

|| Jai Sri Gurudev |

Sri AdichunchanagiriShikshana Trust (R)

SJB INSTITUTE OF TECHNOLOGY



Question Bank Module 2

Subject Name: Exploratory Data Analytics

Subject Code: 23CSE422

By

Faculty Name: Mrs. Shilpashree S

Designation: Assistant Professor

Semester: IV



Department of Computer Science & Engineering

Aca. Year: Even Sem /2024-25

Reference and Textbook Information

This document is designed as a supplementary study aid for the “**Exploratory Data Analytics - 23CSE422**” module, utilizing the prescribed textbook, “**Hands-On Exploratory Data Analysis with Python**” by **Suresh Kumar Mukhiya** and **Usman Ahmed**. The content presented here is based on the concepts and information provided in the textbook, with additional explanations, examples, and elaborations intended to enhance student understanding. This document is provided for educational purposes, in alignment with the intended use of the prescribed course textbook. This document is a study aid, based upon the prescribed text book, and is intended for the students of this course.

Bibliographic information

Title	Hands-On Exploratory Data Analysis with Python: Perform EDA Techniques to Understand, Summarize, and Investigate Your Data
Authors	Suresh Kumar Mukhiya , Usman Ahmed
Publisher	Packt Publishing, Limited, 2020
ISBN	1789537258, 9781789537253
Length	352 pages

Module 2 Question Bank with Scheme of Evaluation

Data Transformation: Merging database - style dataframes

Sl.No.	Question with Scheme of Evaluation	Marks
1.	List and explain the various data transformation techniques, highlighting their objectives.	10 Marks
SoE	Evaluation Scheme: <ul style="list-style-type: none"> • Explanation of Data Deduplication (1 Mark) • Explanation of Key Restructuring (1 Mark) • Explanation of Data Cleansing (1 Mark) • Explanation of Data Validation (1 Mark) • Explanation of Format Revisioning (1 Mark) • Explanation of Data Derivation (1 Mark) • Explanation of Data Aggregation (1 Mark) • Explanation of Data Integration (1 Mark) • Explanation of Data Filtering (1 Mark) • Explanation of Data Joining (1 Mark) 	
2.	How would you merge two pandas DataFrames using the <code>concat()</code> method? Demonstrate with an example where you combine data along rows and columns.	10 Marks
SoE	Evaluation Scheme: <ol style="list-style-type: none"> 1. Introduction to <code>concat()</code> Method (2 Marks) <ul style="list-style-type: none"> ○ Briefly explain the purpose of the <code>concat()</code> method (1 Mark). ○ Mention the role of <code>axis</code> argument (0.5 Marks). ○ Mention the <code>ignore_index</code> argument and its effect (0.5 Marks). 2. Example: Concatenating Along Rows (3 Marks) <ul style="list-style-type: none"> ○ Code example using <code>pd.concat([dataFrame1, dataFrame2], ignore_index=True)</code> to combine data along rows (2 Marks). ○ Output explanation: Describe the output and explain how the rows from both DataFrames are combined (1 Mark). 3. Example: Concatenating Along Columns (3 Marks) <ul style="list-style-type: none"> ○ Code example using <code>pd.concat([dataFrame1, dataFrame2], axis=1)</code> to combine data along columns 	

	<p>(2 Marks).</p> <ul style="list-style-type: none"> Output explanation: Describe the output and explain how the columns from both DataFrames are combined (1 Mark). <p>4. Explanation of Differences Between Row and Column Concatenation (2 Marks)</p> <ul style="list-style-type: none"> Explanation of what changes when <code>axis=0</code> (row-wise) and <code>axis=1</code> (column-wise) are used (2 Marks). 	
3.	<p>How can you merge DataFrames when the two DataFrames have different student IDs, where some students are in both courses, and some are only in one course? Explain the methods and demonstrate with examples.</p>	10 Marks
SoE	<p>Evaluation Scheme:</p> <ol style="list-style-type: none"> Introduction to Merging and Key Concepts (2 Marks) <ul style="list-style-type: none"> Mention the purpose of merging DataFrames based on common keys (1 Mark). Define <code>StudentID</code> as the key column for merging (1 Mark). Merge Methodology (Inner Join) (2 Marks) <ul style="list-style-type: none"> Explain the concept of an inner join and how it combines rows based on common values (1 Mark). Provide an example where <code>merge()</code> is used with <code>how='inner'</code> (1 Mark). Merge Methodology (Outer Join) (2 Marks) <ul style="list-style-type: none"> Explain the concept of an outer join and its usefulness when handling students who might only appear in one course (1 Mark). Provide an example where <code>merge()</code> is used with <code>how='outer'</code> (1 Mark). Handling Missing Data (2 Marks) <ul style="list-style-type: none"> Discuss what happens to students who are not in both courses and how missing data is handled in the merged DataFrame (2 Marks). Demonstration with Provided DataFrames (2 Marks) <ul style="list-style-type: none"> Provide code examples for merging the two DataFrames using both inner and outer joins. Ensure the resulting DataFrames show the correct combined data, including missing values where applicable (2 Marks). Conclusion and Observations (2 Marks) <ul style="list-style-type: none"> Summarize how merging works for both common and non-common student IDs (1 Mark). 	

