

U.S.N.								
--------	--	--	--	--	--	--	--	--

B.M.S. College of Engineering, Bengaluru-560019

Autonomous Institute Affiliated to VTU

April 2025 Semester End Make-Up Examinations

Programme: B.E.

Semester: V

Branch: Computer Science and Engineering

Duration: 3 hrs.

Course Code: 23CS5PCDEV

Max Marks: 100

Course: Data Exploration and Visualization

Instructions: 1. Answer any FIVE full questions, choosing one full question from each unit.
2. Missing data, if any, may be suitably assumed.

UNIT - I			CO	PO	Marks
1	a)	Discuss the generic steps involved in Exploratory data analysis (EDA) process.	CO1	PO I	06
	b)	Imagine you are a data analyst working for an educational institute. You have been given a dataset containing information about students' performance in various subjects. The dataset includes the following columns: Student ID, Name, Gender, Age, Grade: The grade level of the student (e.g., 10th, 11th), Math Score, Science Score, English Score, Participation in Extracurricular Activities: Indicates if the student participates in extracurricular activities (e.g., Yes, No), Parent's Education Level. Classify the types of data in the given dataset and identify their measurement scales.	CO2	PO2	06
	c)	Demonstrate any four types of analysis that could be done for each of the following domains using EDA techniques. Justify your answer. i. Professional Sports ii. Healthcare	CO2	PO2	08
OR					
2	a)	Define Exploratory data analysis. Point out the primary aim and explain the significance of the same.	CO1	PO1	06
	b)	Demonstrate any three types of analysis that could be done for Marketing domain using EDA techniques. Justify your answer.	CO2	PO2	06
	c)	You have been provided with a dataset from a health and fitness app that tracks users' activities, dietary habits, and other lifestyle factors. The goal is to help the company understand different types of data, including their measurement scales, to improve the design of the app and offer personalized recommendations. The dataset includes the following variables: • User Demographics: Age, Gender, Height, Weight, City.	CO2	PO2	08

Important Note: Completing your answers, compulsorily draw diagonal cross lines on the remaining blank pages. Revealing of identification, appeal to evaluator will be treated as malpractice.

		<ul style="list-style-type: none"> Physical Activity Data: Daily Step Count, Hours of Sleep, Calories Burned, Exercise Frequency. Dietary Data: Daily Water Intake, Daily Caloric Intake, Food Categories. Health Metrics: Blood Pressure, Cholesterol Level, Blood Sugar Level. User Feedback: Overall Health Rating, Exercise Enjoyment Rating, Sleep Quality Rating <p>Classify the types of data in the given dataset and identify their measurement scales.</p>																																																										
		UNIT - II																																																										
3	a)	Discuss the process of data deduplication and the steps involved in the same with an example program.	<i>COI</i>	<i>POI</i>	06																																																							
	b)	<p>Write a Python program to perform the append and concatenation (with axis=1) operations on the following dataframes. Show the appropriate output.</p> <pre>data1 = { 'Name': ['Alice', 'Bob', 'Charlie'], 'Age': [25, 30, 35] } data2 = { 'Name': ['David', 'Eva'], 'Age': [40, 45] }.</pre>	<i>COI</i>	<i>POI</i>	06																																																							
	c)	<p>A retail company has been collecting data from its sales transactions, which includes the following information.</p> <ol style="list-style-type: none"> TransactionID: Unique identifier for each transaction. CustomerID: Identifier for the customer who made the purchase. ProductID: Identifier for the product purchased. Quantity: Quantity of the product purchased in the transaction. Price: Price per unit of the product. TotalAmount: Total amount for the transaction (Quantity * Price). <p>However, some of the records in the dataset are duplicates as shown below due to system errors, such as customers placing multiple orders for the same product, or the same transaction being recorded more than once.</p> <table border="1"> <thead> <tr> <th>TransactionID</th> <th>CustomerID</th> <th>ProductID</th> <th>Quantity</th> <th>Price</th> <th>TotalAmount</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>1</td> <td>101</td> <td>1001</td> <td>1</td> <td>20</td> <td>20</td> </tr> <tr> <td>1</td> <td>2</td> <td>102</td> <td>1002</td> <td>2</td> <td>15</td> <td>30</td> </tr> <tr> <td>2</td> <td>3</td> <td>101</td> <td>1001</td> <td>1</td> <td>20</td> <td>20</td> </tr> <tr> <td>3</td> <td>4</td> <td>104</td> <td>1003</td> <td>5</td> <td>30</td> <td>150</td> </tr> <tr> <td>4</td> <td>5</td> <td>103</td> <td>1002</td> <td>2</td> <td>15</td> <td>30</td> </tr> <tr> <td>5</td> <td>6</td> <td>102</td> <td>1002</td> <td>3</td> <td>15</td> <td>45</td> </tr> <tr> <td>6</td> <td>7</td> <td>101</td> <td>1001</td> <td>1</td> <td>20</td> <td>20</td> </tr> </tbody> </table> <p>Write a Python program to do the following and show subsequent output after each step.</p> <ol style="list-style-type: none"> Identify duplicate rows. Removing Duplicates Instead of simply removing the duplicates, replace duplicate value by aggregating data such as summing up the quantities and total amount and average the prices for the duplicate records. 	TransactionID	CustomerID	ProductID	Quantity	Price	TotalAmount	0	1	101	1001	1	20	20	1	2	102	1002	2	15	30	2	3	101	1001	1	20	20	3	4	104	1003	5	30	150	4	5	103	1002	2	15	30	5	6	102	1002	3	15	45	6	7	101	1001	1	20	20	<i>CO2</i>	<i>PO2</i>	08
TransactionID	CustomerID	ProductID	Quantity	Price	TotalAmount																																																							
0	1	101	1001	1	20	20																																																						
1	2	102	1002	2	15	30																																																						
2	3	101	1001	1	20	20																																																						
3	4	104	1003	5	30	150																																																						
4	5	103	1002	2	15	30																																																						
5	6	102	1002	3	15	45																																																						
6	7	101	1001	1	20	20																																																						

