

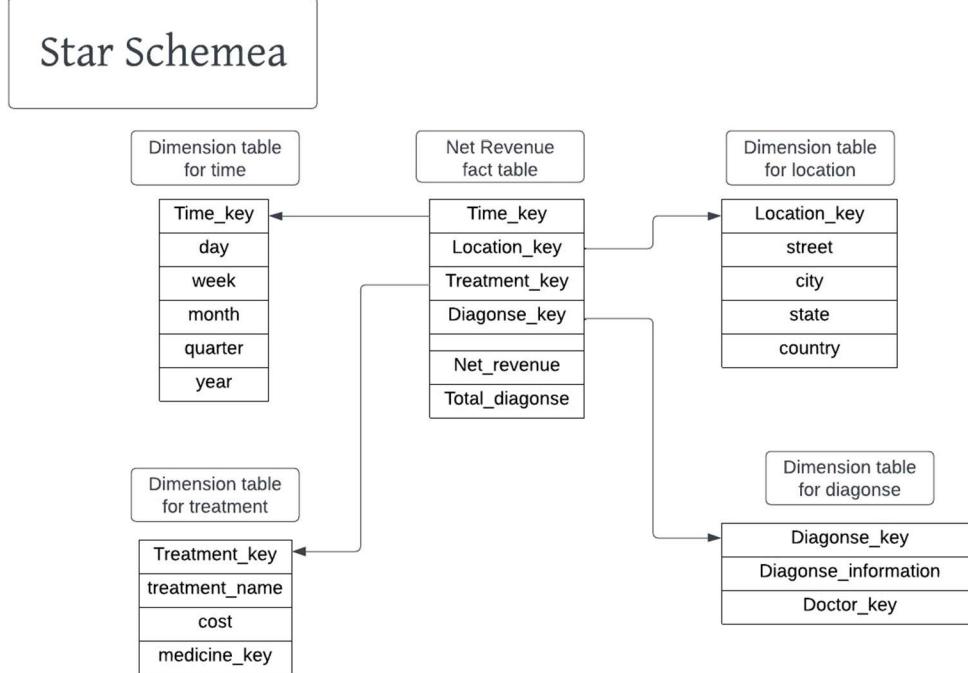
Topic: Dental Clinical Data System

Goal: Our goal is to analyze the clinical data through revenue generation of each branch through Treatment analysis of each patient in every branch

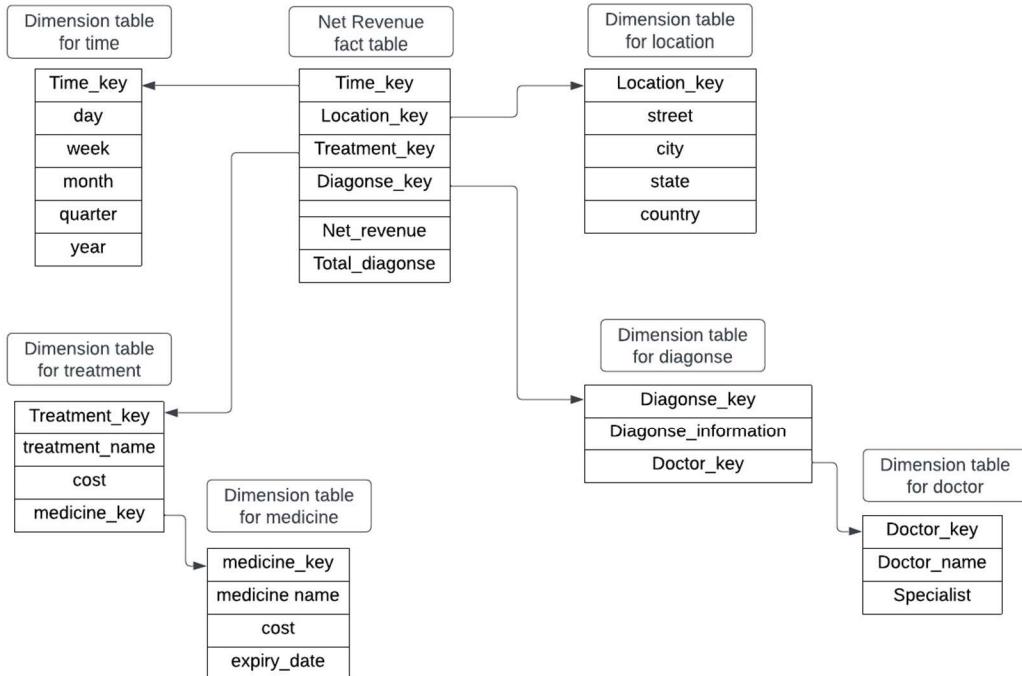
We are creating dental clinical system to keep a proper track of patients under services. We also need them for the follow up sessions of every patient. This system will increase the patients under our services at different location thus increasing the profit. And we can also get to know the profit of different branches. Thus, to achieve the goal of better decision-making using dimensions like Time, Location, Treatment, Diagnosis and facts included are: Revenue generated, Profit, Expense, No. of treatments.

Information Diagram:

Business Process: Treatment Data Analysis			
Time	Location	Treatment	Disease Diagnosed
Day	Street	Treatment key	Disease name
Week	City	Treatment name	Disease information
Month	State	Treatment cost	Doctor key
Quarter	Country	Medicines used	
years			
Facts: Revenue generated, Profit, Expense, No. of treatments.			



Snowflake Schema



1. Justify whether slowly changing dimension modeling is required for the problem selected?

Ans:

Slowly Changing Dimension (SCD) Modeling: SCD modeling is typically used in data warehousing and business intelligence when you have data that changes over time, such as historical customer addresses or product prices. The choice of whether to implement SCD modeling for a particular problem depends on the nature of the problem. Here are some justifications for and against using SCD modeling:

Justification for SCD Modeling:

- **Historical Analysis:** If your problem requires historical analysis, SCD modeling is crucial. For instance, analyzing how customer demographics change over time or tracking product price fluctuations over the years.
- **Auditing and Compliance:** When you need to maintain an audit trail and ensure compliance with historical data changes, SCD modeling helps you keep a record of all changes.

Justification against SCD Modeling:

- **Static Data:** If the data in your problem domain doesn't change over time, there's no need for SCD modeling. For example, if you're dealing with reference data that remains constant, like a list of countries and their attributes, you don't need SCD modeling.

2. Justify why Fact-less fact table modeling is not required for the problem selected? Give the examples of fact less fact tables with its type.

Ans:

Fact-less Fact Table Modeling: Fact-less fact tables are used to represent events or facts that don't contain any measures or numeric values. These tables are primarily used to establish relationships between dimensions. Whether fact-less fact tables are required depends on the problem you're dealing with. Here are some justifications for not using fact-less fact table modeling:

Justification against Fact-less Fact Table Modeling:

- **Lack of Relationships:** If your problem doesn't involve establishing relationships between dimensions, there's no need for fact-less fact tables. For instance, in a simple sales reporting system, you might not need fact-less fact tables if you're only interested in sales figures without intricate dimension relationships.

Examples of Fact-less Fact Tables:

- **Date Dimension Bridge:** A fact-less fact table could be used to establish relationships between date dimensions. For example, to track holidays, weekends, or special events for analytical purposes without containing any measures.
- **Student Enrollment:** In an educational context, a fact-less fact table could be used to track student enrollments, helping to identify enrollment trends and relationships between students and courses.

3. Justify your approach to deal with large dimension tables for the problem selected.

Ans:

Large Dimension Table Management: When dealing with large dimension tables, you may encounter performance and storage challenges. To justify your approach for managing large dimension tables in your problem domain, consider the following:

- **Data Partitioning:** Use data partitioning techniques to break down large dimension tables into smaller, more manageable pieces. For instance, you can partition a time dimension table by year or month.
- **Indexing:** Implement appropriate indexing strategies to optimize query performance for large dimensions. This could involve using clustered and non-clustered indexes.
- **Aggregations and Summarizations:** Consider pre-computing aggregations and summaries to reduce the complexity and size of your dimension tables, especially when dealing with historical data.
- **Caching:** Implement caching mechanisms to store frequently accessed parts of large dimension tables in memory for faster retrieval.

4. Justify the use of junk dimension and surrogate key for the problem selected.

Ans:

Junk Dimension and Surrogate Key Usage: Junk dimensions are used to combine several low-cardinality attributes into a single dimension to simplify the data warehouse structure. Surrogate keys are used as unique identifiers for dimension members. Their use depends on the specific

problem domain:

- **Junk Dimension Justification:** If your problem domain contains multiple low-cardinality attributes that don't need to be analyzed separately, but their combinations are meaningful, you can use a junk dimension. For example, in a retail environment, you could use a junk dimension to combine different discount types and payment methods.
- **Surrogate Key Justification:** Surrogate keys are valuable when you need to maintain a stable and unique identifier for dimension members, especially when dealing with slowly changing dimensions. These keys help in tracking and auditing changes to dimension attributes over time. For example, in a customer dimension, a surrogate key could be used to uniquely identify customers, even if their names or addresses change.