	<ul style="list-style-type: none"> ○ Mention the importance of choosing the right type of join for the task (1 Mark). 	
4.	Describe the difference between concatenating DataFrames along axis=0 and axis=1. (10 Marks)	10 Marks
SoE	Introduction to concatenation (2 Marks) Explanation of axis=0 concatenation (3 Marks) Explanation of axis=1 concatenation (3 Marks) Real-life example demonstrating both (1 Mark) Conclusion summarizing the differences (1 Mark)	
5.	Explain how the merge function works in pandas, with an example for each type of join (inner, left, right, and outer).	10 Marks
SoE	Introduction to merging (2 Marks) Explanation of the different join types (6 Marks) <ul style="list-style-type: none"> ● Inner Join (1 Mark) ● Left Join (1 Mark) ● Right Join (1 Mark) ● Outer Join (1 Mark) Example code for each join type (1 Mark) Conclusion summarizing when to use each join (1 Mark)	
6.	Using pandas, merge two DataFrames on the index and explain the significance of using left_index=True and right_index=True. DataFrame 1 (df1)	10 Marks

	<table><tr><th>Product</th><th>Price</th><th>Quantity</th></tr><tr><td>A</td><td>10</td><td>100</td></tr><tr><td>B</td><td>15</td><td>150</td></tr><tr><td>C</td><td>20</td><td>200</td></tr><tr><td>D</td><td>25</td><td>250</td></tr></table> <p>DataFrame 2 (df2)</p> <table><tr><th>Discount</th><th>Sales</th></tr><tr><td>5</td><td>500</td></tr><tr><td>10</td><td>1000</td></tr><tr><td>15</td><td>1500</td></tr><tr><td>20</td><td>2000</td></tr></table>	Product	Price	Quantity	A	10	100	B	15	150	C	20	200	D	25	250	Discount	Sales	5	500	10	1000	15	1500	20	2000	
Product	Price	Quantity																									
A	10	100																									
B	15	150																									
C	20	200																									
D	25	250																									
Discount	Sales																										
5	500																										
10	1000																										
15	1500																										
20	2000																										
SoE	Introduction to merging on index (2 Marks)																										

	Code example of merging on index (4 Marks) Explanation of the left_index=True and right_index=True arguments (2 Marks) Conclusion on when merging on index is useful (2 Marks)	
7.	Describe the concept of hierarchical indexing in pandas and demonstrate how the stack and unstack methods are used to transform data with an example.	10 Marks
SoE	<ul style="list-style-type: none"> ● Introduction to Hierarchical Indexing (2 Marks): <ul style="list-style-type: none"> ○ Clear explanation of hierarchical indexing and its use in pandas. ● Explanation of the stack() Method (2 Marks): <ul style="list-style-type: none"> ○ Definition and purpose of the stack() method. ● Explanation of the unstack() Method (2 Marks): <ul style="list-style-type: none"> ○ Definition and purpose of the unstack() method. ● Example Demonstration (3 Marks): <ul style="list-style-type: none"> ○ Complete code example illustrating stacking and unstacking with clear comments. ● Missing Data Handling (1 Mark): <ul style="list-style-type: none"> ○ Explanation of NaN values appearing during unstacking. <p>Marks Breakdown Summary: Explanation of hierarchical indexing (2 Marks) Description of the stack method and its function (2 Marks) Description of the unstack method and its function (2 Marks) Example code demonstrating both stacking and unstacking (3 Marks) Discussion of missing data handling during unstacking (1 Mark)</p>	
	Alternate Question	
	Using pandas, demonstrate the process of stacking and unstacking a DataFrame. Discuss the potential issues when unstacking data and how to handle them. Explanation of stacking and unstacking (3 Marks) Step-by-step demonstration with an example (3 Marks) Explanation of missing data (NaN) during unstacking (2 Marks) Conclusion or real-life application (2 Marks)	
8.	Describe the process of data deduplication in a DataFrame. Include the use of the duplicated() and drop_duplicates() methods.	10 Marks

