

Test: FT4

Course Code & Title: 21CSS303T-Data Science

Year& Sem: III Year /VI Sem

Date: 29-04-2025

Duration: Two periods

Max.Marks:50

**SET –**  
**Answer Key**

**Part – A (10 x 1 = 10 Marks)**

S.No	Question	Marks
1	a) Data Collection → Data Cleaning → Data Transformation → Data Analysis	1
2	b) Replaces all NaN values with 0	1
3	d) merge()	1
4	a) To combine datasets horizontally or vertically	1
5	b) Remove non-numeric values or replace them with NaN	1
6	b) To add text annotations to specific points on the plot	1
7	c) Pair plot	1
8	a) plt.style.use('seaborn-darkgrid')	1
9	b) sns.histplot()	1
10	d) Plots a scatter plot matrix grouped by species	1

Q. No	Part – B (4 x 5 = 20 Marks) Instructions: Answer <b>ANY FOUR</b>	Marks
11	<b>Discuss the general programming tips to deal with large data sets.</b> <ul style="list-style-type: none"> <li>• Don't reinvent the wheel. Use tools and libraries developed by others</li> <li>• Get the most out of your hardware. Your machine is never used to its full potential; with simple adaptations you can make it work harder.</li> <li>• Reduce the computing need. Slim down your memory and processing needs as much as possible.</li> </ul>	5
12	<b>When merging two DataFrames in pandas that have columns with the same name, how can you ensure the column names are distinguishable?</b> Use the suffixes parameter in the merge() function to add distinguishing suffixes to overlapping column names. <pre>import pandas as pd df1 = pd.DataFrame({'ID': [1, 2], 'Value': [10, 20]}) df2 = pd.DataFrame({'ID': [1, 2], 'Value': [30, 40]}) merged_df = pd.merge(df1, df2, on='ID', suffixes=('_left', '_right')) print(merged_df)</pre>	5

13	<p><b>Given the dataset data = {'Ages': [3, 18, 22, 10, 25, 29, 34, 14, 40, 45, 50, 55, 60, 12, 65, 70, 75, 80, 85]}, categorize the continuous Ages values into the groups of children, young, middle, and elder. Define appropriate age ranges for each category and implement the conversion.</b></p> <pre>import pandas as pd data = {'Ages': [3, 18, 22, 10, 25, 29, 34, 14, 40, 45, 50, 55, 60, 12, 65, 70, 75, 80, 85]} df = pd.DataFrame(data) bins = [0, 12, 24, 59, 100] labels = ['Child', 'Young', 'Middle', 'Elder'] df['Category'] = pd.cut(df['Ages'], bins=bins, labels=labels) print(df)</pre>	5
14	<p><b>Compare a box plot and a histogram, highlighting their use cases and strengths.</b></p> <p><b>Box Plot:</b></p> <ul style="list-style-type: none"> <li>Displays the distribution of data and highlights outliers.</li> <li>Ideal for comparing multiple datasets.</li> </ul> <p><b>Histogram:</b></p> <ul style="list-style-type: none"> <li>Shows the frequency distribution of data values.</li> <li>Useful for understanding the shape of the data (e.g., skewness).</li> </ul>	5
15	<p><b>How can you control the line properties (e.g., color, style, and width) of a chart in Matplotlib. Write the python code and explain.</b></p> <pre>import matplotlib.pyplot as plt x = [1, 2, 3, 4, 5] y = [10, 20, 15, 25, 30] plt.plot(x, y, color='red', linestyle='--', linewidth=2) plt.title("Line Properties Example") plt.xlabel("X-axis") plt.ylabel("Y-axis") plt.show()</pre> <ul style="list-style-type: none"> <li>color: Sets the line color.</li> <li>linestyle: Controls the style of the line (e.g., dashed, solid).</li> <li>linewidth: Adjusts the thickness of the line.</li> </ul>	5

Q. No	<b>Part – C (2 x 10 = 20 Marks)</b> Instructions: Answer ALL questions.	Marks
16 a	<ul style="list-style-type: none"> <li><b>Missing Data:</b> Fill missing sales values with the median (robust to outliers). Drop rows if there are very few missing values. Example Code: df['Sales'] = df['Sales'].fillna(df['Sales'].median())</li> <li><b>Irregular Formats:</b> Convert all dates into a uniform format (YYYY-MM-DD) using pd.to_datetime. Example Code: df['Date'] = pd.to_datetime(df['Date'], errors='coerce')</li> <li><b>Duplicate Records:</b> Remove rows where Product, Region, and Date are duplicated, keeping the first occurrence: Example Code: df = df.drop_duplicates(subset=['Product', 'Region', 'Date'], keep='first')</li> <li><b>Irrelevant Data:</b> Drop unnecessary or irrelevant columns like Transaction ID Example Code:</li> </ul>	10