Name: Shivprasad CP

Roll No: 9696

Batch: A

Experiment 2

**Implement SQL queries for OLAP operations:
Part 1**

```
create table dim_time(time_key serial primary key, days date, weeks int, months int, quarter int, years int);
select * from dim_time;

insert into dim_time(days, weeks, months, quarter, years)
values ('2016-03-08', 33, 8, 3, 2016),
('2018-06-27', 27, 26, 6, 2018),
('2019-04-30', 30, 18, 4, 2019),
('2015-06-01', 1, 22, 6, 2015),
('2015-12-06', 6, 49, 12, 2015),
('2012-12-21', 21, 51, 12, 2012);

create table dim_location(location_key serial primary key, street varchar, city varchar, states varchar, country varchar);
select * from dim_location;

insert into dim_location(street, city, states, country)
values ('gully chowl', 'Navi Mumbai', 'Maharashtra', 'India'),
('Film city road', 'Mumbai', 'Maharashtra', 'India'),
('Sector 24', 'Noida', 'U.P', 'India'),
('Aul market road', 'Patamundi', 'Orrisa', 'India'),
('Shivaji Chowk', 'Pune', 'Maharashtra', 'India');

insert into dim_location(street, city, states, country)
values ('New market', 'Noida', 'U.P', 'India');

create table dim_treatment(treatment_key serial primary key, treatment_name varchar, costs int, medicine_info varchar);
select * from dim_treatment;

insert into dim_treatment(treatment_name, costs, medicine_info)
values('Root Canal', 2000, 'Yes'),
('Braces', 2500, 'No'),
('Teeth whitening', 3000, 'Yes'),
('Root Canal', 2500, 'Yes'),
('Wisdom extract', 4000, 'No'),
('Root Canal', 3000, 'Yes');
```

```

create table dim_diagonse(diagonse_key serial primary key, diagonse_info varchar, Doctor_key int);
select * from dim_diagonse;

insert into dim_diagonse(diagonse_info, Doctor_key)
values('Root Canal', 20),
('Braces', 21),
('Teeth whitening', 40),
('Root Canal', 72),
('Wisdom extract', 34),
('Root Canal', 40);

create table fact_revenue(time_key int references dim_time(time_key),
                         location_key int references dim_location(location_key),
                         treatment_key int references dim_treatment(treatment_key),
                         diagonse_key int references dim_diagonse(diagonse_key),
                         net_revenue decimal(19,4), total_diagonse int,
                         primary key(time_key, location_key, treatment_key, diagonse_key))
                         );

insert into fact_revenue values(1,1,1,1,20000.50, 10),
(2,2,2,10000.50, 5),
(3,3,3,15000.50, 8),
(4,4,4,22000.50, 12),
(5,5,5,24000.50, 15),
(6,6,6,18000.50, 9);

select * from fact_revenue;

drop table fact_revenue;

-- Roll Up Operations
select city, states, sum(net_revenue) from dim_location inner join
fact_revenue on dim_location.location_key = fact_revenue.location_key
group by rollup(states, city) order by states,city;

--Cube operations

select years, quarter, total_diagonse, sum(net_revenue) from fact_revenue natural inner join
dim_time
group by cube(years, quarter, total_diagonse);

--Slice Operations
select states, city, sum(net_revenue) from dim_location inner join fact_revenue on
fact_revenue.location_key = dim_location.location_key where states = 'Maharashtra' group by states,
city;

--Dice Operations
select states, city, sum(net_revenue) from dim_location inner join fact_revenue on
fact_revenue.location_key = dim_location.location_key where states = 'Maharashtra' and city =
'Mumbai'
group by states, city;

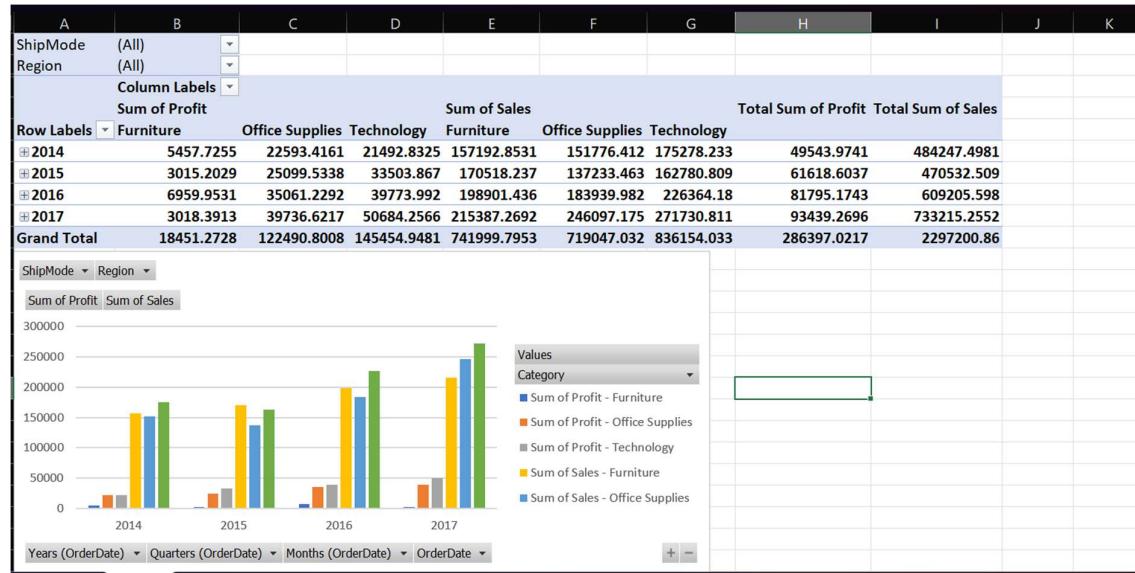
```

Part 2:

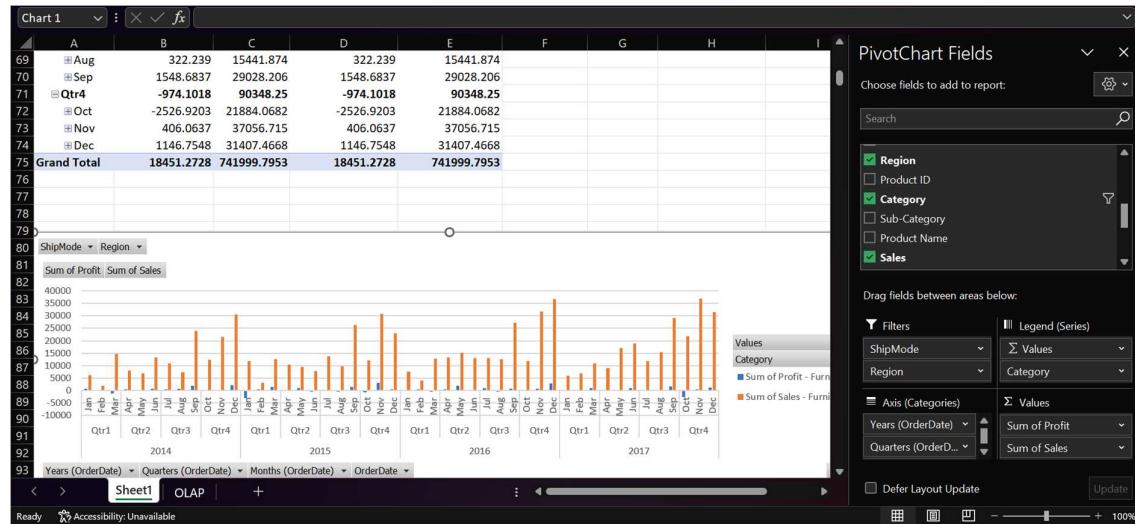
Perform data analysis with visualization for the following:

a) To view monthly, quarterly, yearly profit, sales of each category, region wise

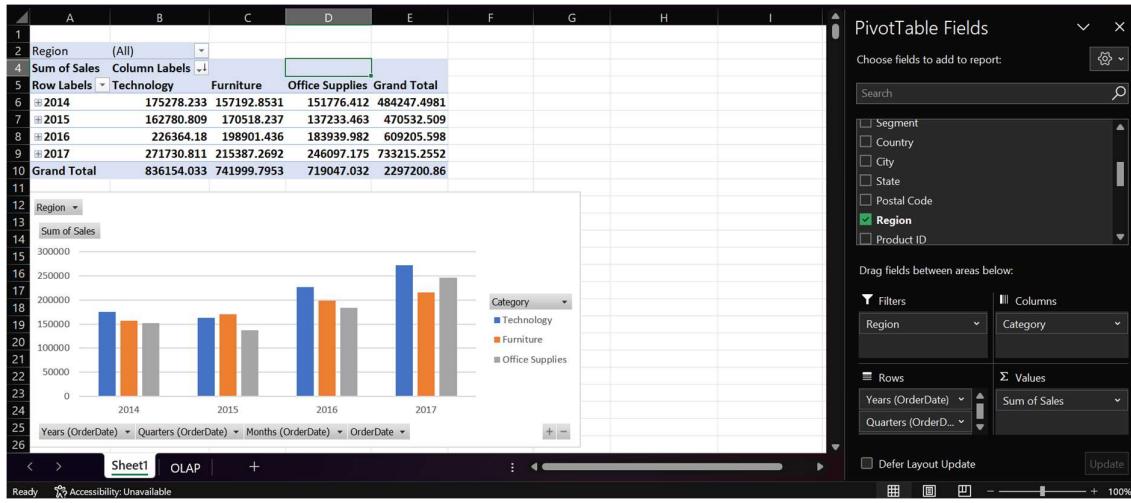
b) Comparison of sales and profit on various years.