SoE	<ol style="list-style-type: none"> 1. Introduction to Data Deduplication (2 Marks): Brief explanation of the need for deduplication in datasets and its importance. 2. Explanation of duplicated() Method (3 Marks): <ul style="list-style-type: none"> Describe how it returns a Boolean series indicating which rows are duplicates. Example: Show how duplicated() can identify duplicates. 3. Explanation of drop_duplicates() Method (3 Marks): <ul style="list-style-type: none"> Describe the behavior of drop_duplicates() to remove duplicate rows. Mention how the first or last occurrence can be retained. 4. Detecting Duplicates in Specific Columns (1 Mark): <ul style="list-style-type: none"> Explain how drop_duplicates() can be used to identify duplicates in a subset of columns. 5. Practical Example (1 Mark): Provide a real-life example where deduplication would be crucial. 6. Conclusion (1 Mark): Summarize the importance of data deduplication in ensuring data quality. 	
	Alternate Questions	
	<ol style="list-style-type: none"> 1. Explain how to perform data deduplication based on specific columns in a DataFrame. Provide an example with a practical use case. Introduction to Column-Specific Deduplication (2 Marks): Explain the necessity of checking duplicates in specific columns for accurate data analysis. Using drop_duplicates() with Subset (4 Marks): <ol style="list-style-type: none"> Describe how to apply drop_duplicates() on a subset of columns. Example: Provide code to demonstrate column-specific deduplication. Handling Missing Data with Column-Specific Deduplication (2 Marks): <ol style="list-style-type: none"> Discuss how missing values are treated when identifying duplicates in selected columns. Practical Example (1 Mark): Give a real-world scenario where deduplication based on specific columns is required (e.g., removing duplicate email addresses). 	

	<p>Conclusion (1 Mark): Conclude the importance of deduplication based on columns in cleaning the dataset.</p> <hr/> <p>2. What are the consequences of not performing data deduplication in a dataset? Explain with an example. (10 Marks)</p> <p>Introduction to Data Deduplication (2 Marks): Define data deduplication and explain its importance.</p> <p>Consequences of Not Removing Duplicates (4 Marks):</p> <ul style="list-style-type: none"> ○ Discuss how duplicates can lead to skewed analysis and incorrect results. ○ Example: Demonstrate a case where duplicates cause miscalculations in analysis. <p>Impact on Data Quality and Integrity (2 Marks): Explain how duplicates reduce data quality and integrity.</p> <p>Real-World Example (1 Mark): Provide a scenario where ignoring deduplication caused problems (e.g., in customer databases).</p> <p>Conclusion (1 Mark): Summarize the importance of deduplication for accurate data analysis.</p>	
9.	<p>Explain the <code>replace()</code> method in pandas. How can it be used to replace single and multiple values in a DataFrame? Provide examples.</p>	6 Marks
SoE	<p>Introduction to <code>replace()</code> Method (1 Mark): Briefly define the <code>replace()</code> method in pandas.</p> <p>Single Value Replacement (2 Marks):</p> <ul style="list-style-type: none"> ● Explain how a single value can be replaced in a DataFrame. ● Provide a code example. <p>Multiple Value Replacements (2 Marks):</p> <ul style="list-style-type: none"> ● Explain how multiple values can be replaced at once. ● Provide a code example with different values being replaced. <p>Conclusion (1 Mark): Summarize the importance and usefulness of the <code>replace()</code> method in data preprocessing.</p>	
	<p>Alternate Questions</p>	

	<p>1. Demonstrate how you would replace specific values in a DataFrame with NaN and other values using pandas. (6 Marks)</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> • Introduction to Value Replacement (1 Mark): Define the need for replacing values in a DataFrame (e.g., missing data or invalid values). • Replacing a Single Value with NaN (2 Marks): <ul style="list-style-type: none"> ○ Show how to replace a single value with NaN. ○ Provide an example with -786 replaced by NaN. • Replacing Multiple Values (2 Marks): <ul style="list-style-type: none"> ○ Demonstrate replacing multiple values with NaN and other specific values. ○ Provide an example where -786 is replaced by NaN and 0 is replaced by 2. • Conclusion (1 Mark): Explain why replacing values (especially with NaN) is essential for data cleaning. <p>2. Why would you replace certain values in a DataFrame? Explain with an example how you would replace specific values with NaN and other values with different numbers. (6 Marks)</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> • Introduction to the Importance of Replacing Values (1 Mark): Explain the rationale behind replacing values in a DataFrame (e.g., handling missing or erroneous data). • Replacing Values with NaN (2 Marks): <ul style="list-style-type: none"> ○ Demonstrate replacing a specific value with NaN. ○ Provide an example with replacing -786 with NaN. • Replacing Multiple Values (2 Marks): <ul style="list-style-type: none"> ○ Show how to replace multiple values with NaN and other values with different numbers (e.g., replacing 0 with 2). • Conclusion (1 Mark): Summarize the utility of the <code>replace()</code> method in ensuring data consistency and cleanliness. 	
10.	Describe the reasons for missing data in a dataset. Provide examples.	10 Marks
SoE	Definition of NaN (2 Marks): Define NaN as an indicator of missing values.	

