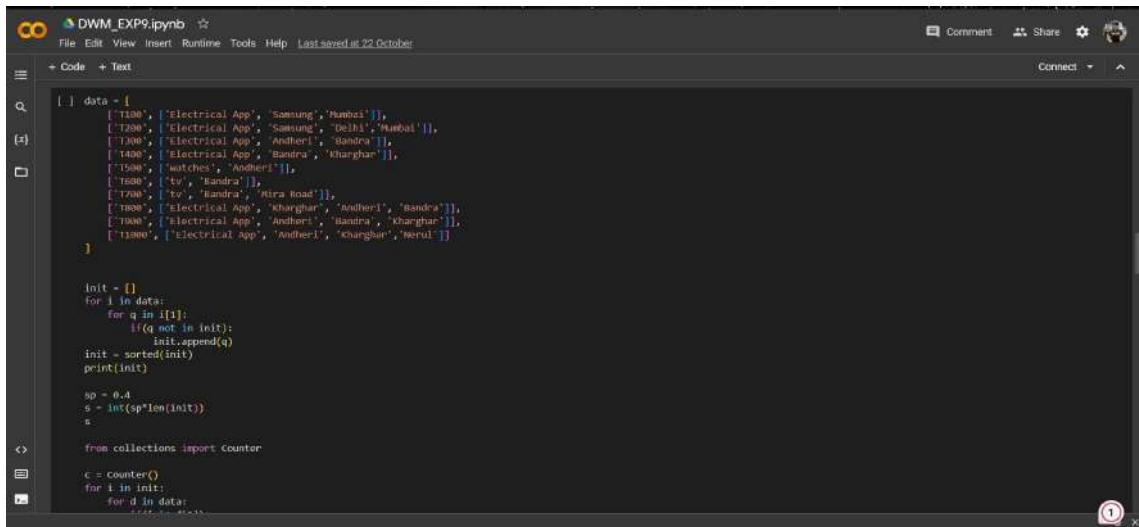
Consider the above example; lift refers to the increase in the ratio of the sale of chocolates when you sell biscuits. The mathematical equations of lift are given below.

Lift = (Confidence (Biscuits - chocolates)) / (Support (Biscuits))

$$= 50/10 = 5$$

It means that the probability of people buying both biscuits and chocolates together is five times more than that of purchasing the biscuits alone. If the lift value is below one, it requires that the people are unlikely to buy both items together. The larger the value, the better the combination.

Output:



```

DWM_EXP9.ipynb ☆
File Edit View Insert Runtime Tools Help Last saved at 22 October
Comment Share Connect ▾

+ Code + Text
☰ Q (x) □ <> C <> M

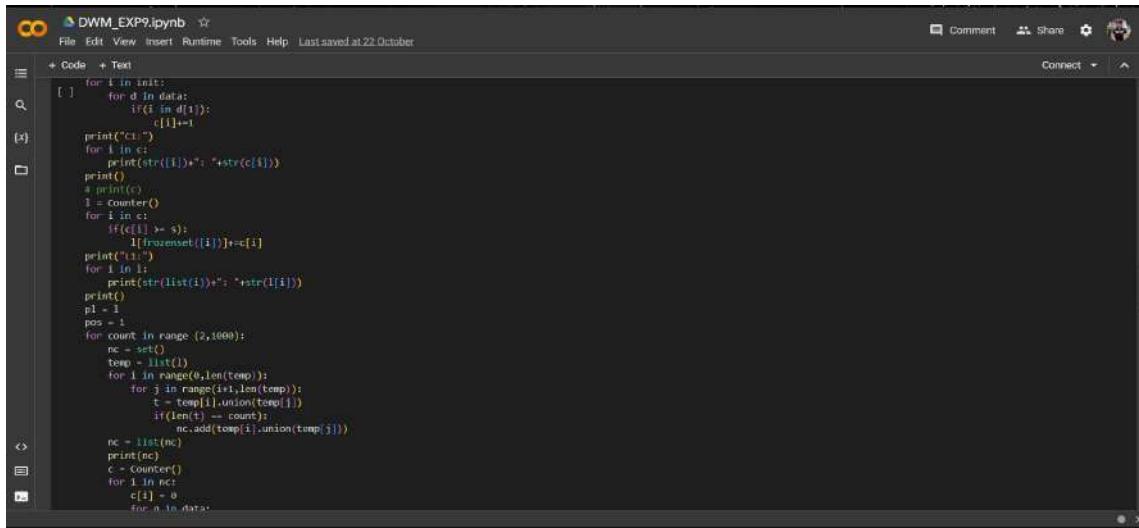
[ ] data = [
    ['T100', 'Electrical App', 'Samsung', 'Mumbai'],
    ['T200', 'Electrical App', 'Samsung', 'Delhi', 'Mumbai'],
    ['T300', 'Electrical App', 'Andheri', 'Bandra'],
    ['T400', 'Electrical App', 'Bandra', 'Kharghar'],
    ['T500', 'Watches', 'Andheri'],
    ['T600', 'tv', 'Bandra'],
    ['T700', 'tv', 'Bandra', 'Mira Road'],
    ['T800', 'Electrical App', 'Kharghar', 'Andheri', 'Bandra'],
    ['T900', 'Electrical App', 'Andheri', 'Bandra', 'Kharghar'],
    ['T1000', 'Electrical App', 'Andheri', 'Kharghar']
]

init = []
for i in data:
    for q in i[1]:
        if(q not in init):
            init.append(q)
print(init)

sp = 0.4
s = int(sp*len(init))
n = 

<> from collections import Counter
c = Counter()
for i in init:
    for d in data:
        if(d[1] == i):
            c[d[1]] += 1
print(c)
for i in c:
    print(str(i) + " : " + str(c[i]))
print()
l = Counter()
for i in init:
    l[i] = 0
    for j in data:
        if(j[1] == i):
            l[i] += 1
print(l)
for i in l:
    print(str(i) + " : " + str(l[i]))
print()
pl = 1
pos = 1
for count in range(1,1000):
    nc = set()
    temp = list()
    for i in range(0,len(temp)):
        for j in range(i+1,len(temp)):
            t = temp[i].union(temp[j])
            if(len(t) == count):
                nc.add(temp[i].union(temp[j]))
    nc = list(nc)
    print(nc)
    c = Counter()
    for i in nc:
        c[i] = 0
        for n in data:
            if(n[1] == i):
                c[i] += 1
    print(c)
    for i in nc:
        if(c[i] >= s):
            print(i)
    pl += 1
    pos = 1