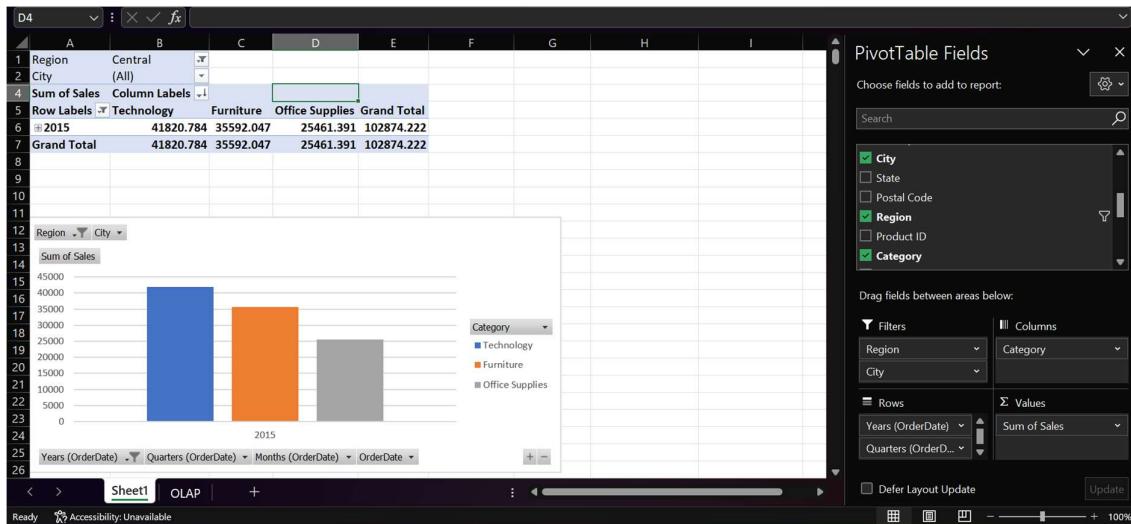
c) Comparison of sales in various months for product category =furniture.



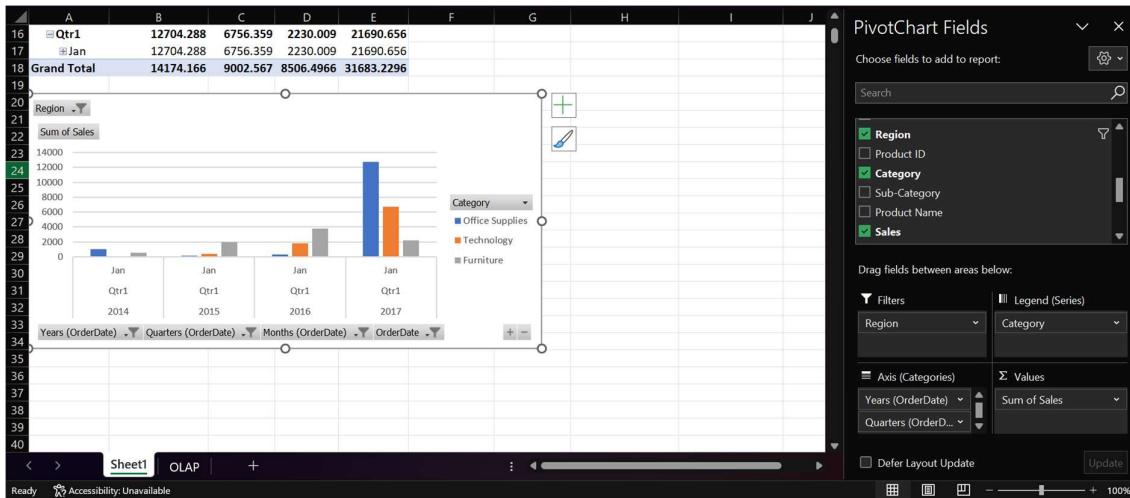
d) Need to know which product has more demand on which location?



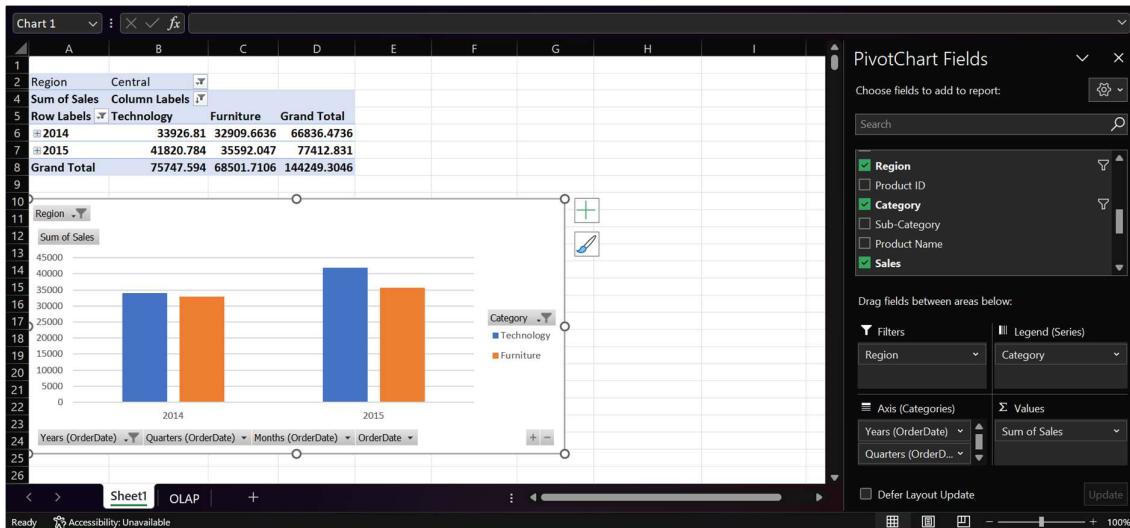
e) What is sale of product category wise, city wise for year=2015 ?



f) What is sales of product categories of Jan month of each year?



g) What is trend of sales on year 2014 and 2015 for product category furniture and technology?



Part 3:

List down the data Extraction Transformation and Loading processes applicable for your DW system.

Extract:

- Data Source Identification: Different sources of dental clinical data, which could include electronic health records, appointment systems, billing systems, patient information, etc.
- Data Extraction: Extract relevant data from the identified sources using appropriate methods, such as database queries, flat file exports, etc.
- Data Profiling: Analyze and profile the extracted data to understand its quality, structure, and potential issues.

Transform:

- Data Cleansing: Cleaning the extracted data by identifying missing values, inconsistencies, errors, and duplicates.

- Data Transformation: Convert data into a common format and standardize units, terminologies, and codes to ensure consistency.
- Data Integration: Integrate data from various sources into a unified format, considering data types, relationships, and hierarchies.
- Data Validation: Validate the transformed data to ensure that it meets the defined quality standards and business rules.

Loading:

- Staging: Store the cleaned and transformed data in a staging area, separate from the data warehouse, to facilitate further validation.
- Data Warehouse Loading: Load the validated and transformed data into the data warehouse.

Experiment No: 03

TE AI&DS

Date of Performance:

Roll No: 9696

Aim: Apply data Exploration techniques on given data to organize data (Tutorial)

CO3: Apply data exploration and Data preprocessing techniques to organize and prepare data for data mining

Rubrics for assessment of Experiment:

Sr. No	Parameters	Exceed Expectations(EE)	Meet Expectations (ME)	Below Expectations (BE)
1	Timeline (2)	Early or on time (2)	One session late (1)	More than one session late (0)
2	Preparedness (2)	Knows the basic theory related to the experiment very well. (2)	Managed to explain the theory related to the experiment. (1)	Not aware of the theory to the point. (1)
3	Effort (3)	Done expt on their own. (3)	Done expt with help from other. (2)	Just managed. (1)
4	Documentation(2)	Lab experiment is documented in proper format and maintained neatly. (2)	Documented in proper format but some formatting guidelines are missed. (1)	Experiments not written in proper format (0.5)
5	Result (1)	Specific conclusion.(1)	Partially specific conclusion. (0.5)	Not specific at all. (0)

Assessment Marks:

Timeline(2)	Preparedness(2)	Effort(3)	Documentation(2)	Result(1)	Total(10)

Tutorial 1. DWM Experiment 3.

1]

→ 18, 20, 21, 21, 24, 25, 25, 26, 27, 27, 29,
29, 29, 29, 35, 38, 38, 40, 40, 40, 40, 41, 45, 50,
51, 57, 75.

a) mode = ~~28~~ 33.15

median = 29.

b)

modes are 29 & 40

∴ Data's modality = bimodal

c) mid range = $(18 + 75)/2$.

mid range = 46.5

d)

→ Since the data is already sorted

$$Q_1 = (27+1)/4 = 7^{\text{th}} \text{ value}$$

$$Q_3 = 3 * (27+1)/4 = 21^{\text{st}} \text{ value.}$$

$$Q_1 = 25$$

$$Q_3 = 40$$

c) minimum = 18 $Q_1 = \frac{n+50}{100} = \frac{27+1}{2} = 13.5 \approx 14^{\text{th}}$
 $Q_1 = 25$ $Q_3 = 40$ $IQR = Q_3 - Q_1 = 40 - 25 = 15$
Median = 29 Lower limit = $Q_1 - 1.5 \times IQR$
 $Q_3 = 40$ $= 25 - 1.5 \times 15 = 25$
Maximum = 78 Upper limit = $Q_3 + 1.5 \times IQR$
 $= 40 + 1.5 \times 15 = 62.5$

g) A quantile-quantile ($Q-Q$) plot is used to compare the quantiles of data distribution with quantiles of a theoretical distribution like a normal distribution.
It helps to determine if the data follows a specific distribution. A quantile plot on the other hand, typically refers to a plot of the quantiles of data itself.