	<p>Reasons for Missing Data (4 Marks):</p> <ul style="list-style-type: none"> • Data retrieval issues (1 Mark). • Data joining issues (1 Mark). • Data collection errors (1 Mark). • Shape changes or reindexing issues (1 Mark). <p>Examples (3 Marks):</p> <ul style="list-style-type: none"> • Example from the DataFrame of stores and fruits (1.5 Marks). • Explanation of how each scenario leads to missing data (1.5 Marks). <p>Conclusion (1 Mark): Summarize the importance of handling missing data.</p>	
	Alternate Questions	
	<p>1. Explain how to detect missing data in a pandas DataFrame. Illustrate your answer with code examples.</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> • Detecting Missing Data (5 Marks): <ul style="list-style-type: none"> ○ Explanation of <code>isnull()</code> and <code>notnull()</code> methods (2 Marks). ○ Code examples of both methods (3 Marks). • Counting Missing Data (3 Marks): <ul style="list-style-type: none"> ○ Explanation of using <code>sum()</code> and <code>count()</code> methods (2 Marks). ○ Code example to count missing values (1 Mark). • Conclusion (2 Marks): Discuss why detecting missing data is crucial. <p>2. What are the methods available in pandas to handle missing data? Describe each with an example. (10 Marks)</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> • Methods for Handling Missing Data (7 Marks): 	

	<ul style="list-style-type: none"> ○ Explanation of <code>dropna()</code> for removing rows and columns (3 Marks). ○ Explanation of <code>dropna(how='all')</code> for dropping rows with all NaN values (2 Marks). ○ Explanation of <code>dropna(thresh=5, axis=1)</code> for removing columns with excessive NaN values (2 Marks). ● Code Examples (3 Marks): <ul style="list-style-type: none"> ○ Provide a code example for each of the methods described above. <p>3. How can missing values be counted and visualized in a pandas DataFrame? Illustrate with an example. (10 Marks)</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> ● Counting Missing Values (4 Marks): <ul style="list-style-type: none"> ○ Explanation of how to use <code>isnull().sum()</code> (2 Marks). ○ Code example showing the result (2 Marks). ● Visualizing Missing Values (3 Marks): <ul style="list-style-type: none"> ○ Mention visualization tools (if applicable), though not required for direct code. ● Practical Example (3 Marks): <ul style="list-style-type: none"> ○ Code example counting missing values in a DataFrame. <p>4. Demonstrate how to drop rows and columns with missing data using pandas. Provide code examples and discuss their functionality. (10 Marks)</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> ● Dropping Rows with Missing Data (4 Marks): <ul style="list-style-type: none"> ○ Explanation of <code>dropna()</code> and the behavior when applied to rows (2 Marks). ○ Code example (2 Marks). ● Dropping Columns with Missing Data (4 Marks): <ul style="list-style-type: none"> ○ Explanation of using <code>axis=1</code> for columns (2 Marks). ○ Code example (2 Marks). ● Conclusion (2 Marks): Summarize when it is appropriate to drop missing data and the impact on the dataset. 	
--	--	--

11.	Explain how NumPy handles missing values (NaN) during mathematical operations. Provide examples to support your answer.	10 Marks
SoE	<p>Definition of NaN in NumPy (2 Marks): Explanation of how NaN is treated in NumPy during operations.</p> <p>Mathematical Operations (4 Marks):</p> <ul style="list-style-type: none"> Describe how functions like <code>mean()</code>, <code>sum()</code>, etc., behave when NaN is present (2 Marks). Example to demonstrate NaN in a NumPy operation (2 Marks). <p>Conclusion (2 Marks): Summarize how NaN affects the overall output of calculations in NumPy.</p> <p>Example from Textbook (2 Marks): Use a relevant example from the provided content for clarity.</p>	
	Alternate Questions	
	<p>1. Compare the handling of NaN values in NumPy and Pandas when performing mathematical operations like mean, sum, and cumulative sum.</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> Difference in Handling NaN (4 Marks): <ul style="list-style-type: none"> Explain how NumPy returns NaN when NaN values are present, while Pandas ignores them (2 Marks). Example of <code>mean()</code> in both NumPy and Pandas (2 Marks). Sum Operation (2 Marks): Describe how Pandas treats NaN as 0 for summing. Cumulative Sum (2 Marks): Explain the difference in cumulative summing between NumPy and Pandas. Practical Example (2 Marks): Provide an example of a dataset and the output from both libraries. <p>2. In Pandas, how are NaN values handled during operations like summing and averaging? Illustrate with code examples.</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> Handling of NaN in Pandas (4 Marks): 	