```



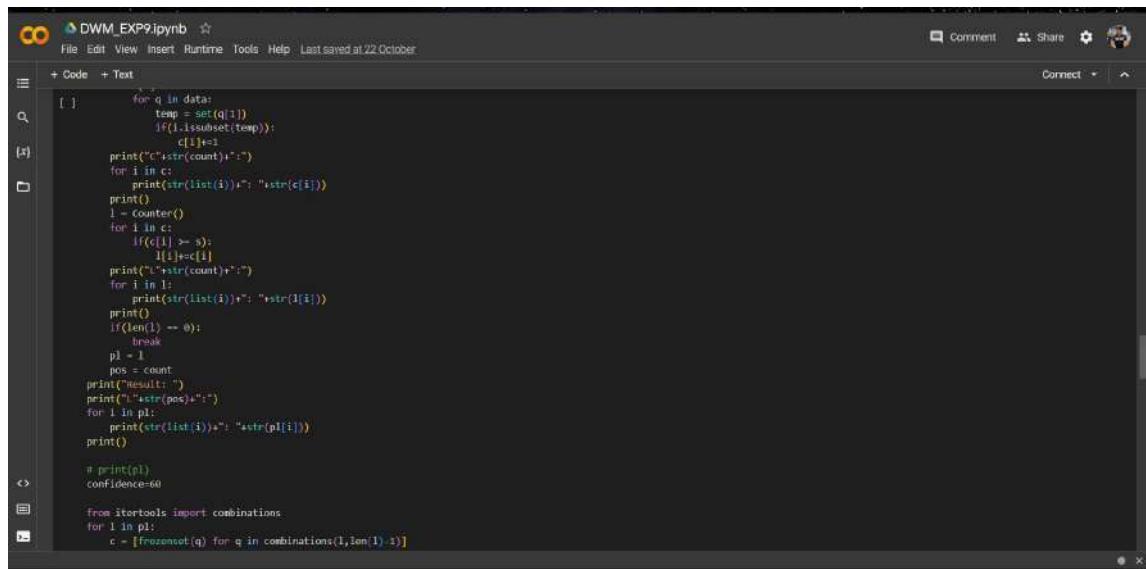
```

DWM_EXP9.ipynb ☆
File Edit View Insert Runtime Tools Help Last saved at 22 October
Comment Share Connect ▾

+ Code + Text
☰ Q (x) □ <> C <> M

[ ] for i in init:
    for d in data:
        if(d[1] == i):
            c[i] += 1
print("C1")
for i in c:
    print(str(i) + " : " + str(c[i]))
print()
l = Counter()
for i in init:
    l[i] = 0
    for j in data:
        if(j[1] == i):
            l[i] += 1
print(l)
for i in l:
    print(str(i) + " : " + str(l[i]))
print()
pl = 1
pos = 1
for count in range(1,1000):
    nc = set()
    temp = list()
    for i in range(0,len(temp)):
        for j in range(i+1,len(temp)):
            t = temp[i].union(temp[j])
            if(len(t) == count):
                nc.add(temp[i].union(temp[j]))
    nc = list(nc)
    print(nc)
    c = Counter()
    for i in nc:
        c[i] = 0
        for n in data:
            if(n[1] == i):
                c[i] += 1
    print(c)
    for i in nc:
        if(c[i] >= s):
            print(i)
    pl += 1
    pos = 1