WZ

2].

age	frequency	c.f.
1 - 5	300	300
6 - 10	550	850
11 - 20	450	1300
21 - 50	1200	2500
51 - 80	800	3300
81 - 110	65	3365

$$n = 110, \quad n/2 = 110/2 = 55.$$

55th item lies in the class of 51 - 80.

$$\therefore L_1 = 51.$$

$$n = 110.$$

$$(\sum f_{\text{freq}})_L = 2500.$$

$$(\sum f_{\text{freq}})_m = 800.$$

$$\text{median} = L_1 + \left(\frac{n/2 - (\sum f_{\text{freq}})_L}{\sum f_{\text{freq}}}_m \right) \times \text{width}$$

$$= 51 + \left(\frac{55 - 2500}{800} \right) \times 29$$

$$= 51 - 29.07 = 21.93$$

3)
→ a)

$$\text{mean} = \frac{(23 + 23 + 27 + 27 + 39 + 41 + 47 + 49 + 50 + 52 + 54 + 54 + 56 + 57 + 58 + 58 + 60 + 61)}{18}$$

$$= 45.06$$

$$\text{Median} = 50$$

$$\text{standard deviation} = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$= 2(23 - 45.06)^2 + 2(27 - 45.06)^2 + \\ (39 - 45.06)^2 + (41 - 45.06)^2 \\ + (47 - 45.06)^2 + (49 - 45.06)^2 \\ + (50 - 45.06)^2 + \\ \cancel{50} - (52 - 45.06)^2 + 2(54 - 45.06)^2 \\ + (56 - 45.06)^2 + (57 - 45.06)^2.$$

$$\text{Standard deviation} = \sqrt{\frac{2970.94}{18}} = 12.84$$

Ans Jadi:

$$\text{Mean} = \frac{516.1}{18} = 28.67$$

medium:

$$N=18 \Rightarrow \frac{N}{2} = \frac{18}{2} = 9$$

Median :- $\frac{N+1}{2}$ (9th+10th) = $\frac{30.2+33.6}{2} = 31.9$

fat	fat - 28.67	(fat - 28.67) ²
9.5	-21.17	448.16
24.5	-4.17	17.38
5.8	-22.87	523.03
15.8	-12.87	165.67
33.4	4.73	22.37
28.9	-4.77	22.75
31.4	2.73	7.45
29.2	0.53	0.2809
30.2	1.53	2.3409
33.6	4.93	24.3049
44.5	15.83	280.83
28.8	0.13	0.0169
31.4	2.73	7.4529
33.2	4.83	20.52
32.1	3.43	11.76
34.9	6.23	38.81
40.2	11.53	132.94
35.7	7.1	50.41

$$\sigma = \sqrt{\frac{1746.1827}{18}}$$

$$\approx 9.84$$

$$\Sigma = 1746.1827$$

b) box plot

for age

$$Q_1 = \frac{N \times 25}{100} = \frac{18 \times 25}{100} = 4.5 \leftarrow 5^{\text{th}} = 39.$$

$$Q_2 = \frac{N \times 50}{100} = \frac{9}{2} = \frac{80+52}{2} = 51$$

$$Q_3 = \frac{N \times 75}{100} = \frac{13.8}{2} = 57$$

$$\text{IQR} = Q_3 - Q_1 = 57 - 39 = 18$$

$$\begin{aligned}\text{Lower limit} &= Q_1 - \text{IQR} \times 1.5 \\ &= 39 - 18 \times 1.5 \\ &= 12\end{aligned}$$

$$\begin{aligned}\text{Upper limit} &= Q_3 + \text{IQR} \times 1.5 \\ &= 57 + 18 \times 1.5 \\ &= 84\end{aligned}$$

No outliers

for fat:-

5.8, 7.5, 15.8, 23.9, 24.5, 28.8, 29.2, 30.2, 31.4, 31.9,
32.1, 33.2, 33.4336, 34.9, 35.7, 40.2, 49.5

$$Q_1 = 4.5 \approx 5^{\text{th}} = 24.5$$

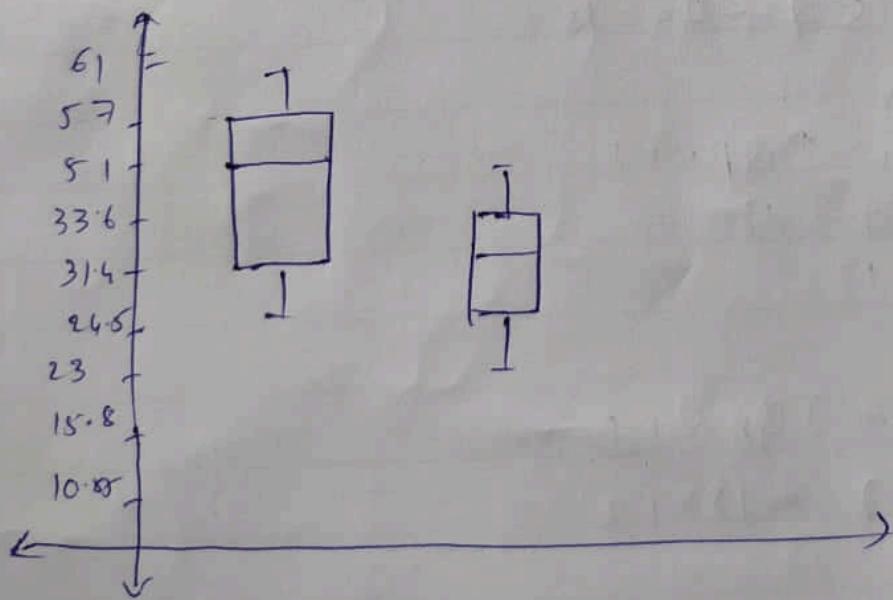
$$Q_2 = \frac{Q_1 + Q_3}{2} = \frac{24.5 + 31.4}{2} = 27.9$$

$$Q_3 = 14^{\text{th}} = 33.6$$

$$\text{IQR} = 33.6 - 24.5 \\ = 9.1$$

$$\text{Lower limit} = 24.5 - (9.1 \times 1.5) \\ = 10.85$$

$$\text{Upper limit} = 33.6 + (9.1 \times 1.5) \\ = 47.25$$



4)

→

Age	Ward	Gender	Tst	Health Progress	Fau
56	Ward A	F	P	Good	13 674
76	Ward B	M	N	Better	12343
23	Ward B	M	N	Beth	6542
47	Ward C	F	N	Bust	3459.

a) Nominal attributes - Ward.

$$d(P, j) = \frac{P - M}{P}$$

$$d(2, 1) = \frac{1 - 0}{1} = 1$$

$$d(3, 1) = \frac{1 - 0}{1} = 1$$

$$d(3, 2) = \frac{1 - 1}{1} = 0$$

$$d(4, 1) = \frac{1 - 1}{1} = 0$$

$$d(4, 2) = \frac{1 - 0}{1} = 1$$

$$d(4, 3) = \frac{1 - 0}{1} = 1$$

$$d(i, j) = \begin{bmatrix} 0 & & & \\ 1 & 0 & 0 & \\ 0 & 0 & 0 & \\ 0 & 1 & 0 & \end{bmatrix}$$

11) Asymmetric Test

- 1 $P \rightarrow i$
- 2 $N \rightarrow o$
- 3 $N \rightarrow o$
- 4 $N \rightarrow o$

		obj		
		1	1	0
P	1	1	0	s
	0	1	0	
		0	s	t

$$d(i, j) = \frac{a+s}{q+r+s}$$

$$d(2, 1) = \frac{1+0}{0+1+0} = 1$$

$$d(3, 1) = \frac{0+1}{0+0+1} = 1$$

$$d(3, 2) = \frac{0}{0} = 0$$

$$d(4, 1) = \frac{0+1}{0+0+1} = 1$$

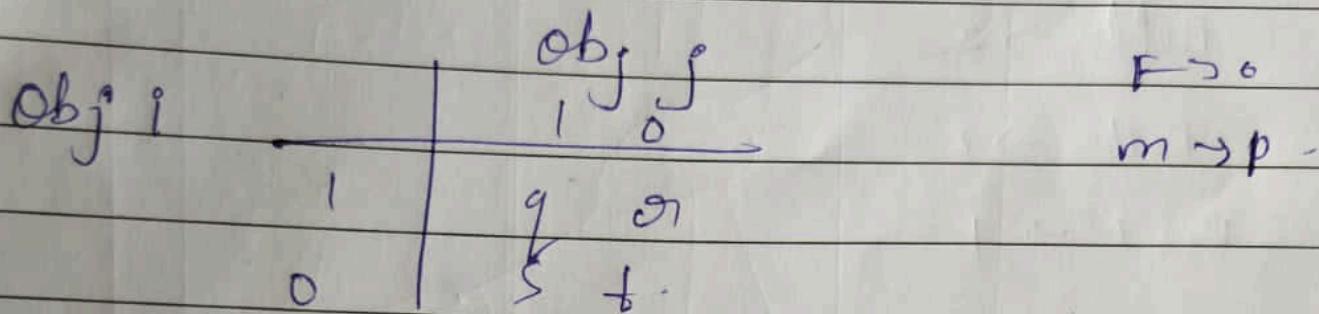
$$d(4, 2) = \frac{0}{0} = 0$$

$$d(4, 3) = \frac{0}{0} = 0$$

$$d(i, j) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

more similarity -

Symmetric.