	<ul style="list-style-type: none"> Describe how Pandas treats NaN values as 0 during sum and average computations (2 Marks). Code examples for sum and mean (2 Marks). Cumulative Sum (2 Marks): Provide a code example for cumulative summing in Pandas and explain how NaN is treated. Edge Case of All NaNs (2 Marks): Discuss how Pandas handles cases where all values are NaN in the operation. Conclusion (2 Marks): Summarize how Pandas ensures that operations proceed without errors due to NaN. 	
12.	Explain the purpose and working of the <code>fillna()</code> method in Pandas. Illustrate its impact on statistical metrics with an example.	10 Marks
SoE	<ul style="list-style-type: none"> Introduction to Missing Values (1 Mark): Define missing values and their impact on data analysis. Definition of <code>fillna()</code> (2 Marks): Explain the method and its syntax. Example with Explanation (4 Marks): <ul style="list-style-type: none"> Show a DataFrame with NaN values and its filled version. Provide code snippets and outputs. Discussion on Statistical Impact (2 Marks): Explain how replacing NaN values affects statistical metrics. Conclusion (1 Mark): Summarize key insights. <p>Marks Breakdown:</p> <ul style="list-style-type: none"> Introduction: 1 Mark Definition of <code>fillna()</code>: 2 Marks Example and Explanation: 4 Marks Impact on Statistics: 2 Marks Conclusion: 1 Mark 	
	Alternate Question	
	Demonstrate how replacing NaN values with 0 can influence the mean calculation of a DataFrame. Use code snippets and output examples to	

	<p>support your explanation.</p> <ul style="list-style-type: none"> ● Problem Context (1 Mark): Explain why NaN values are significant in analysis. ● Initial Mean Calculation (3 Marks): Show code and output for the original DataFrame mean. ● Modified Mean Calculation (3 Marks): Show code and output for the filled DataFrame mean. ● Comparison and Insights (2 Marks): Discuss the differences in results and explain why this occurs. ● Conclusion (1 Mark): Summarize the findings and their implications. <p>Marks Breakdown:</p> <ul style="list-style-type: none"> ● Problem Context: 1 Mark ● Initial Mean: 3 Marks ● Modified Mean: 3 Marks ● Comparison: 2 Marks ● Conclusion: 1 Mark 	
13.	<p>Explain and differentiate between forward and backward filling techniques for handling NaN values in a DataFrame. Illustrate with code examples and outputs.</p>	10 Marks
SoE	<ul style="list-style-type: none"> ● Definition of NaN Values and Filling Techniques (2 Marks): <ul style="list-style-type: none"> ○ Define missing values and explain why filling is necessary. ○ Introduce forward and backward filling. ● Explanation of Forward Filling (3 Marks): <ul style="list-style-type: none"> ○ Define the method. ○ Provide code and output. ● Explanation of Backward Filling (3 Marks): <ul style="list-style-type: none"> ○ Define the method. ○ Provide code and output. ● Comparison and Context (2 Marks): <ul style="list-style-type: none"> ○ Highlight differences between the techniques. ○ State when each method is preferable. <p>Marks Breakdown:</p> <ul style="list-style-type: none"> ● Definition: 2 Marks ● Forward Filling: 3 Marks ● Backward Filling: 3 Marks 	

	<ul style="list-style-type: none"> • Comparison: 2 Marks 	
	Alternate Question	
	Discuss how the <code>fillna()</code> method with forward and backward filling impacts data imputation. Provide suitable examples and outputs.	
	<ol style="list-style-type: none"> 1. Introduction to Data Imputation (1 Mark): Explain the need for filling NaN values. 2. Forward Filling Explanation (3 Marks): <ul style="list-style-type: none"> ○ Define the method. ○ Provide code and output examples. 3. Backward Filling Explanation (3 Marks): <ul style="list-style-type: none"> ○ Define the method. ○ Provide code and output examples. 4. Impact on Data (2 Marks): <ul style="list-style-type: none"> ○ Discuss how forward and backward filling methods influence data integrity and analysis results. ○ Include observations from examples. 5. Conclusion (1 Mark): Summarize the importance of choosing the right imputation method. <p>Marks Breakdown:</p> <ul style="list-style-type: none"> • Introduction: 1 Mark • Forward Filling: 3 Marks • Backward Filling: 3 Marks • Impact: 2 Marks • Conclusion: 1 Mark 	
14.	Explain how the <code>interpolate()</code> function in pandas is used to handle missing values in a dataset. Illustrate with an example and calculation steps.	10 Marks
SoE	<ul style="list-style-type: none"> • Introduction to Interpolation (2 Marks): <ul style="list-style-type: none"> ○ Define interpolation and its purpose in handling missing values. ○ Mention that pandas supports the <code>interpolate()</code> method. • Explanation of Linear Interpolation (3 Marks): 	