		OR																																																			
4	a)	<p>Discuss the importance of random sampling of a dataset. Differentiate between random sampling with replacement and without replacement.</p>	CO1	PO1	06																																																
	b)	<p>A retail company has collected data on the ages of its customers to analyze purchasing behavior across different age groups. To perform a more meaningful analysis, the company wants to classify the customers into distinct age groups (bins) rather than working with individual age values.</p> <p>The dataset contains the following columns:</p> <ul style="list-style-type: none"> CustomerID: Unique identifier for each customer. Age: Age of the customer. Spending Amount: Total amount spent by the customer in the last month. <p>Write a Python program to categorize customers into predefined age groups such as "18-25", "26-35", "36-45", etc. and find the average of their spending behavior within each group.</p>	CO2	PO2	06																																																
	c)	<p>A company collects data from two different sources: one containing information about its employees, and another containing the details of the projects that each employee is working on.</p> <p>DataFrame- df_employees</p> <table border="1"> <thead> <tr> <th></th><th>EmployeeID</th><th>Name</th><th>Department</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>Alice</td><td>HR</td></tr> <tr> <td>1</td><td>2</td><td>Bob</td><td>Finance</td></tr> <tr> <td>2</td><td>3</td><td>Charlie</td><td>IT</td></tr> <tr> <td>3</td><td>4</td><td>David</td><td>Sales</td></tr> <tr> <td>4</td><td>5</td><td>Eva</td><td>Marketing</td></tr> </tbody> </table> <p>DataFrame- df_projects</p> <table border="1"> <thead> <tr> <th></th><th>EmployeeID</th><th>ProjectID</th><th>ProjectName</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>101</td><td>Project A</td></tr> <tr> <td>1</td><td>2</td><td>102</td><td>Project B</td></tr> <tr> <td>2</td><td>4</td><td>103</td><td>Project C</td></tr> <tr> <td>3</td><td>5</td><td>104</td><td>Project D</td></tr> <tr> <td>4</td><td>6</td><td>105</td><td>Project E</td></tr> </tbody> </table> <p>These two data sources need to be merged to gain a complete picture of each employee and their associated projects.</p> <p>Analyse the types of joins that needs to be performed to do the following. Write necessary code along with suitable output.</p> <ol style="list-style-type: none"> If an employee doesn't have a corresponding project, they will be excluded from the result. A join that returns all rows from the df_employees and the matching rows from the right df_projects. If there is no match, the result will contain NaN for columns from the right DataFrame. 		EmployeeID	Name	Department	0	1	Alice	HR	1	2	Bob	Finance	2	3	Charlie	IT	3	4	David	Sales	4	5	Eva	Marketing		EmployeeID	ProjectID	ProjectName	0	1	101	Project A	1	2	102	Project B	2	4	103	Project C	3	5	104	Project D	4	6	105	Project E	CO2	PO2	08
	EmployeeID	Name	Department																																																		
0	1	Alice	HR																																																		
1	2	Bob	Finance																																																		
2	3	Charlie	IT																																																		
3	4	David	Sales																																																		
4	5	Eva	Marketing																																																		
	EmployeeID	ProjectID	ProjectName																																																		
0	1	101	Project A																																																		
1	2	102	Project B																																																		
2	4	103	Project C																																																		
3	5	104	Project D																																																		
4	6	105	Project E																																																		