```



```

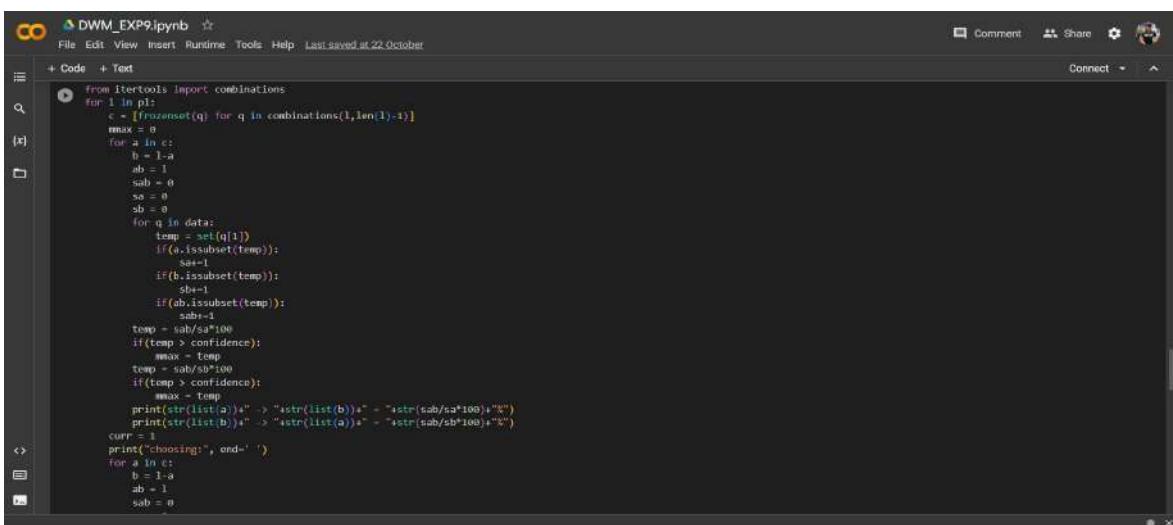
DWM_EXP9.ipynb ☆
File Edit View Insert Runtime Tools Help Last saved at 22 October
Comment Share Connect ▾

+ Code + Text
[ ] for q in data:
    temp = set(q[1])
    if(l.issubset(temp)):
        c[1]+=1
    print("c"+str(count)+"")
for i in c:
    print(str(list(i))+" "+str(c[i]))
print()
l = Counter()
for i in c:
    if(c[1]>=s):
        l[i]=c[1]
    print("l"+str(count)+"")
for i in l:
    print(str(list(i))+" "+str(l[i]))
print()
if(len(l) == 0):
    break
p1 = 1
pos = count
print("Result:")
print("c"+str(pos)+"")
for l in p1:
    print(str(list(l))+" "+str(p1[l]))
print()

# print(p1)
confidence=60

from itertools import combinations
for l in p1:
    c = [frozenset(q) for q in combinations(l,len(l)-1)]

```



```

DWM_EXP9.ipynb ☆
File Edit View Insert Runtime Tools Help Last saved at 22 October
Comment Share Connect ▾

+ Code + Text
[ ] from itertools import combinations
for l in p1:
    c = [frozenset(q) for q in combinations(l,len(l)-1)]
    rmax = 0
    for a in c:
        b = l-a
        sab = 1
        sa = 0
        sb = 0
        for q in data:
            temp = set(q[1])
            if(a.issubset(temp)):
                sa+=1
            if(b.issubset(temp)):
                sb+=1
            if(ab.issubset(temp)):
                sab-=1
        temp = sab/sa*100
        if(temp > confidence):
            rmax = temp
        temp = sab/sb*100
        if(temp > confidence):
            rmax = temp
        print(str(list(b))+" -> "+str(sab/sa*100)+"%")
        print(str(list(b))+" -> "+str(list(a))+" -> "+str(sab/sb*100)+"%")
    curr = 1
    print("choosing:", end=' ')
    for a in c:
        b = l-a
        sab = 1
        sa = 0
        sb = 0