	<ul style="list-style-type: none"> ○ Discuss the default linear interpolation method. ○ Highlight how it calculates values using the known points before and after the NaN sequence. ● Code Example and Output (3 Marks): <ul style="list-style-type: none"> ○ Provide the example with the ser3 series. ○ Show the output after interpolation. ● Calculation Steps (2 Marks): <ul style="list-style-type: none"> ○ Clearly explain the steps to calculate missing values in the given example. <p>Marks Breakdown:</p> <ul style="list-style-type: none"> ● Introduction: 2 Marks ● Linear Interpolation: 3 Marks ● Code Example: 3 Marks ● Calculation: 2 Marks 	
	Alternate Questions	
	<p>Describe the process of linear interpolation in pandas using the interpolate() method. How does it calculate missing values in a series? Provide an example to explain your answer.</p> <p>Scheme of Evaluation</p> <ol style="list-style-type: none"> 1. Definition and Purpose (1 Mark): <ul style="list-style-type: none"> ○ Define linear interpolation and its role in filling missing values. 2. Interpolation Process (4 Marks): <ul style="list-style-type: none"> ○ Explain how the method uses values before and after NaN sequences. ○ Discuss dividing the difference evenly among missing entries. 3. Detailed Example (3 Marks): <ul style="list-style-type: none"> ○ Present the ser3 example with output. ○ Explain the calculations for each interpolated value. 4. Conclusion (2 Marks): <ul style="list-style-type: none"> ○ Summarize the effectiveness of interpolation for data cleaning. ○ Mention its utility in time series data for complex interpolations. <p>Marks Breakdown:</p> <ul style="list-style-type: none"> ● Definition: 1 Mark ● Interpolation Process: 4 Marks 	

	<ul style="list-style-type: none"> • Example: 3 Marks • Conclusion: 2 Marks 	
15.	Explain the use of the <code>index.map()</code> and <code>rename()</code> methods in pandas for renaming axis indexes. Provide examples to illustrate your answer.	8 Marks
SoE	<ul style="list-style-type: none"> • Introduction to Axis Index Renaming (1 Marks): <ul style="list-style-type: none"> ○ Explain why renaming indexes or columns is necessary in data transformation. • Explanation of <code>index.map()</code> (3 Marks): <ul style="list-style-type: none"> ○ Describe how <code>index.map()</code> works. ○ Mention that it directly modifies the original DataFrame. ○ Provide an example (e.g., converting index to uppercase). • Explanation of <code>rename()</code> (3 Marks): <ul style="list-style-type: none"> ○ Describe how the <code>rename()</code> method transforms indexes or columns. ○ Highlight its ability to preserve the original DataFrame. ○ Provide an example (e.g., converting columns to uppercase and indexes to title case). • Comparison and Use Cases (1 Marks): <ul style="list-style-type: none"> ○ Compare the two methods, noting when to use each. ○ Highlight the non-modifying nature of <code>rename()</code>. <p>Marks Breakdown:</p> <ul style="list-style-type: none"> • Introduction: 1 Marks • <code>index.map()</code> Explanation: 3 Marks • <code>rename()</code> Explanation: 3 Marks • Comparison: 1 Marks 	
16.	Define discretization and binning. How do these techniques help in data analysis?	10 Marks
SoE	<ul style="list-style-type: none"> • Definition of Discretization (3 Marks) 	

	<ul style="list-style-type: none"> ○ Define discretization as converting continuous data into discrete intervals. Mention its importance in simplifying analysis. (3 Marks) ● Definition of Binning (3 Marks) <ul style="list-style-type: none"> ○ Define binning as the method used to discretize data into intervals or bins. Discuss different types of binning (e.g., equal-width, equal-frequency). (3 Marks) ● How Discretization and Binning Aid Data Analysis (2 Marks) <ul style="list-style-type: none"> ○ Explain how these techniques help reduce complexity and make data easier to analyze and visualize. (2 Marks) ● Example (2 Marks) <ul style="list-style-type: none"> ○ Provide a practical example of binning (e.g., discretizing heights or age data into bins). (2 Marks) <p>Marks Breakdown</p> <ul style="list-style-type: none"> ● Definition of Discretization (3 Marks) ● Definition of Binning (3 Marks) ● How Discretization and Binning aid Data Analysis (2 Marks) ● Example (2 Marks) 	
	Alternate Questions	
	<ol style="list-style-type: none"> 1. Explain how to use the cut () method in Pandas to perform binning. Provide an example. (10 Marks) Evaluation Scheme: <ul style="list-style-type: none"> ● Overview of cut () method (3 Marks) ● Explanation of syntax and parameters (4 Marks) ● Example of cut () method with explanation (3 Marks) 2. Discuss the difference between open and closed intervals with examples. How do these affect the output when using the cut () method in Pandas? (10 Marks) Evaluation Scheme: <ul style="list-style-type: none"> ● Definition and Explanation of Open and Closed Intervals (4 Marks) ● Examples of Open and Closed Intervals (3 Marks) 	

	<ul style="list-style-type: none"> • How Open/Closed Intervals Affect the cut () Method (3 Marks) <p>3. How does the qcut () method in Pandas differ from the cut () method? Provide examples to support your answer. (10 Marks)</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> • Definition and Explanation of qcut () (3 Marks) • Comparison with cut () (3 Marks) • Example of qcut () and Explanation (4 Marks) <p>4. Explain how to customize bin labels in the cut () method. Provide an example to illustrate your explanation. (10 Marks)</p> <p>Evaluation Scheme:</p> <ul style="list-style-type: none"> ○ Explanation of Customizing Labels in cut () (3 Marks) ○ Example of Custom Labels (4 Marks) ○ Explanation of the Output (3 Marks) 	
17.	<p>Define outliers and explain their impact on data analysis. Describe the steps to detect and filter outliers in a dataset. Illustrate your answer with an example using sales data, where TotalPrice is used as a metric for outlier detection.</p>	10 Marks
SoE	<p>Definition of Outliers (2 Marks):</p> <ul style="list-style-type: none"> • Provide a clear definition. • Discuss the reasons for their occurrence. <p>Impact of Outliers (2 Marks):</p>	

