

PRACTICAL 4

ALL PRACTICALS:

4.2.8. Titanic Dataset Analysis and Data Cleaning - 4

You are provided with the Titanic dataset containing information about passengers on the Titanic. Your task is to write Python code to answer the following questions based on the dataset.

1. Get the number of survivors by gender (Sex).
2. Get the number of non-survivors by gender (Sex).
3. Get the number of survivors by embarkation location (Embarked_S).
4. Get the number of non-survivors by embarkation location (Embarked_S).
5. Calculate the percentage of children (Age < 18) who survived.
6. Calculate the percentage of adults (Age >= 18) who survived.
7. Get the median age of survivors.
8. Get the median age of non-survivors.
9. Get the median fare of survivors.
10. Get the median fare of non-survivors.

The Titanic dataset contains columns as shown below,

Pas sen ger id	Sur vive d	Pcla ss	Na me	Sex	Age	Sib Sp	Par ch	Tick et	Fare	Cab in	Em bark ed

Sample Test Cases

titanicDat...

```
8 # 1. Get the number of survivors by gender
9 survivors_by_gender = data[data['Survived'] == 1]['Sex'].value_counts()
10 print(survivors_by_gender)
11
12 # 2. Get the number of non-survivors by gender
13 non_survivors_by_gender = data[data['Survived'] == 0]
14 ['Sex'].value_counts()
15 print(non_survivors_by_gender)
16
17 # 3. Get the number of survivors by embarked location (Embarked_S)
18 survivors_by_embarked_s = data[data['Survived'] == 1]
```

Average time0.094 sMaximum time0.094 s1 out of 1 shown test case(s) passed

Test case 164 ms

Expected output

Actual output

female: 233

male: 109

Name: Sex, dtype: int64

male: 468

female: 81

Name: Sex, dtype: int64

TerminalTest cases

4.2.7. Titanic Dataset Analysis and Data Cleaning - 3

You are provided with the Titanic dataset containing information about passengers on the Titanic. Your task is to write Python code to answer the following questions based on the dataset.

1. Calculate the survival rate by class.
2. Calculate the survival rate by embarkation location (Embarked_S).
3. Calculate the survival rate by family size (FamilySize).
4. Calculate the survival rate by being alone (IsAlone).
5. Get the average fare by passenger class (Pclass).
6. Get the average age by passenger class (Pclass).
7. Get the average age by survival status (Survived).
8. Get the average fare by survival status (Survived).
9. Get the number of survivors by class (Pclass).
10. Get the number of non-survivors by class (Pclass).

The Titanic dataset contains columns as shown below,

Pas sen ger id	Sur vive d	Pcla ss	Na me	Sex	Age	Sib Sp	Par ch	Tick et	Fare	Cab in	Em bark ed

Sample Test Cases

titanicDat...

```
10 data = pd.read_csv('Titanic-Dataset.csv')
11 data['FamilySize'] = data['SibSp'] + data['Parch']
12 data['IsAlone'] = np.where(data['FamilySize'] > 0, 0, 1)
13 data = pd.get_dummies(data, columns=['Embarked'], drop_first=True)
14
15 print(data.groupby('Pclass')['Survived'].mean())
16
17 # 2. Calculate the survival rate by embarked location (Embarked_S)
18 print(data.groupby('Embarked_S')['Survived'].mean())
19
```

Average time0.130 sMaximum time0.130 s1 out of 1 shown test case(s) passed

Test case 1630 ms

Expected output

Actual output

Pclass:

1: 0.629630

2: 0.472826

3: 0.242363

Name: Survived, dtype: float64

Embarked_S:

0.629630

0.472826

0.242363

TerminalTest cases

Upload files · Harshada02 · cre...Course

mitaoe.codetantra.com/secure/course.jsp?euclid=6773e3f2f1f9c5320ca6bcd85#/contents/6773e455f1f9c5320ca6bd19/6773e4e4f1f9c5320ca6bdce/67e378a17028471d94dad0a6

CODETANTRAHome202401120055@mitaoe.ac.inSupportLogout

4.2.6. Titanic Dataset Analysis and Data Cleaning - 20/134

You are provided with the Titanic dataset containing information about passengers on the Titanic. Your task is to write Python code to answer the following questions based on the dataset.

- Create a new column 'IsAlone' which is 1 if the passenger is alone (FamilySize = 0), otherwise 0.
- Convert the 'Sex' column to numeric values (male: 0, female: 1).
- One-hot encode the 'Embarked' column, dropping the first category.
- Get the mean age of passengers.
- Get the median fare of passengers.
- Get the number of passengers by class.
- Get the number of passengers by gender.
- Get the number of passengers by survival status.
- Calculate the survival rate of passengers.
- Calculate the survival rate by gender.