```

```

+ Code + Text
print("choosing: ", end=" ")
for a in c:
    b = 1-a
    ab = 1
    sa = 0
    sb = 0
    for q in data:
        temp = set(q[1])
        if(a.issubset(temp)):
            sa+=1
        if(b.issubset(temp)):
            sb+=1
        if(ab.issubset(temp)):
            sab+=1
    temp = sab/sa*100
    if(temp >= confidence):
        print(curr, end = ' ')
    curr += 1
    temp = sab/sb*100
    if(temp >= confidence):
        print(curr, end = ' ')
    curr += 1
print()
print()

['Andheri', 'Bandra', 'Delhi', 'Electrical App', 'Kharhgar', 'Mira Road', 'Mumbai', 'Nerul', 'Samsung', 'tv', 'watches']

C1:
['Andheri']: 5
['Bandra']: 6
['Delhi']: 1
['Electrical App']: 7
['Kharhgar']: 4
['Mira Road']: 1
['Mumbai']: 2
['Nerul']: 1
['Samsung']: 2
['tv']: 2
['watches']: 1

L1:
['Andheri']: 5
['Bandra']: 6
['Electrical App']: 7
['Kharhgar']: 4

(frozenset(['Kharhgar', 'Andheri']), frozenset(['Bandra', 'Kharhgar']), frozenset(['Bandra', 'Andheri']), frozenset(['Kharhgar', 'Electrical App']), frozenset(['Electrical App', 'Andheri']), frozenset(['Kharhgar', 'Electrical App']), frozenset(['Bandra', 'Kharhgar', 'Electrical App']))

C2:
['Kharhgar', 'Andheri']: 3
['Bandra', 'Kharhgar']: 3
['Bandra', 'Andheri']: 3
['Kharhgar', 'Electrical App']: 4
['Electrical App', 'Andheri']: 4
['Bandra', 'Electrical App']: 4

L2:
['Kharhgar', 'Electrical App']: 4
['Electrical App', 'Andheri']: 4
['Bandra', 'Electrical App']: 4

(frozenset(['Bandra', 'Electrical App', 'Andheri']), frozenset(['Kharhgar', 'Electrical App', 'Andheri']), frozenset(['Bandra', 'Kharhgar', 'Electrical App']))

```

```

+ Code + Text
[ ] [ 'Kharhgar', 'Electrical App']; 4
[ ] [ 'Electrical App', 'Andheri']; 4
[ ] [ 'Bandra', 'Electrical App']; 4

(frozenset(['Bandra', 'Electrical App', 'Andheri']), frozenset(['Kharhgar', 'Electrical App', 'Andheri']), frozenset(['Bandra', 'Kharhgar', 'Electrical App']))

C3:
[ ] [ 'Bandra', 'Electrical App', 'Andheri']; 3
[ ] [ 'Kharhgar', 'Electrical App', 'Andheri']; 3
[ ] [ 'Bandra', 'Kharhgar', 'Electrical App']; 3

L3:
Result:
L2:
[ ] [ 'Kharhgar', 'Electrical App']; 4
[ ] [ 'Electrical App', 'Andheri']; 4
[ ] [ 'Bandra', 'Electrical App']; 4

[ 'Kharhgar'] -> [ 'Electrical App'] = 100.0%
[ 'Electrical App'] -> [ 'Kharhgar'] = 57.14285714285714%
[ 'Electrical App'] -> [ 'Kharhgar'] = 57.14285714285714%
[ 'Kharhgar'] -> [ 'Electrical App'] = 100.0%
choosing: 1 4

[ ] [ 'Electrical App'] -> [ 'Andheri'] = 57.14285714285714%
[ 'Andheri'] -> [ 'Electrical App'] = 80.0%
[ 'Andheri'] -> [ 'Electrical App'] = 80.0%
[ 'Electrical App'] -> [ 'Andheri'] = 57.14285714285714%
choosing: 2 3

[ ] [ 'Bandra'] -> [ 'Electrical App'] = 66.66666666666666%
[ 'Electrical App'] -> [ 'Bandra'] = 57.14285714285714%
[ 'Electrical App'] -> [ 'Bandra'] = 57.14285714285714%
[ 'Bandra'] -> [ 'Electrical App'] = 66.66666666666666%
choosing: 1 4

```

EXPERIMENT - 10

Aim: Implementation of Page rank/HITS algorithm

Theory:

PageRank and HITS (Hyperlink-Induced Topic Search) are two important algorithms used in web search and information retrieval to rank web pages. They both aim to determine the importance or relevance of web pages based on their connections and links to other pages. Here's an overview of both algorithms:

1. PageRank Algorithm:

PageRank is an algorithm developed by Larry Page and Sergey Brin, the co-founders of Google. It is used to rank web pages based on the importance of their incoming links. The basic idea behind PageRank is that important pages are more likely to be linked to by other pages.