	<ul style="list-style-type: none">● Explain how outliers affect analysis and modeling. <p>Steps for Outlier Detection (4 Marks):</p> <ul style="list-style-type: none">● Load dataset (1 Mark).● Calculate derived metrics like TotalPrice (1 Mark).● Detect and filter outliers using thresholds (2 Marks). <p>Example (2 Marks):</p> <ul style="list-style-type: none">● Provide a real-life example using the sales data																																																			
	Alternate Questions																																																			
	<p>Question 2</p> <p>You are provided with a dataset of sales transactions. Perform the following tasks:</p> <ul style="list-style-type: none">a. Calculate the TotalPrice by multiplying UnitPrice and Quantity.b. Detect outliers where TotalPrice exceeds 3,000,000.c. Filter and display rows where TotalPrice is greater than 6,741,112. <p>Dataset :</p> <table><tr><th>TransactionID</th><th>Quantity</th><th>UnitPrice</th><th>CustomerID</th><th>Date</th></tr><tr><td>1</td><td>10</td><td>150000</td><td>101</td><td>2025-01-01</td></tr><tr><td>2</td><td>5</td><td>750000</td><td>102</td><td>2025-01-02</td></tr><tr><td>3</td><td>2</td><td>400000</td><td>103</td><td>2025-01-03</td></tr><tr><td>4</td><td>8</td><td>550000</td><td>104</td><td>2025-01-04</td></tr><tr><td>5</td><td>12</td><td>70000</td><td>105</td><td>2025-01-05</td></tr><tr><td>6</td><td>15</td><td>100000</td><td>106</td><td>2025-01-06</td></tr><tr><td>7</td><td>20</td><td>65000</td><td>107</td><td>2025-01-07</td></tr><tr><td>8</td><td>3</td><td>1200000</td><td>108</td><td>2025-01-08</td></tr><tr><td>9</td><td>9</td><td>450000</td><td>109</td><td>2025-01-09</td></tr></table>	TransactionID	Quantity	UnitPrice	CustomerID	Date	1	10	150000	101	2025-01-01	2	5	750000	102	2025-01-02	3	2	400000	103	2025-01-03	4	8	550000	104	2025-01-04	5	12	70000	105	2025-01-05	6	15	100000	106	2025-01-06	7	20	65000	107	2025-01-07	8	3	1200000	108	2025-01-08	9	9	450000	109	2025-01-09	
TransactionID	Quantity	UnitPrice	CustomerID	Date																																																
1	10	150000	101	2025-01-01																																																
2	5	750000	102	2025-01-02																																																
3	2	400000	103	2025-01-03																																																
4	8	550000	104	2025-01-04																																																
5	12	70000	105	2025-01-05																																																
6	15	100000	106	2025-01-06																																																
7	20	65000	107	2025-01-07																																																
8	3	1200000	108	2025-01-08																																																
9	9	450000	109	2025-01-09																																																

	10 7 950000 110 2025-01-10	
	Scheme of Evaluation: <ol style="list-style-type: none"> 1. Code to Calculate TotalPrice (3 Marks) 2. Outlier Detection (3 Marks) 3. Filtering Outliers (2 Marks) 4. Explanations (2 Marks) 	
18.	Explain the concept of permutation and random sampling. Demonstrate how to randomly reorder the rows of a DataFrame using Pandas and NumPy.	8 Marks
SoE	<ul style="list-style-type: none"> • Definition of Permutation and Random Sampling (1 Marks): <ul style="list-style-type: none"> ○ Define both concepts with relevance to data analysis. • Explanation of <code>np.random.permutation()</code> (2 Marks): <ul style="list-style-type: none"> ○ Describe how this function generates a random order of indices. • Explanation of <code>take()</code> Function in Pandas (2 Marks): <ul style="list-style-type: none"> ○ Discuss how <code>take()</code> applies the random order to reorder the DataFrame. • Python Code Demonstration (2 Marks): <ul style="list-style-type: none"> ○ Include properly indented and correct code. • Conclusion/Output Explanation (1 Mark): <ul style="list-style-type: none"> ○ Explain how the final DataFrame is permuted based on the sampler. 	
	Alternate Questions	
	<p>Using the <code>numpy.random.permutation()</code> function, write a Python program to permute the rows of a 10x8 DataFrame. Explain the process step-by-step. (8 Marks)</p> <p>Scheme of Evaluation:</p> <p>Creation of DataFrame (1 Marks):</p> <ul style="list-style-type: none"> • Code to create a 10x8 DataFrame using NumPy and Pandas. <p>Using <code>np.random.permutation()</code> (3 Marks):</p>	