	<pre>df = df.drop(columns=['Transaction ID'])</pre> <ul style="list-style-type: none"> <li> <b>Outliers:</b>  Identify outliers in <code>Sales</code> using the interquartile range (IQR)  Example Code: <pre>Q1 = df['Sales'].quantile(0.25) Q3 = df['Sales'].quantile(0.75) IQR = Q3 - Q1 lower_bound = Q1 - 1.5 * IQR upper_bound = Q3 + 1.5 * IQR df = df[(df['Sales'] &gt;= lower_bound) &amp; (df['Sales'] &lt;= upper_bound)]</pre> </li> <li> <b>Categorical Inconsistencies:</b>  Standardize inconsistent product names using a mapping dictionary:  Example Code: <pre>product_mapping = {'Appl': 'Apple', 'Bananaa': 'Banana'} df['Product'] = df['Product'].replace(product_mapping)</pre> </li> <li> <b>Merging:</b>  Load the profit margins dataset and merge with the sales data on Product and Region  Example Code: <pre>profit_data = pd.read_csv('profit_margins.csv') df = pd.merge(df, profit_data, on=['Product', 'Region'], how='inner')</pre> </li> <li> <b>Final Quality Checks</b> <ul style="list-style-type: none"> <li> <b>Ensure all columns have consistent data types:</b>   <pre>df['Sales'] = df['Sales'].astype(float) df['Date'] = pd.to_datetime(df['Date'])</pre> </li> <li> <b>Verify no missing or inconsistent values remain:</b>   <pre>print(df.isnull().sum())</pre> </li> </ul> </li> </ul>	
16 b	<p><b>Output of <code>pivot_df = df.pivot(index='Date', columns='Product', values='Sales')</code></b>  The pivot function reshapes the DataFrame by specifying:</p> <ul style="list-style-type: none"> <li>index: Rows of the resulting DataFrame (Date here).</li> <li>columns: Columns of the resulting DataFrame (Product here).</li> <li>values: Data to fill the cells (Sales here).</li> </ul> <p><b>Output:</b></p> <pre>       Product  Apple  Banana Date 2023-01-01   100    150 2023-01-02   200     50 </pre> <p><b>Discussion:</b></p> <ul style="list-style-type: none"> <li>The rows are indexed by Date.</li> <li>The columns are determined by unique values in Product.</li> <li>The values in the cells are taken from the Sales column.</li> </ul> <p><b>2. Output of <code>stacked_df = df.stack()</code></b>  The stack function compresses columns into a hierarchical index at the row level.</p> <p><b>Output:</b></p> <pre> 0 Date    2023-01-01 </pre>	10

	<pre> Product      Apple Region       North Sales        100 1 Date       2023-01-01 Product      Banana Region       North Sales        150 2 Date       2023-01-02 Product      Apple Region       South Sales        200 3 Date       2023-01-02 Product      Banana Region       South Sales        50 dtype: object </pre> <p><b>Discussion:</b></p> <ul style="list-style-type: none"> <li>Each row is identified by a combination of the original row index (e.g., 0, 1, etc.) and the column name (e.g., Date, Product, Region, Sales).</li> <li>The DataFrame is reshaped into a Series with a multi-level index.</li> </ul> <p><b>3. Output of <code>stacked_pivot = pivot_df.stack()</code></b>  The stack function on pivot_df moves the columns (Product) back into the row index.</p> <p><b>Output:</b></p> <pre> Date      Product 2023-01-01 Apple    100           Banana    150 2023-01-02 Apple    200           Banana     50 dtype: int64 </pre> <p><b>Discussion:</b></p> <ul style="list-style-type: none"> <li>The columns (Apple, Banana) are turned into a new hierarchical index level under Date.</li> <li>The resulting structure is a Series, with the multi-level index representing the combination of Date and Product.</li> </ul>	
17 a	<p><b>Functionalities of the Seaborn Library</b></p> <p>Seaborn is a Python library built on top of Matplotlib that simplifies creating informative and aesthetically pleasing statistical graphics. It provides high-level interfaces for creating attractive and complex visualizations.</p> <p><b>Key Features:</b></p> <ol style="list-style-type: none"> <li>1. Theme Customization: Automatically styles plots for aesthetics.</li> <li>2. Dataset-Oriented: Works efficiently with DataFrames and arrays.</li> <li>3. Built-In Statistical Analysis: Includes options for regression, distribution fitting, and more.</li> <li>4. Integration with Pandas: Seamless handling of DataFrame columns.</li> <li>5. Wide Range of Plot Types: Includes pair plots, box plots, violin plots, heatmaps, and more.</li> </ol> <p><b>Examples of Key Visualizations</b></p> <p><b>1. Pair Plot</b></p> <p>A pair plot is useful for visualizing pairwise relationships in a dataset, especially numerical features. It provides scatterplots for relationships and histograms for univariate distributions.</p>	10