		<p>3. A join that returns all rows from the df_projects and the matching rows from df_employees. If there is no match, the result will contain NaN for columns from the left DataFrame.</p> <p>4. A join that returns all rows from both DataFrames, with matching records from both sides when available. If there is no match, NaN is filled in for the columns from the DataFrame without a match.</p>			
		UNIT - III			
5	a)	Explain the most common measures for analyzing the distribution frequency. Demonstrate the same with a Python program.	CO1	PO1	06
	b)	<p>You are a data analyst working for a retail company that tracks sales data across different regions and products. You have a dataset containing information about sales figures and profit margins for various products in different regions.</p> <p>Data set: 'Region': ['North', 'North', 'South', 'South', 'North', 'South'], 'Product': ['A', 'B', 'A', 'B', 'A', 'B'], 'Sales': [100, 150, 200, 120, 180, 90], 'Profit_Margin': [0.2, 0.3, 0.25, 0.15, 0.18, 0.12]</p> <p>Write a program to calculate the total sales and maximum profit margin for each product across all regions.</p>	CO2	PO2	06
	c)	<p>Write a Python program to do the following.</p> <p>i.Create a dataframe for the following data.</p> <ol style="list-style-type: none"> 1. Date: The date of the transaction or record. 2. Category: The category of the product sold (e.g., Electronics, Clothing, Food). 3. Region: The region where the sale took place (e.g., North, South, East, West). 4. Sales: The total sales amount for each transaction. 5. Units Sold: The number of units sold in the transaction <p>ii.Create a pivot table that shows the total Sales by Category and Region.</p> <p>iii.Create a cross-tabulation that shows the Units Sold by Category and Region.</p> <p>iv.Print the average sales per unit sold for each product category across different regions.</p> <p>v.Print the region that has the highest sales for the Electronics category.</p>	CO2	PO2	08
		OR			
6	a)	Illustrate the various characteristics of the time series dataset with suitable example.	CO1	PO1	06
	b)	<p>As a data analyst for a retail company that wants to understand the relationship between the amount spent by customers and the time they spend on the company's website. This information will aid in optimizing the website experience and marketing strategies.</p> <p>1. What type of analysis technique can be applied?</p>	CO2	PO2	06

		2. Write a panda's code to explore the relationship between the time spent on the website (in minutes) and the amount spent by customers (in dollars) to derive insights for the retail company.																												
	c)	We have a dataset with information about employees in an organization. The dataset includes Employee ID , Department , Training Hours (hours spent in training during the last quarter), Job Performance Score (a composite score based on efficiency, quality of work, and job satisfaction), and Job Satisfaction (on a scale from 1 to 10). <table border="1"> <thead> <tr> <th>Employee ID</th> <th>Department</th> <th>Training Hours</th> <th>Job Performance Scores</th> <th>Job Satisfaction</th> </tr> </thead> <tbody> <tr> <td>101</td> <td>Sales</td> <td>40</td> <td>85</td> <td>8</td> </tr> <tr> <td>102</td> <td>Marketing</td> <td>50</td> <td>90</td> <td>9</td> </tr> <tr> <td>103</td> <td>HR</td> <td>20</td> <td>75</td> <td>7</td> </tr> <tr> <td>104</td> <td>IT</td> <td>60</td> <td>92</td> <td>9</td> </tr> </tbody> </table> a. Write a Python program to do the following. i. Create a data frame for the data shown. ii. Find the correlation between Training Hours and Job Performance Score for each department. iii. Find the overall correlation between Training Hours and Job Performance Score across all departments. b. Interpret the output for the given queries. i. Find the correlation between Training Hours and Job Performance Score for each department. ii. Find the overall correlation between Training Hours and Job Performance Score across all departments.	Employee ID	Department	Training Hours	Job Performance Scores	Job Satisfaction	101	Sales	40	85	8	102	Marketing	50	90	9	103	HR	20	75	7	104	IT	60	92	9	CO2	PO2	08
Employee ID	Department	Training Hours	Job Performance Scores	Job Satisfaction																										
101	Sales	40	85	8																										
102	Marketing	50	90	9																										
103	HR	20	75	7																										
104	IT	60	92	9																										
		UNIT - IV																												
7	a)	Illustrate the linear scale and logarithmic scale with respect to axes representation. Explain how the following data points could be represented using both the scales. 1,3,16,10,31,6,100	CO1	PO1	06																									