gender

F → 0

m → 1

N → 1

P → 0

$$d(2,1) = \frac{1+0}{0+1+0+0} = 1$$

$$d(3,1) = 1$$

$$d(3,2) = 0$$

$$d(4,1) = 0$$

$$d(4,2) = 1$$

$$d(4,3) = 1$$

$$d(i,j) = \begin{cases} 0 & \\ 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{cases}$$

ordinal attribute

Health progress

Good $\rightarrow 3$

Better $\rightarrow 2$

| Better $\rightarrow 1$

Best $\rightarrow 1$

$$\text{if } \frac{\text{anif} - 1}{\text{anif} - 1} = 1$$

$$d(2,1) = 0.5$$

$$d(3,1) = 0.5$$

$$\text{Best} = \frac{1-1}{3-1} = 0$$

$$d(3,2) \approx$$

$$\text{Better} = \frac{2-1}{3-1} = \frac{1}{2} = 0.5$$

$$\text{Good} = \frac{3-1}{3-1} = 1$$

$$df \left[\begin{array}{ccc} 0 & & \\ 0.5 & 0 & \\ 0.5 & 0 & 0 \\ 1 & 0.5 & 0.5 \\ & & 0 \end{array} \right]$$

Numeric attributes

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}|$$

$$\begin{aligned} d(2, 1) &= |12343 - 13674| + |76 - 50| \\ &\approx 1331 + 20 \\ &= 1351 \end{aligned}$$

$$\begin{aligned} d(3, 1) &= |6542 - 13674| + |23 - 566| \\ &\approx 7132 + 33 \\ &= 7165 \end{aligned}$$

$$\begin{aligned} d(3, 2) &= |6542 - 12343| + |23 - 76| \\ &\approx 7132 + 33 \\ &= 7165 \end{aligned}$$

$$d(3, 2) = 5854$$

$$d(4, 1) = 10224$$

$$d(4, 2) = 8913$$

$$d(4, 3) = 3107$$

$$d(i, j) = \begin{bmatrix} 0 & & & \\ 1351 & 0 & & \\ 7165 & 5854 & 0 & \\ 10224 & 8913 & 3107 & 0 \end{bmatrix}$$

New Engineer.

5]
→

a)

$$\text{Euclidean Distance} = \sqrt{(32-20)^2 + (10-0)^2 + (40-31)^2 + (20-5)^2}$$
$$= 27.037$$

b) Manhattan Distance = $|32-20| + |10-0| + |40-31| + |20-5|$

$$= 46 //$$

Experiment No: 04	TE AI&DS
Date of Performance:	Roll No: 9696
Aim: Apply data exploration and Data preprocessing techniques to organize data for data mining	
Related CO3: Apply data exploration and Data preprocessing techniques to organize and prepare data for data mining	

Rubrics for assessment of Experiment:

Sr. No	Parameters	Exceed Expectations(EE)	Meet Expectations (ME)	Below Expectations (BE)
1	Timeline (2)	Early or on time (2)	One session late (1)	More than one session late (0)
2	Preparedness (2)	Knows the basic theory related to the experiment very well. (2)	Managed to explain the theory related to the experiment. (1)	Not aware of the theory to the point. (1)
3	Effort (3)	Done expt on their own. (3)	Done expt with help from other. (2)	Just managed. (1)
4	Documentation(2)	Lab experiment is documented in proper format and maintained neatly. (2)	Documented in proper format but some formatting guidelines are missed. (1)	Experiments not written in proper format (0.5)
5	Result (1)	Specific conclusion.(1)	Partially specific conclusion. (0.5)	Not specific at all. (0)

Assessment Marks:

Timeline(2)	Preparedness(2)	Effort(3)	Documentation(2)	Result(1)	Total(10)

Theory:

Data Exploration:

Measure of Central Tendency

A measure of central tendency (also referred to as measures of centre or central location) is a summary measure that attempts to describe a whole set of data with a single value that represents the middle or centre of its distribution.

There are three main measures of central tendency: the mean, the median and the mode. Each of these measures describes a different indication of the typical or central value in the distribution.

What is the mean?

The mean is the sum of the value of each observation in a dataset divided by the number of observations. This is also known as the arithmetic average.

Looking at the retirement age distribution again:

54, 54, 54, 55, 56, 57, 57, 58, 58, 60, 60

The mean is calculated by adding together all the values ($54+54+54+55+56+57+57+58+58+60+60 = 623$) and dividing by the number of observations (11) which equals 56.6 years.

The population mean is indicated by the Greek symbol μ (pronounced ‘mu’). When the mean is calculated on a distribution from a sample it is indicated by the symbol \bar{x} (pronounced X-bar).

What is the median?

The median is the *middle value* in distribution when the values are arranged in ascending or descending order.

The median divides the distribution in half (there are 50% of observations on either side of the median value). In a distribution with an odd number of observations, the median value is the middle value.

Looking at the retirement age distribution (which has 11 observations), the median is the middle value, which is 57 years:

54, 54, 54, 55, 56, 57, 57, 58, 58, 60, 60

When the distribution has an even number of observations, the median value is the mean of

the two middle values. In the following distribution, the two middle values are 56 and 57, therefore the median equals 56.5 years:

52, 54, 54, 54, 55, 56, 57, 57, 58, 58, 60, 60

The median cannot be identified for categorical nominal data, as it cannot be logically ordered.

Standard Deviation (Dispersion)

The **standard deviation** is a statistic that measures the dispersion of a dataset relative to its mean and is calculated as the square root of the variance.

Formula

$$\sigma = \sqrt{\frac{\sum(x_i - \mu)^2}{N}}$$

σ = population standard deviation

N = the size of the population

x_i = each value from the population

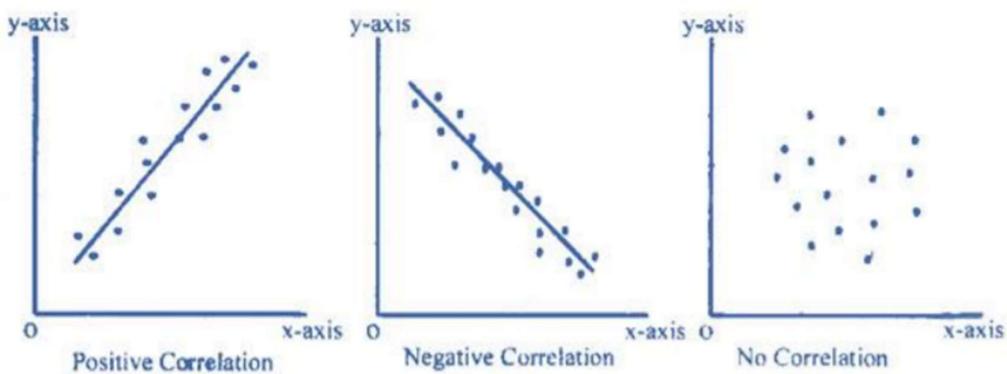
μ = the population mean

CORRELATION

Definition: The Correlation is a statistical tool used to measure the relationship between two or more variables, i.e. the degree to which the variables are associated with each other, such that the change in one is accompanied by the change in another.

Types of Correlation:

1. **Positive Correlation** A correlation in the same direction is called a positive correlation. If one variable increases the other also increases and when one variable decreases the other also decreases. For example, the length of an iron bar will increase as the temperature increases. • Price and Supply • Sales and Expenditure on Advertisement • Yield and Fertilizer Applied
2. **Negative Correlation** Correlation in the opposite direction is called a negative correlation. Here if one variable increases the other decreases and vice versa. For example, the volume of gas will decrease as the pressure increases, or the demand for a particular commodity increases as the price of such commodity decreases. Examples: • Price and Demand • Yield and Weed
3. **No Correlation or Zero Correlation** If there is no relationship between the two variables such that the value of one variable changes and the other variable remains constant, it is called no or zero correlation.



Methods of Determining Correlation:

- Pearson's Coefficient of Correlation.
- Spearman's Rank Correlation Coefficient;
- Kendall

Pearson r correlation

Pearson r correlation is the most widely used correlation statistic to measure the degree of the relationship between linearly related variables. For example, in the stock market, if we want to measure how two stocks are related to each other, Pearson r correlation is used to measure the degree of relationship between the two. The point-biserial correlation is conducted with the Pearson correlation formula except that one of the variables is dichotomous. The following formula is used to calculate the Pearson r correlation:

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

r_{xy} = Pearson r correlation coefficient between x and y

n = number of observations

x_i = value of x (for i th observation)

y_i = value of y (for i th observation)

Data preprocessing

Handling Missing Values

One of the important stages of data mining is preprocessing, where we prepare the data for mining. Real-world data tends to be incomplete, noisy, and inconsistent and an important task when preprocessing the data is to fill in missing values, smooth out noise and correct inconsistencies.

If we specifically look at dealing with missing data, there are several techniques that can be used. Choosing the right technique is a choice that depends on the problem domain — the data's domain and our goal for the data mining process.

So, there are several techniques which can be used for handling missing values.

1. Ignore the data row

This is usually done when the class label is missing (assuming your data mining goal is classification), or many attributes are missing from the row (not just one). However, you'll obviously get poor performance if the percentage of such rows is high.

2. Use a global constant to fill in for missing values

Decide on a new global constant value, like “unknown“, “N/A” or minus infinity, that will be used to fill all the missing values.

This technique is used because sometimes it just doesn't make sense to try and predict the missing value.

3. Use attribute mean

Replace missing values of an attribute with the mean (or median if its discrete) value for that attribute in the database.

4. Use a data mining algorithm to predict the most probable value

The value can be determined using regression, inference based tools using Bayesian formalism, decision trees, clustering algorithms (K-Mean\Median etc.).