- PageRank treats the web as a graph, where web pages are nodes, and hyperlinks between them are edges.
- Each page is assigned an initial PageRank value. Typically, all pages are given equal initial PageRank values.
- PageRank is calculated iteratively. At each iteration, each page distributes a fraction of its PageRank to the pages it links to.
- The PageRank of a page is determined by the PageRank of the pages that link to it. A link from a page with high PageRank is more valuable.
- Pages that have more incoming links and are linked from important pages have higher PageRank values.
- The PageRank of a page converges to a stable value after many iterations, and this value represents the page's importance.

Mathematically, the PageRank of a page 'P' is calculated as follows:

$$\begin{aligned} PR(P) = & (1 - d) + d * (PR(Q1)/L(Q1) + PR(Q2)/L(Q2) + \dots + \\ & PR(Qn)/L(Qn)) \end{aligned}$$

Where:

- $PR(P)$ is the PageRank of page P.
- d is a damping factor (usually set to 0.85) to account for the probability that a user might randomly jump to any page.
- $PR(Q1), PR(Q2), \dots, PR(Qn)$ are the PageRank values of pages that link to page P.

- $L(Q_1), L(Q_2), \dots, L(Q_n)$ are the number of outbound links on pages Q_1, Q_2, \dots, Q_n .

2. HITS (Hyperlink-Induced Topic Search) Algorithm:

HITS, also known as "Hubs and Authorities," is an algorithm that assesses web page authority and hub status. It was developed by Jon Kleinberg.

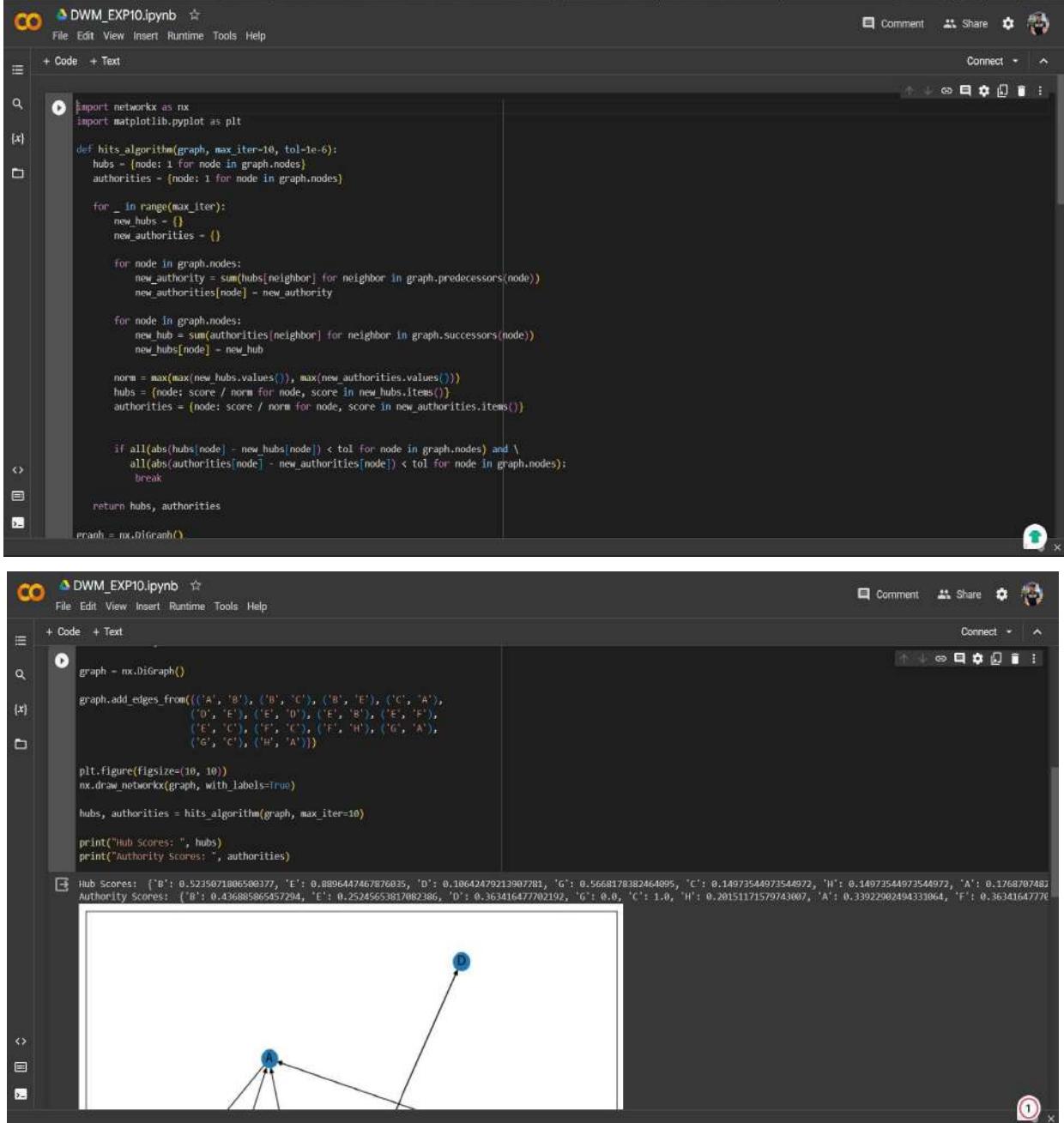
- In the HITS algorithm, each web page is assigned two scores: a hub score and an authority score.
- Hub pages are those that link to many authority pages, while authority pages are those linked to by many hub pages.
- The HITS algorithm iteratively updates the hub and authority scores until they converge to stable values.
- The final hub and authority scores provide a way to rank pages based on their role as hubs and authorities in the web graph.