		b)	<p>Consider the following monthly sales dataset for 4 regions (North, South, East, West) across 6 months (January to June). The dataset contains the total sales in dollars for each region each month.</p> <table border="1"> <thead> <tr> <th>Month</th><th>North</th><th>South</th><th>East</th><th>West</th></tr> </thead> <tbody> <tr> <td>January</td><td>50,000</td><td>55,000</td><td>40,000</td><td>45,000</td></tr> <tr> <td>February</td><td>60,000</td><td>50,000</td><td>45,000</td><td>50,000</td></tr> <tr> <td>March</td><td>55,000</td><td>52,000</td><td>47,000</td><td>55,000</td></tr> <tr> <td>April</td><td>70,000</td><td>65,000</td><td>60,000</td><td>65,000</td></tr> <tr> <td>May</td><td>80,000</td><td>75,000</td><td>70,000</td><td>72,000</td></tr> <tr> <td>June</td><td>85,000</td><td>80,000</td><td>78,000</td><td>78,000</td></tr> </tbody> </table> <p>Identify the various plots and charts that are could be used to visualize the given data under following circumstances. Justify your answer.</p> <ol style="list-style-type: none"> To compare the total sales in each region across the months. To show individual data points for each month's sales across regions. To visualize the sales performance across regions and months in a more holistic way, where colour intensity represents sales volume. 	Month	North	South	East	West	January	50,000	55,000	40,000	45,000	February	60,000	50,000	45,000	50,000	March	55,000	52,000	47,000	55,000	April	70,000	65,000	60,000	65,000	May	80,000	75,000	70,000	72,000	June	85,000	80,000	78,000	78,000	CO2	PO2	06
Month	North	South	East	West																																					
January	50,000	55,000	40,000	45,000																																					
February	60,000	50,000	45,000	50,000																																					
March	55,000	52,000	47,000	55,000																																					
April	70,000	65,000	60,000	65,000																																					
May	80,000	75,000	70,000	72,000																																					
June	85,000	80,000	78,000	78,000																																					
		c)	<p>Explore the various plots and charts used to visualize associations among two or more quantitative variables. Analyze the visualization strategies that could be adopted to address the following relationships.</p> <ol style="list-style-type: none"> Relationship of two variables. Same types of variables with two different types of scales, position and size. More than three to four quantitative variables. 	CO2	PO2	08																																			
OR																																									
	8	a)	<p>List the limitations of Histogram and Density Plot. Explain how these limitations are handled using Empirical Cumulative Distribution Functions with an example.</p>	1	1	06																																			
		b)	<p>Consider a business scenario where a company wants to analyze the revenue composition across different product categories. Propose a specific visualization method that breaks down revenue proportions, highlighting the parts of the total. Justify your choice and explain how this visualization could guide decision-making within the company.</p>	CO2	PO2	06																																			
		c)	<p>Describe how confidence intervals are used to convey uncertainty in curve fits. What information do confidence intervals provide, and how can they be visually integrated into a curve fit visualization?</p>	1	1	08																																			
UNIT - V																																									
	9	a)	<p>Write a Python program to sort an array of words from longest to shortest length.</p>	1	1	06																																			
		b)	<p>Consider the following data set.</p>	CO2	PO2	06																																			

		<pre>data = { 'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'], 'Age': [24, 27, 22, 32, 29], 'City': ['New York', 'Los Angeles', 'Chicago', 'Houston'], 'Score': [85, 90, 78, 88, 92] }</pre> <p>Write python code by creating a dataframe of the above dataset to do the following.</p> <ol style="list-style-type: none"> Access the first name Select Name and City from first three rows. Print rows where age>25 			
	c)	Write a Python program that demonstrates how to create an Employee database with ID, Name, Age, Department as attributes using SQLite, insert data into it and then transfer the data into a Pandas dataframe.	CO2	PO2	08
OR					
10	a)	Write a Python program demonstrating how to reshape the data using melt() and pivot() methods. Illustrate the purpose of the same.	CO1	PO1	06
	b)	Write a Python program demonstrating the use of Web APIs and processing the retrieved data using pandas.	CO2	PO2	06
	c)	Write a Python program to create a Student dataframe with ID, Name, Age, Grade and Score. Find the missing values if any and fill the same using mean, mode and median for whichever data applicable. Justify your choice of filling the missing values.	CO2	PO2	08