Feature Scaling- Normalization and Standardization

Normalization

Normalization is a scaling technique in which values are shifted and rescaled so that they end up ranging between 0 and 1. It is also known as Min-Max scaling.

Here's the formula for normalization:

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Here, Xmax and Xmin are the maximum and the minimum values of the feature respectively.

- When the value of X is the minimum value in the column, the numerator will be 0, and hence X' is 0
- On the other hand, when the value of X is the maximum value in the column, the numerator is equal to the denominator and thus the value of X' is 1
- If the value of X is between the minimum and the maximum value, then the value of X' is between 0 and 1

Standardization (z-score normalization)

Standardization is another scaling technique where the values are centered around the mean with a unit standard deviation. This means that the mean of the attribute becomes zero and the resultant distribution has a unit standard deviation.

Here's the formula for standardization:

$$X' = \frac{X - \mu}{\sigma}$$

μ is the mean of the feature values and σ is the standard deviation of the feature values. Note that in this case, the values are not restricted to a particular range.

Normalization is good to use when you know that the distribution of your data does not follow a Gaussian distribution. This can be useful in algorithms that do not assume any distribution of the data like K-Nearest Neighbors and Neural Networks.

Standardization, on the other hand, can be helpful in cases where the data follows a Gaussian distribution. However, this does not have to be necessarily true. Also, unlike normalization, standardization does not have a bounding range. So, even if you have outliers in your data, they will not be affected by standardization.

Binning

Data binning, bucketing is a data pre-processing method used to minimize the effects of small observation errors. The original data values are divided into small intervals known as bins and then they are replaced by a general value calculated for that bin. This has a smoothing effect on the input data and may also reduce the chances of overfitting in the case of small datasets.

There are 2 methods of dividing data into bins:

Equal Frequency Binning: bins have an equal frequency.

Input: [5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215]

Output:

[5, 10, 11, 13]
[15, 35, 50, 55]
[72, 92, 204, 215]

Equal Width Binning : bins have equal width with a range of each bin are defined as [min + w], [min + 2w] [min + nw] where w = (max – min) / (no of bins).

Input: [5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215]
[5, 10, 11, 13, 15, 35, 50, 55, 72]

[92]

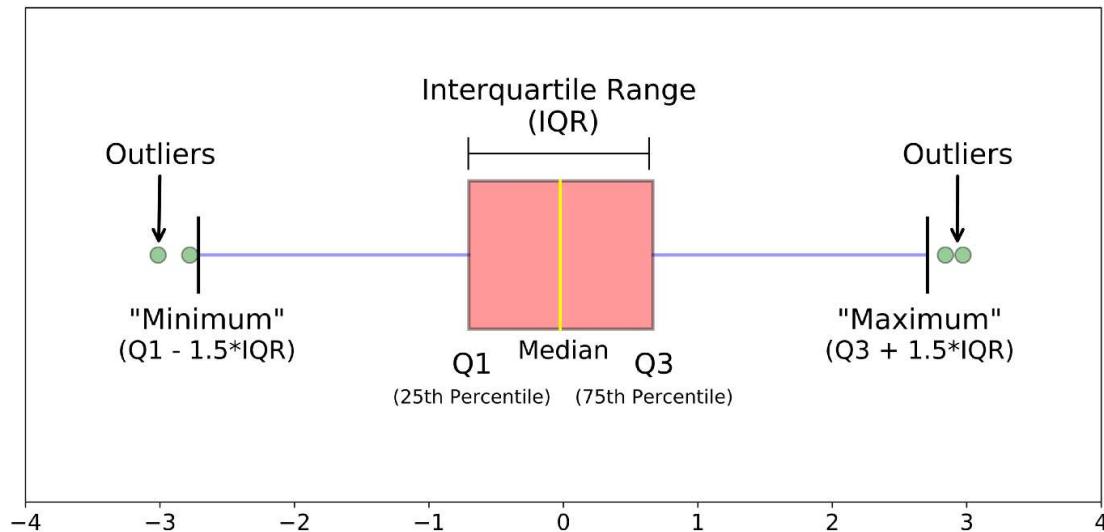
[204, 215]

Data visualization

BOX PLOT:

A boxplot is a standardized way of displaying the distribution of data based on a five number summary (“minimum”, first quartile (Q1), median, third quartile (Q3), and “maximum”). It

can tell you about your outliers and what their values are. It can also tell you if your data is symmetrical, how tightly your data is grouped, and if and how your data is skewed.

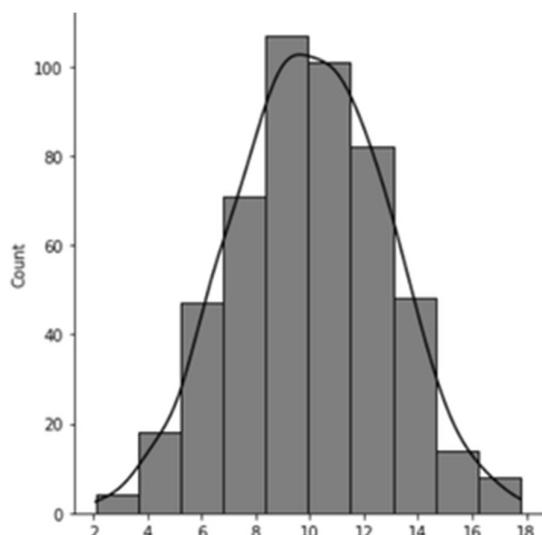


Histogram

A Histogram is a variation of a bar chart in which data values are grouped together and put into different classes. This grouping enables you to see how frequently data in each class occur in the dataset.

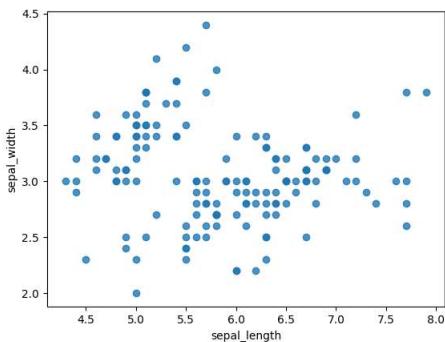
The histogram graphically shows the following:

- Frequency of different data points in the dataset.
- Location of the centre of data.
- The spread of dataset.



Scatter plot:

Scatter plots are a **commonly used data visualisation tool within data science**. They allow us to plot two numerical variables, as points, on a two-dimensional graph. From these plots, we can understand if there is a relationship between the two variables, and what the strength of that relationship is.



Implementation:

- 1) Load the libraries
- 2) Download the data set from classroom (realEstate_trans_full.csv)
- 3) Read the file –select appropriate file read function according to data type of file.
- 4) Display attributes in the data set-10 samples.
- 5) Describe the attributes name, count no of values, and find min, max, mean,data type, range, quartile, percentile.
- 6) Display correlation between any two attributes.
- 7) Download the dataset from classroom (realEstate_trans_less_dirty.csv)
- 8) Identify missing values fill them with mean.
- 9) Implement data normalization and standardization on real estate data.
- 10) Implement binning techniques on real estate data.

Code:

```
# data normalization
import pandas as pd

csv_read = pd.read_csv('D:/DWM/Exp 4/realEstate_trans_imputed.csv')

def normalize(col):
    return (col - col.min()) / (col.max() - col.min())

def standardize(col):
    return (col - col.mean()) / col.std()

cols = ['price_mean', 'sq_ft']

for col in cols:
    csv_read['n_' + col] = normalize(csv_read[col])
    csv_read['s_' + col] = standardize(csv_read[col])
```

```
with open('D:/DWM/Exp 4/realEstate_trans_imputed2.csv', 'w') as write_csv:  
    write_csv.write(csv_read.to_csv(sep=',', index=False))
```

data smoothing

```
import numpy as np  
import pandas as pd  
  
# read the data  
csv_read = pd.read_csv('D:/DWM/Exp 4/realEstate_trans_imputed2.csv')  
  
# create bins for the price that are based on the  
# linearly spaced range of the price values  
  
bins = np.linspace(  
    csv_read['price_mean'].min(),  
    csv_read['price_mean'].max(),  
    6  
)  
  
# and apply the bins to the data  
csv_read['b_price'] = np.digitize(  
    csv_read['price_mean'],  
    bins  
)  
  
# print out the counts for the bins  
counts_b = csv_read['b_price'].value_counts()  
print(counts_b.sort_index())  
  
# and write to a file  
  
with open('D:/DWM/Exp 4/realEstate_trans_binning.csv', 'w') as write_csv:  
    write_csv.write(csv_read.to_csv(sep=',', index=False))  
  
# OutPut :-  
# 1      350  
# 2      480  
# 3      118  
# 4       26  
# 5        4  
# 6        3  
# Name: b_price, dtype: int64
```

handling missing value:

```
import pandas as pd  
# read the data
```

```

csv_read =
pd.read_csv('/home/universe/Desktop/9689/realEstate_trans_less_dirty -
realEstate_trans_less_dirty.csv')

# impute mean in place of NaNs

csv_read['price_mean'] =
csv_read['price'].fillna(csv_read.groupby('zip')['price'].transform('mean'))

# impute median in place of NaNs
csv_read['price_median'] =
csv_read['price'].fillna(csv_read.groupby('zip')['price'].transform('median'))

# and write to a file
with open('/home/universe/Desktop/9689/realEstate_trans_imputed.csv', 'w') as
write_csv:
    write_csv.write(csv_read.to_csv(sep=',', index=False))