Mathematically, the HITS algorithm works as follows:

1. Initialize all hub and authority scores to 1.
2. Iteratively update the hub and authority scores based on the links between pages.
3. Normalize the hub and authority scores to ensure they sum to 1.

HITS algorithm can be used to find not only the most authoritative pages but also to find good hubs that link to these authoritative pages.

Both PageRank and HITS are foundational algorithms in web search and information retrieval. PageRank is used by Google in its search engine, and HITS provides valuable insights into the link structure of the web. These algorithms have evolved and are part of a broader set of techniques used to rank web pages and discover relevant content on the internet.

Output:


The screenshot shows two code cells in a Jupyter Notebook interface. The top cell contains the implementation of the HITS algorithm, and the bottom cell visualizes the graph and displays the calculated hub and authority scores.

```

 1 #import networkx as nx
 2 import matplotlib.pyplot as plt
 3
 4 def hits_algorithm(graph, max_iter=10, tol=1e-6):
 5     hubs = {node: 1 for node in graph.nodes}
 6     authorities = {node: 1 for node in graph.nodes}
 7
 8     for _ in range(max_iter):
 9         new_hubs = {}
10         new_authorities = {}
11
12         for node in graph.nodes:
13             new_authority = sum(hubs[neighbor] for neighbor in graph.predecessors(node))
14             new_authorities[node] = new_authority
15
16             for node in graph.nodes:
17                 new_hub = sum(authorities[neighbor] for neighbor in graph.successors(node))
18                 new_hubs[node] = new_hub
19
20         norm = max(max(new_hubs.values()), max(new_authorities.values()))
21         hubs = {node: score / norm for node, score in new_hubs.items()}
22         authorities = {node: score / norm for node, score in new_authorities.items()}
23
24         if all(abs(hubs[node] - new_hubs[node]) < tol for node in graph.nodes) and \
25             all(abs(authorities[node] - new_authorities[node]) < tol for node in graph.nodes):
26             break
27
28     return hubs, authorities
29
30 graph = nx.DiGraph()
  
```

The bottom cell contains the following code to create a directed graph and visualize it:

```

graph = nx.DiGraph()

graph.add_edges_from([('A', 'B'), ('B', 'C'), ('B', 'E'), ('C', 'A'),
                     ('D', 'E'), ('E', 'D'), ('E', 'B'), ('E', 'F'),
                     ('E', 'C'), ('F', 'C'), ('F', 'H'), ('G', 'A'),
                     ('G', 'C'), ('H', 'A')])

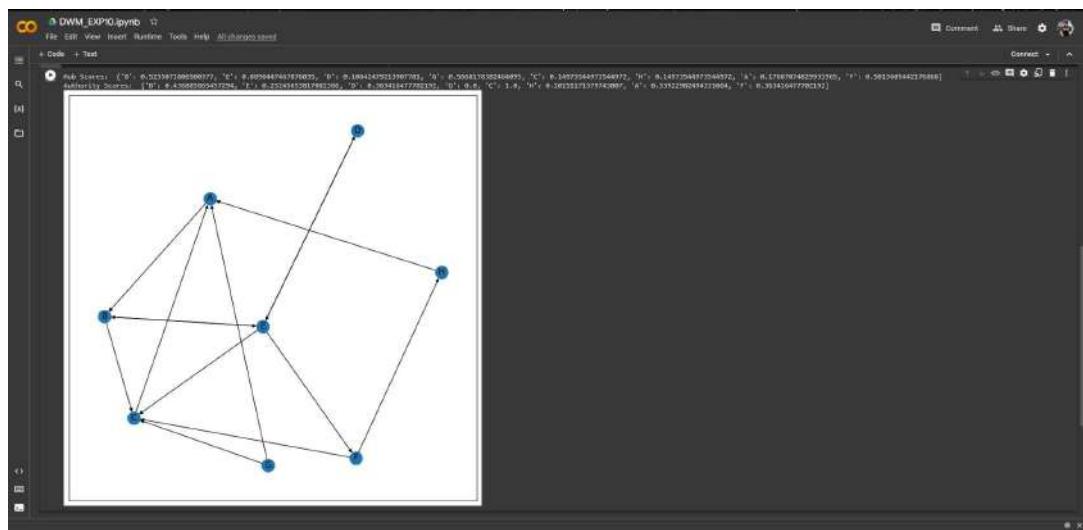
plt.figure(figsize=(10, 10))
nx.draw_networkx(graph, with_labels=True)

hubs, authorities = hits_algorithm(graph, max_iter=10)

print("Hub Scores: ", hubs)
print("Authority Scores: ", authorities)
  