	<p><b>Code Example:</b></p> <pre>import seaborn as sns import pandas as pd  # Example Dataset data = sns.load_dataset('iris')  # Pair Plot sns.pairplot(data, hue='species')</pre> <p><b>Functionality:</b></p> <ul style="list-style-type: none"> <li>• Displays scatter plots between every pair of numerical columns.</li> <li>• Includes diagonal histograms to visualize the distribution of each feature.</li> <li>• Uses hue to color the data points based on a categorical column (species).</li> </ul> <p><b>2. Box Plot</b></p> <p>A box plot summarizes the distribution of a dataset through five-number summary statistics: minimum, first quartile (Q1), median, third quartile (Q3), and maximum. It also highlights potential outliers.</p> <p><b>Code Example:</b></p> <pre># Box Plot sns.boxplot(x='species', y='sepal_width', data=data)</pre> <p><b>Functionality:</b></p> <ul style="list-style-type: none"> <li>• Displays distributions and compares groups (e.g., species) for a numerical column (sepal_width).</li> <li>• Identifies outliers as points outside the whiskers.</li> <li>• Can be enhanced with swarm plots to overlay individual data points.</li> </ul> <p><b>3. Histogram</b></p> <p>A histogram visualizes the distribution of a single numerical variable by grouping data into bins.</p> <p><b>Code Example:</b></p> <pre># Histogram sns.histplot(data['sepal_length'], kde=True, bins=20)</pre> <p><b>Functionality:</b></p> <ul style="list-style-type: none"> <li>• Shows the frequency of data points within specified bins.</li> <li>• Optionally overlays a kernel density estimate (KDE) curve for a smoothed representation of the distribution.</li> <li>• Parameters like bins control the granularity of the visualization.</li> </ul>	
17 b	<p><b>Creating and Interpreting 3D Surface Plots</b></p> <p><b>Example Program ( 5 marks for the correct usage of add_subplot, polt.surface and set_xlabel methods)</b></p> <p>3D surface plots are a type of visualization used to represent three-dimensional data where the z-axis corresponds to the dependent variable, and the x and y axes represent independent variables. These plots are useful for exploring relationships between variables and identifying patterns or trends.</p> <pre>import numpy as np import matplotlib.pyplot as plt from mpl_toolkits.mplot3d import Axes3D  # Generate data</pre>	10

```

x = np.linspace(-5, 5, 100)
y = np.linspace(-5, 5, 100)
X, Y = np.meshgrid(x, y)
Z = np.sin(np.sqrt(X**2 + Y**2))

# Create the 3D plot
fig = plt.figure(figsize=(8, 6))
ax = fig.add_subplot(111, projection='3d')
surface = ax.plot_surface(X, Y, Z, cmap='viridis', edgecolor='none')

# Add labels and a color bar
ax.set_xlabel('X-axis')
ax.set_ylabel('Y-axis')
ax.set_zlabel('Z-axis')
plt.colorbar(surface, ax=ax, shrink=0.5, aspect=10)
plt.title('3D Surface Plot Example')
plt.show()

```

### Scenarios Where 3D Surface Plots Are Beneficial ( 5 Mark)

#### 1. Data Exploration in Engineering:

- Analyzing stress, temperature, or pressure distribution over a 2D plane.
- Example: Surface temperature of a material under specific conditions.

#### 2. Optimization Problems:

- Visualizing cost functions in machine learning or operations research.
- Example: Understanding the shape of loss functions during model training.

#### 3. Geographic and Environmental Data:

- Representing terrain elevation or pollution levels.
- Example: A surface plot of altitude over a geographical region.

#### 4. Physics and Mathematics:

- Illustrating functions or mathematical surfaces.
- Example: Visualizing wave functions or potential fields.

#### 5. Economics and Finance:

- Exploring relationships between variables like interest rates, risk, and returns.
- Example: 3D visualization of portfolio optimization.