```

mean, median and correlation

```

import pandas as pd
import numpy as np

dataframe =
pd.read_csv('/home/universe/Desktop/9689/realEstate_trans_less_dirty -
realEstate_trans_less_dirty.csv')
# print(dataframe.columns)
dtype = (dataframe.dtypes)
# print(dtype)

# print(dataframe.head(10))
# print(dataframe[['beds', 'baths', 'sq_ft']].describe)

#Mean
print(dataframe[['beds', 'baths', 'sq_ft']].mean())
#Mode
print(dataframe[['beds', 'baths', 'sq_ft']].mode())
#Median
print(dataframe[['beds', 'baths', 'sq_ft']].median())
#Standrad Deviation
print(dataframe[['beds', 'baths', 'sq_ft']].std())

# Coorelation
print(dataframe[['beds', 'baths', 'sq_ft']].corr(method = "pearson"))

```

```
print("\nKendall")
print(dataframe[['beds','baths', 'sq_ft']].corr(method = "kendall"))

print("\nSpearman")
print(dataframe[['beds','baths', 'sq_ft']].corr(method = "spearman"))

dataframe['price_mean'] =
dataframe['price'].fillna(dataframe.groupby('zip')['price'].transform('mean'))
dataframe['price_median'] =
dataframe['price'].fillna(dataframe.groupby('zip')['price'].transform('median'))
```

Post lab:

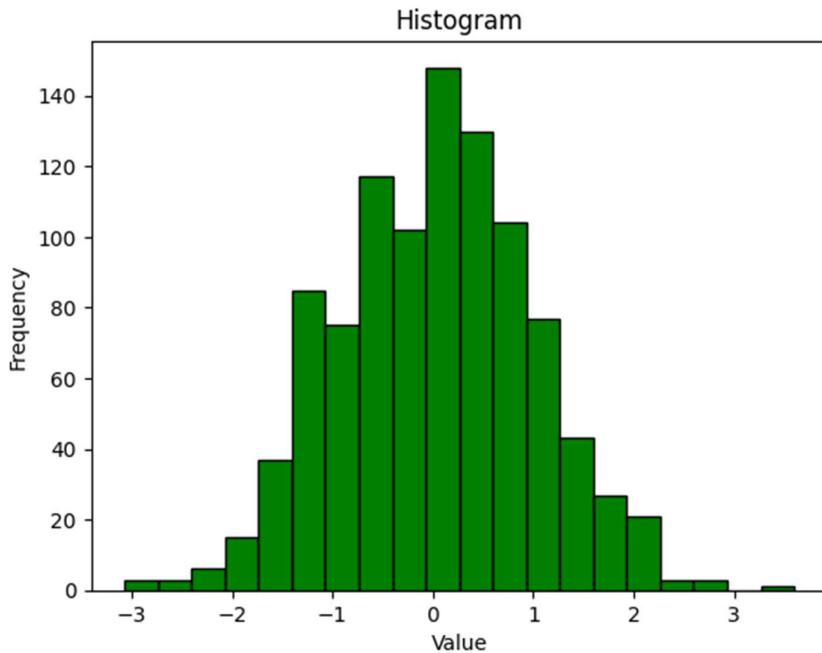
1. Apply data visualization for statistical description of data – in the form of histogram, scatter plot and box plot.

Code:

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

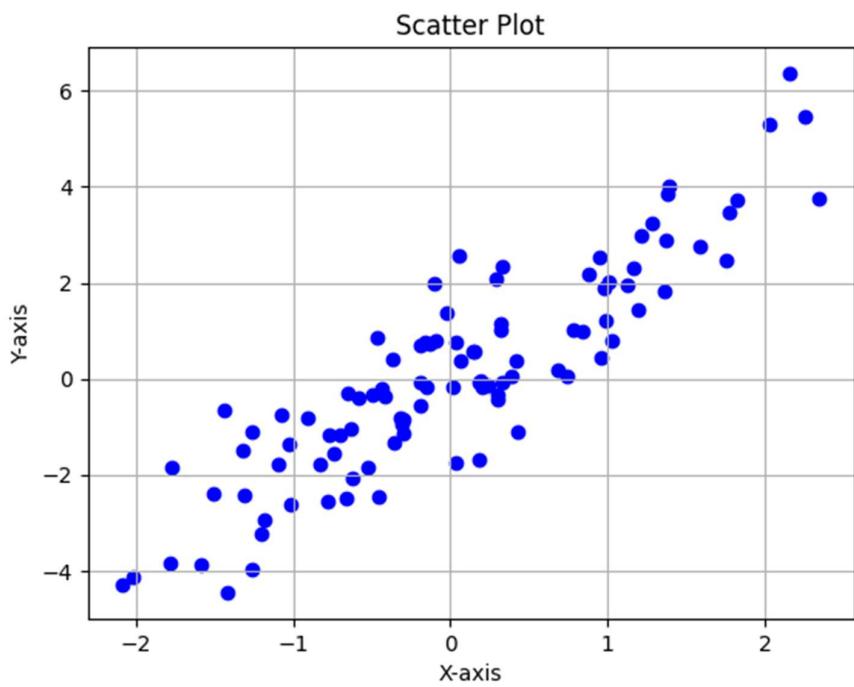
# Generate some example data
data = np.random.randn(1000) # Replace with your own dataset

# Create a histogram
plt.hist(data, bins=20, color='green', edgecolor='black')
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.title('Histogram')
plt.show()
```



```
# Scatter Plot
# Generate some example data
x = np.random.randn(100)
y = 2 * x + np.random.randn(100) # Replace with your own dataset

# Create a scatter plot
plt.scatter(x, y, c='blue', marker='o')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Scatter Plot')
plt.grid(True)
plt.show()
```



Conclusion:

I am able to understand the different data preprocessing technique and able to implement the code.

Experiment No: 05	TE AI&DS
Date of Performance:	Roll No: 9696
Aim: To Implement Naïve Bayesian classification algorithm to build classification model and predict class for unseen data.	
Related CO5: Implement Classification, Clustering and Association mining techniques to extract knowledge	
Objective: To learn how classification is used for Prediction.	

Rubrics for assessment of Experiment:

Sr. No	Parameters	Exceed Expectations(EE)	Meet Expectations (ME)	Below Expectations (BE)
1	Timeline (2)	Early or on time (2)	One session late (1)	More than one session late (0)
2	Preparedness (2)	Knows the basic theory related to the experiment very well. (2)	Managed to explain the theory related to the experiment. (1)	Not aware of the theory to the point. (1)
3	Effort (3)	Done expt on their own. (3)	Done expt with help from other. (2)	Just managed. (1)
4	Documentation(2)	Lab experiment is documented in proper format and maintained neatly. (2)	Documented in proper format but some formatting guidelines are missed. (1)	Experiments not written in proper format (0.5)
5	Result (1)	Specific conclusion.(1)	Partially specific conclusion. (0.5)	Not specific at all. (0)

Assessment Marks:

Timeline(2)	Preparedness(2)	Effort(3)	Documentation(2)	Result(1)	Total(10)

Theory:

Naive Bayes classifiers are a family of “probabilistic classifiers” based on [Bayes’ theorem](#) with strong independence between the features. They are among the simplest Bayesian network models and are capable of achieving high accuracy levels.

Bayes theorem states mathematically as:

$$P(A|B) = (P(B|A) * P(A)) / P(B)$$

where A and B are events and $P(B) \neq 0$.

$P(A|B)$ is a conditional probability: the probability of event A occurring given that B is true.
 $P(B|A)$ is also a conditional probability: the probability of event B occurring given that A is true.

$P(A)$ and $P(B)$ are the probabilities of observing A and B respectively without any given conditions.

A and B must be different events.

Bayes Theorem

- Based on prior knowledge of conditions that may be related to an event, Bayes theorem describes the probability of the event
- conditional probability can be found this way
- Assume we have a Hypothesis(H) and evidence(E),
According to Bayes theorem, the relationship between the probability of Hypothesis before getting the evidence represented as $P(H)$ and the probability of the hypothesis after getting the evidence represented as $P(H|E)$ is:

$$P(H|E) = P(E|H) * P(H) / P(E)$$

- **Prior probability** = $P(H)$ is the probability before getting the evidence
Posterior probability = $P(H|E)$ is the probability after getting evidence
- In general,

$$P(class|data) = (P(data|class) * P(class)) / P(data)$$

Bayes Theorem Example

Assume we have to find the probability of the randomly picked card to be king given that it is a face card.

There are 4 Kings in a Deck of Cards which implies that $P(King) = 4/52$
as all the Kings are face Cards so $P(Face|King) = 1$

there are 3 Face Cards in a Suit of 13 cards and there are 4 Suits in total so $P(Face) = 12/52$

Therefore,

$$P(King|face) = P(face|king) * P(king) / P(face) = 1/3$$

We can frame classification as a conditional classification problem with Bayes Theorem as follows:

- $P(y_i | x_1, x_2, \dots, x_n) = P(x_1, x_2, \dots, x_n | y_i) * P(y_i) / P(x_1, x_2, \dots, x_n)$
The prior $P(y_i)$ is easy to estimate from a dataset, but the conditional probability of the observation based on the class $P(x_1, x_2, \dots, x_n | y_i)$ is not feasible unless the number of examples is extraordinarily large, e.g. large enough to effectively estimate the probability distribution for all different possible combinations of values.

Example

```
#Weather Dataset
Outlook Temp Humidity Windy Play
Rainy Hot High f no
Rainy Hot High t no
Overcast Hot High f yes
Sunny Mild High f yes
Sunny Cool Normal f yes
Sunny Cool Normal t no
Overcast Cool Normal t yes
Rainy Mild High f no
Rainy Cool Normal f yes
Sunny Mild Normal f yes
Rainy Mild Normal t yes
Overcast Mild High t yes
Overcast Hot Normal f yes
Sunny Mild High t no
```

Calculate Prior Probability of Classes P(y)