	<ul style="list-style-type: none"> • Generate and explain the sampler array. <p>Applying Permutation with <code>take()</code> (3 Marks):</p> <ul style="list-style-type: none"> • Demonstrate and explain the reordering of rows. <p>Output Explanation (1 Marks):</p> <ul style="list-style-type: none"> • Compare the permuted DataFrame with the original. 	
19.	<p>Explain the process of random sampling without replacement. Demonstrate how to randomly sample 3 rows from a DataFrame using the <code>np.random.permutation()</code> and <code>df.take()</code> methods.</p>	10 Marks
SoE	<p>Definition and Explanation of Random Sampling Without Replacement (2 Marks):</p> <ul style="list-style-type: none"> • Clear definition of random sampling without replacement and its application in data analysis. <p>Steps for Random Sampling Without Replacement (3 Marks):</p> <ul style="list-style-type: none"> • Correctly describe the steps: creating a permutation array, slicing, and applying <code>df.take()</code>. <p>Code Demonstration (3 Marks):</p> <ul style="list-style-type: none"> • Correct code that demonstrates random sampling without replacement. <p>Output Explanation (2 Marks):</p> <ul style="list-style-type: none"> • Provide a clear explanation of the output and how the sampled rows were selected. 	
20.	<p>Describe how random sampling with replacement can be performed using <code>numpy.random.randint()</code>. Write a Python program that generates a random sample of 10 elements from a given array, allowing repeated values.</p>	10 Marks

SoE	<p>Definition of Random Sampling With Replacement (2 Marks):</p> <ul style="list-style-type: none"> Define random sampling with replacement and explain its use. <p>Explanation of <code>numpy.random.randint()</code> (3 Marks):</p> <ul style="list-style-type: none"> Describe how <code>numpy.random.randint()</code> works and how it is used for random sampling. <p>Code Demonstration (3 Marks):</p> <ul style="list-style-type: none"> Correct code to generate random samples with replacement from an array. <p>Output Explanation (2 Marks):</p> <ul style="list-style-type: none"> Provide an explanation of the output and how some elements were sampled multiple times. 	
21.	<p>Describe the process of creating dummy variables from a categorical column in a DataFrame. Demonstrate with an example by converting a gender column into dummy variables using <code>pd.get_dummies()</code>.</p>	10 Marks
SoE	<p>Definition and Explanation of Dummy Variables (2 Marks):</p> <ul style="list-style-type: none"> Briefly define dummy variables and explain their purpose in converting categorical data to numeric. <p>Step-by-Step Process (3 Marks):</p> <ul style="list-style-type: none"> Explain the process of creating dummy variables using <code>pd.get_dummies()</code>, including the conversion of categorical values into binary columns. <p>Code Demonstration (3 Marks):</p> <ul style="list-style-type: none"> Correct code to generate dummy variables from a categorical column (gender), as shown in the provided example. <p>Output Explanation (2 Marks):</p> <ul style="list-style-type: none"> Provide a clear explanation of the output and how the DataFrame 	

	is transformed into dummy variables.	
22.	Explain how to add a prefix to the column names when creating dummy variables using <code>pd.get_dummies()</code>. Provide a Python code snippet demonstrating this transformation with the gender column in a DataFrame.	10 Marks
SoE	<ul style="list-style-type: none"> • Explanation of Prefix Usage (2 Marks): <ul style="list-style-type: none"> ○ Describe how the prefix argument in <code>pd.get_dummies()</code> is used to modify column names. • Code Demonstration (4 Marks): <ul style="list-style-type: none"> ○ Correct code that demonstrates the addition of a prefix to the column names when creating dummy variables. • Output Explanation (2 Marks): <ul style="list-style-type: none"> ○ Provide a clear explanation of the transformed output, emphasizing the prefixed column names. 	
23.	<p>Write a Python program to demonstrate the following string operations:</p> <ol style="list-style-type: none"> Create strings using single, double, and triple quotes. Display the first and last characters of a string. Extract a substring from the 3rd to the 12th character and between the 3rd and second-last character. Attempt to update a character in a string and explain why it's not possible. Update the entire string with a new value and display the updated string. Use escape sequences for single quotes, double quotes, backslashes, and newlines in a string. Format a string using default, positional, and keyword formatting. Given a CSV file <code>comments.csv</code>, use string methods to: <ol style="list-style-type: none"> Convert text to lowercase/uppercase. Strip unwanted characters. Replace specific words (e.g., "Wolves" → "Fox"). Extract URLs from Reddit comment posts in the dataset using regular expressions. 	10 Marks
SoE	(9 x 1 = 9) + (1 Mark for Output) = 10 Marks	