```

The output of the visualization cell shows a directed graph with nodes A through H. Node E is highlighted with a blue circle, indicating it is a hub. Node D is also highlighted with a blue circle, indicating it is an authority. The edges are black lines connecting the nodes.

C32_Om_2103163_DWM



Assignment 1

What is Metadata? Explain data warehouse with Example.

- Data warehouse metadata is a piece of information stored in one or more special-purpose metadata repository that include:
 - 1) Information about data warehouse, their structure and their location.
 - 2) Information about refreshment of warehouse cleanup.
 - 3) Information regarding characteristics of components.
- Types of Metadata:
 - 1) Operational Metadata
 - In an enterprise, data for the data warehouse comes from various operational systems.
 - Different source systems contain different data structures having different field lengths and data types.
 - So the information of operational data source is given by operational metadata.

2) Extraction & Transformation Metadata

- This contains the information about the extraction of data from heterogeneous source system.
- It also contains information about data transformation in data staging area.

3) End-user Metadata

- ⇒
- This particular category provide the end user the flexibility of looking for information in their own way.
- For e.g.:

- 1) A library catalog may be considered metadata. The directory metadata consists of several components representing specific attribute of resource and each item can have one or more values. These components could be the name of author, name of publication, document, the publication date, etc.
- 2) The table of content and index in a book may be as treated metadata for the book.