```
#Frequency table
(Play=Yes) = 9/14 = 0.64
P(Play=No) = 5/14 = 0.36
```

Calculate the Likelihood Table for all features

```
#Likelihood Table#Outlook
```

	Play	Overcast	Rainy	Sunny
Yes	4/9	2/9	3/9	
No	0/5	3/5	2/5	
	—	—	—	
	4/14	5/14	5/14	

```
#Temp
```

	Play	Cool	Mild	Hot
Yes	3/9	4/9	2/9	
No	1/5	2/5	2/5	
	—	—	—	
	4/14	6/14	4/14	

```
#Humidity
```

Play High Normal

Yes 3/9 6/9

No 4/5 1/5

— —

7/14 7/14

#Windy

Play f t

Yes 6/9 3/9

No 2/5 3/5

— —

8/14 6/14

Now, Calculate Posterior Probability for each class using the Naive Bayesian equation. The Class with maximum probability is the outcome of the prediction.

Query: Whether Players will play or not when the weather conditions are [Outlook=Rainy, Temp=Mild, Humidity=Normal, Windy=t]?

Calculation of Posterior Probability:

Since Conditional independence of two random variables, A and B gave C holds just in case

$$P(A, B | C) = P(A | C) * P(B | C)$$

$$P(y=Yes | x) = P(Yes | Rainy, Mild, Normal, t)$$

$$P(Rainy, Mild, Normal, t | Yes) * P(Yes)$$

$$= \frac{P(Rainy, Mild, Normal, t)}{P(Rainy, Mild, Normal, t)}$$

$$P(Rainy | Yes) * P(Mild | Yes) * P(Normal | Yes) * P(t | Yes) * P(Yes)$$

$$=$$

$$P(Rainy) * P(Mild) * P(Normal) * P(t)$$

$$(2/9) * (4/9) * (6/9) * (3/9) * (9/14)$$

$$= \frac{(2/9) * (4/9) * (6/9) * (3/9) * (9/14)}{(5/14) * (6/14) * (7/14) * (6/14)}$$

```

= 0.43

P(y=No|x) = P(No|Rainy,Mild,Normal,t)

P(Rainy,Mild,Normal,t|No) * P(No)
=
_____
P(Rainy,Mild,Normal,t)

P(Rainy|No) * P(Mild|No) * P(Normal|No) * P(t|No) * P(No)
=
_____
P(Rainy) * P(Mild) * P(Normal) * P(t)

(3/5) * (2/5) * (1/5) * (3/5) * (5/14)
=
_____
(5/14) * (6/14) * (7/14) * (6/14)

= 0.31

```

Now, `P(Play=Yes|Rainy,Mild,Normal,t)` has the highest Posterior probability.

Algorithm/steps:

- 1) Load the dataset and convert it into list.
- 2) Separating the data as per class(0,1)-mp[0], mp[1]
- 3) define the test input: test = [2,1,0,1]
- 4) find the prior probability of each class
- 5) Find the conditional probability of test for each class-[yes,no]

```

probYes = 1 // probNO = 1
count = 0
total = 0
//find the length of test
for i in range(len(test)):
    count = 0
    total = 0
    for row in mp[1]: //mp[0] for NO
        if(test[i] == row[i]):
            count += 1
            total += 1
    probYes *= count/total //probNO*= count/total

```

- 6) Display the posterior probability of each class and find maximization

Code with output:

Code:-

```

import numpy as np
import pandas as pd

```

```

import matplotlib.pyplot as plt
import math

def accuracy_score(y_true, y_pred):
    """
        score = (y_true - y_pred) / len(y_true)
    """
    return round(float(sum(y_pred == y_true))/float(len(y_true)) * 100 ,2)

def pre_processing(df):
    """
        partitioning data into features and target
    """
    X = df.drop([df.columns[-1]], axis = 1)
    y = df[df.columns[-1]]

    return X, y

class NaiveBayes:

    """
        Bayes Theorem:
        Likelihood * Class prior probability
        Posterior Probability = -----
        Predictor prior probability
        * p(c)
        P(x|c)
        P(c|x) = -----
        P(x)
    """

    def __init__(self):
        """
            Attributes:
                likelihoods: Likelihood of each feature per class
                class_priors: Prior probabilities of classes
                pred_priors: Prior probabilities of features
                features: All features of dataset
        """
        self.features = list()
        self.likelihoods = {}
        self.class_priors = {}


```

```

self.pred_priors = {}

self.X_train = np.array
self.y_train = np.array
self.train_size = int
self.num_feats = int

def fit(self, X, y):

    self.features = list(X.columns)
    self.X_train = X
    self.y_train = y
    self.train_size = X.shape[0]
    self.num_feats = X.shape[1]

    for feature in self.features:
        self.likelihoods[feature] = {}
        self.pred_priors[feature] = {}

        for feat_val in np.unique(self.X_train[feature]):
            self.pred_priors[feature].update({feat_val: 0})

            for outcome in np.unique(self.y_train):
                self.likelihoods[feature].update({feat_val+' '+outcome:0})
                self.class_priors.update({outcome: 0})

        self._calc_class_prior()
        self._calc_likelihoods()
        self._calc_predictor_prior()

    def _calc_class_prior(self):

        """ P(c) - Prior Class Probability """

        for outcome in np.unique(self.y_train):
            outcome_count = sum(self.y_train == outcome)
            self.class_priors[outcome] = outcome_count / self.train_size

    def _calc_likelihoods(self):

        """ P(x|c) - Likelihood """

        for feature in self.features:

            for outcome in np.unique(self.y_train):
                outcome_count = sum(self.y_train == outcome)
                feat_likelihood = self.X_train[feature][self.y_train[self.y_train
                == outcome].index.values.tolist()].value_counts().to_dict()

```

```

        for feat_val, count in feat_likelihood.items():
            self.likelihoods[feature][feat_val + '_' + outcome] =
count/outcome_count

def _calc_predictor_prior(self):
    """ P(x) - Evidence """
    for feature in self.features:
        feat_vals = self.X_train[feature].value_counts().to_dict()

        for feat_val, count in feat_vals.items():
            self.pred_priors[feature][feat_val] = count/self.train_size

def predict(self, X):
    """ Calculates Posterior probability P(c|x) """
    results = []
    X = np.array(X)

    for query in X:
        probs_outcome = {}
        for outcome in np.unique(self.y_train):
            prior = self.class_priors[outcome]
            likelihood = 1
            evidence = 1

            for feat, feat_val in zip(self.features, query):
                likelihood *= self.likelihoods[feat][feat_val + '_' +
outcome]
                evidence *= self.pred_priors[feat][feat_val]

            posterior = (likelihood * prior) / (evidence)
            probs_outcome[outcome] = posterior

        result = max(probs_outcome, key = lambda x: probs_outcome[x])
        results.append(result)

    return np.array(results)

if __name__ == "__main__":
    #Weather Dataset
    print("\nWeather Dataset:")

```

```

df = pd.read_table("Weather_dataset.txt")
#print(df)

#Split features and target
X,y = pre_processing(df)

nb_clf=NaiveBayes()
nb_clf.fit(X, y)

print("Train Accuracy: {}".format(accuracy_score(y, nb_clf.predict(X)))))

#Query:
query = np.array([[['Rainy','Mild', 'Normal', 't']]])
print("Query 1:- {} ---> {}".format(query, nb_clf.predict(query)))

# Output:-
# Weather Dataset:
# Train Accuracy: 92.4
# Query 1:- [['Rainy' 'Mild' 'Normal' 't']] ---> ['Overcast Cool Normal t yes']

```

Part 2:

```

import pandas as pd
from sklearn.naive_bayes import CategoricalNB
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import accuracy_score

# Load Weather Dataset
df = pd.read_csv("weather_dataset.txt")

# Encode the target variable (class labels)
le = LabelEncoder()
y = le.fit_transform(df['Play'])

# Encode the categorical features
X = df.drop('Play', axis=1).apply(le.transform)

# Create and train the Naive Bayes classifier
nb_clf = CategoricalNB()
nb_clf.fit(X, y)

# Predict and calculate accuracy on the training data
y_pred = nb_clf.predict(X)
accuracy = accuracy_score(y, y_pred)

```

```

print("Train Accuracy: {:.2f}%".format(accuracy * 100))

# Query:
query = pd.DataFrame({'Outlook': ['Rainy'], 'Temperature': ['Mild'], 'Humidity': ['Normal'], 'Windy': ['False']})

# Use the same LabelEncoder for query data
query_encoded = query.apply(lambda col: le.transform(col))

prediction = nb_clf.predict(query_encoded)
print("Query 1:", le.inverse_transform(prediction))

```

Output:-
Weather Dataset:
Train Accuracy: 92.4
Query 1:- [['Rainy' 'Mild' 'Normal' 't']] --> ['Overcast Cool Normal t yes']

References:

[Naive Bayes Classification Program in Python from Scratch - japp.io](#)
[Naïve Bayes Algorithm -Implementation from scratch in Python. | by ranga_vamsi | Medium](#)

[ML | Naive Bayes Scratch Implementation using Python - GeeksforGeeks](#)
[How to Develop a Naive Bayes Classifier from Scratch in Python \(machinelearningmastery.com\)](#)

Conclusion:

I am able to understand the navie bayes theorem and am able to implement the code