Ques 23/10/2020

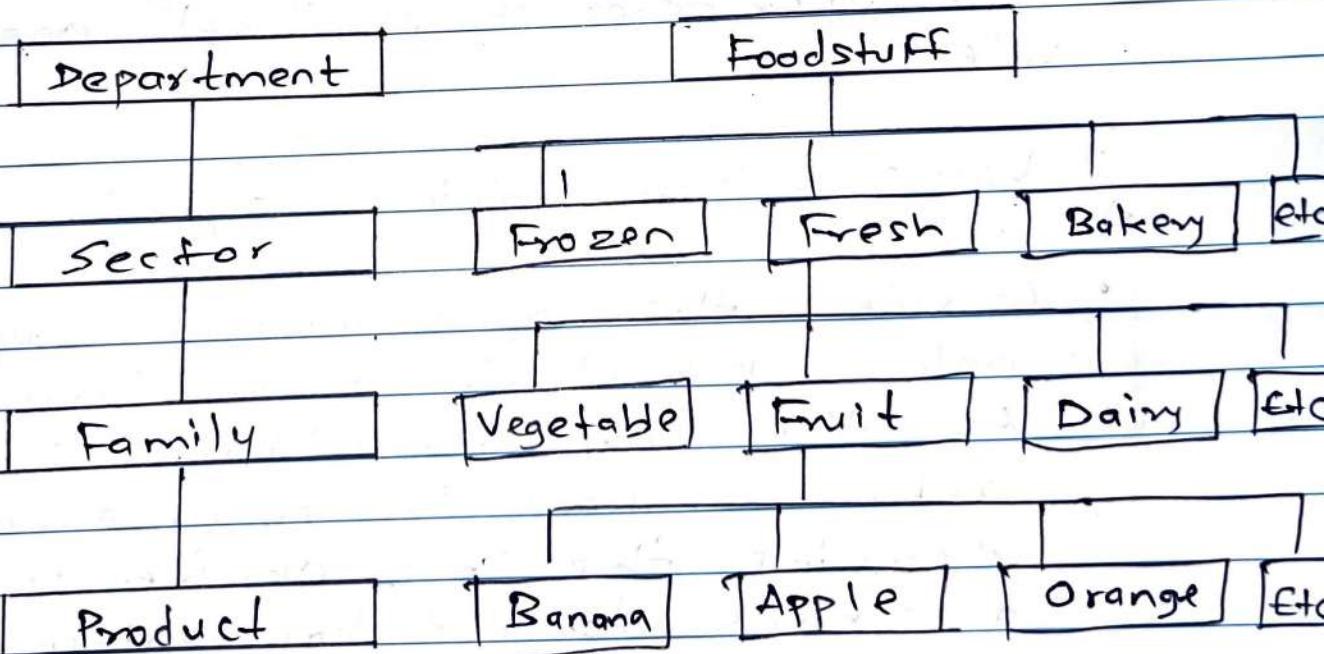
Assignment 2

Q.1 Write short note on

a) Multilevel Association Rules and Multidimensional association rules



- Items are always in the form of hierarchy.
- Items which are at leaf nodes are having lower support.
- An item can be either generalized or specialized as per the described hierarchy of that item and its level can be powerfully present in transactions.
- Rules which combines associations with hierarchy of concepts are called Multilevel Association Rules.



(Hierarchy of Concept)

- Support and confidence of multilevel association rules :

- The support and confidence of an item is affected due to its generalized or specialized value of attributes.
- The support of generalized item is more than the support of specialized item.
- Similarly the support of rules increases from specialized to generalized itemsets.
- If the support is below the threshold value then that rules become invalid.
- Confidence is not affected for general or specialized.

- Two approaches of multilevel association rules :

- 1) Using uniform minimum support for all levels
- Consider the same minimum support is below for all levels of hierarchy.
 - As the only one minimum support is set, so there is no necessary to examine the items of itemset whose ancestors do not have minimum support.

level 1

min-support = 5%

Milk
 (support = 10%)

Level 2

Cheese

min-support = 5%

(support = 6%)

Butter

(support = 4%)

- 2) Using reduced minimum support at lowest level



- Consider separate minimum support at each level of hierarchy.
- As the every level is having its own minimum support, the support at lowest level reduces.

Level 2

min-sup = 5%

Milk
 [support = 10%]

Level 2

min-sup = 3%

Cheese

[support = 6%]

Butter

[support = 9%]

- Multidimensional Association Rules:

- The rules contains two or more dimensions
 - predicates.
- Inter-dimension association rules
 - ⇒ The rule doesn't have any repeated predicate.
- Hybrid -dimension association rules
 - ⇒ The rule have many occurrences of same predicate , i.e. buys.

gender (x, "Male") \wedge buys (x, "TV")
 \Rightarrow buys (x, "DVD").

- Techniques of Mining MD Associations:

- 1) Using static discretization of quantitative attributes
 - ⇒ Using the concept hierarchy, discretizes the quantitative attributes.
- Convert the numeric values by ranges or categorical values.
- Mining is always more faster on data cube.
- Predicates sets are determined by the cells of n-dimensional cuboid.

2) Quantitative association rules:

- - Numeric attributes are dynamically discretized such that the confidence of the rules mined is maximized.
 - In quantitative association rule, the left hand side of rule contains 2-D quantitative attributes and right hand side of rules contains one categorical attribute.

Aquantitative \rightarrow Aquantitative² \Rightarrow Categorical

3) Distance-based association rules

- - This is a dynamic discretization process that considers the distance between data points.
 - This mining process has only two steps:
 - Perform clustering to find the interval of attributes involved.
 - Obtain association rules by searching for groups of clusters that occur together.

b) Web Usage Mining

- Web usage mining is the type of web mining which predicts about which pages are likely to be visited in near future based on the active user's behavior.
- Such pages can be pre-fetched to reduce access times.
- The usage data records the user's behaviour when the user browses or makes transactions on the web site in order to better understand and serve the needs of users or Web-based applications.
- Automatic discovery of patterns from one or more web servers.
- Organizations often generate and collect large volumes of data can help these organizations to determine:
 - The value of particular customers.
 - Cross marketing strategies across products.
 - The effectiveness of promotional campaigns, etc.
- The first web analysis tools simply provide mechanisms to report user activity as recorded in the servers.
- Using such tools, it was possible to determine such information as:

- The number of accesses to the server.
- The times or time intervals of visits.
- The domain names and the URLs of users of the Web Server.

- A very little or no analysis is provided by these tools of the data relationships among the accessed files and the directories within web spaces.

Q.3 - The tools for discovery and analysis of patterns may be classified into following two categories:

- Pattern Discovery Tools
- Pattern Analysis Tools.

- Web server, web proxies and client applications can quite easily capture web usage data.

- Web server log: It is a file that is created by the server to record all the activities it performs.

- For example, when a user enters URL into the browser's address bar or requests by clicking on a link.

Q3
10/29