The Titanic dataset contains columns as shown below,

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked

Sample Test Cases+

titanicDat...

import numpy as np

Load the Titanic dataset

data = pd.read_csv('Titanic-Dataset.csv')

data['FamilySize'] = data['SibSp'] + data['Parch']

data['Alone'] = np.where(data['FamilySize'] == 0, 1, 0)

2. Convert 'Sex' to numeric (male: 0, female: 1)

data['Sex'] = data['Sex'].map({'male': 0, 'female': 1})

Average time0.106 sMaximum time0.106 s1 out of 1 shown test case(s) passed

Test case 1106 ms

Expected output

29,69911764795882

14,4542

3...491

1...216

2...184

Name: Pclass, dtype: int64

Actual output

29,69911764795882

14,4542

3...491

1...216

2...184

Name: Pclass, dtype: int64

TerminalTest cases

Upload files · Harshada02 · cre...Course

mitaoe.codetantra.com/secure/course.jsp?euclid=6773e3f2f1f9c5320ca6bcd85#/contents/6773e455f1f9c5320ca6bd19/6773e4e4f1f9c5320ca6bdce/67e378a17028471d94dad0a6

CODETANTRAHome202401120055@mitaoe.ac.inSupportLogout

4.2.5. Titanic Dataset Analysis and Data Cleaning0/10

You are provided with the Titanic dataset containing information about passengers on the Titanic. Your task is to write Python code to answer the following questions based on the dataset. For each question, perform necessary data cleaning, transformations, and calculations as required.

- Display the first 5 rows of the dataset.
- Display the last 5 rows of the dataset.
- Get the shape of the dataset (number of rows and columns).
- Get a summary of the dataset (using .info()).
- Get basic statistics (mean, standard deviation, etc.) of the dataset using .describe().
- Check for missing values and display the count of missing values for each column.
- Fill missing values in the 'Age' column with the median age.
- Fill missing values in the 'Embarked' column with the most frequent value (mode).
- Drop the 'Cabin' column due to many missing values.
- Create a new column, 'FamilySize' by adding the 'SibSp' and 'Parch' columns.

The Titanic dataset contains columns as shown below,

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked

Sample Test Cases+

titanicDat...

import pandas as pd

import numpy as np

Load the Titanic dataset

data = pd.read_csv('Titanic-Dataset.csv')

print(data.head())

2. Display the last 5 rows of the dataset

print(data.tail())

Average time0.213 sMaximum time0.213 s1 out of 1 shown test case(s) passed

Test case 1213 ms

Expected output

PassengerId Survived Pclass Fare Cabin Embarked

0 1 0 3 7.2500 S

1 2 1 1 71.2833 C

2 3 1 3 7.9250 S

Actual output

PassengerId Survived Pclass Fare Cabin Embarked

0 1 0 3 7.2500 S

1 2 1 1 71.2833 C

2 3 1 3 7.9250 S

TerminalTest cases

Upload files · Harshada02-crea

Course

mitaoe.codetantra.com/secure/course.jsp?euclid=6773e3f2f1f9c5320ca6bc85#/contents/6773e455f1f9c5320ca6bd19/6773e4e4f1f9c5320ca6bdce/6788eea89f79b6ea9e3ece1

CODETANTRA Home202401120055@mitaoe.ac.inSupportLogout

4.2.4. Most Frequently Sold Product Pairs01:04

Write a Python program that takes the file name of a CSV file as input, reads the data, and performs the following operations:

- The CSV file contains the following columns: Date, Product, Quantity, Price, and City.
- For each date, find all pairs of products that were sold together (i.e., two products sold on the same date).
- Output the product pair/s that was sold most frequently.

Sample Data:

```
Date,Product,Quantity,Price,City
2025-01-01,Product A,5,20,New York
2025-01-01,Product B,3,15,Los Angeles
2025-01-02,Product A,7,20,New York
2025-01-02,Product C,4,30,Chicago
2025-01-03,Product B,2,15,Chicago
2025-01-03,Product A,8,20,Los Angeles
2025-01-04,Product C,6,30,New York
2025-01-04,Product B,5,15,Los Angeles
2025-01-05,Product A,3,20,Chicago
2025-01-05,Product C,10,30,Los Angeles
```

Explanation:
Transactions:

Sample Test Cases

frequentl...sales_dat...

```
11 # write the code
12 product_pairs = []
13
14 for _, group in df.groupby("Date"):
15     products = list(group["Product"].unique())
16     product_pairs.extend(combinations(sorted(products), 2)) # Generate
17     unique product pairs
18
19 # Count occurrences of each product pair
20 pair_counts = Counter(product_pairs)
```

Average time0.020 sMaximum time0.040 s
20.00 ms40.00 ms

1 out of 1 shown test case(s) passed
2 out of 2 hidden test case(s) passed

Test case 140 ms
Expected output
sales_data.csv
Product A and Product B: 2 times
Product A and Product C: 2 times
Actual output
sales_data.csv
Product A and Product B: 2 times
Product A and Product C: 2 times

TerminalTest cases

Upload files · Harshada02-crea

Course

mitaoe.codetantra.com/secure/course.jsp?euclid=6773e3f2f1f9c5320ca6bc85#/contents/6773e455f1f9c5320ca6bd19/6773e4e4f1f9c5320ca6bdce/67861553a50fe56b8dc9a6e

CODETANTRA Home202401120055@mitaoe.ac.inSupportLogout

4.2.3. City that Sold the Most Products04:30

Write a Python program that takes the file name of a CSV file as input, reads the data, and performs the following operations:

- The CSV file contains the columns: Date, Product, Quantity, Price, and City.
- Group the data by City and calculate the total quantity of products sold for each city.
- Find the city that sold the most products (based on the total quantity sold).

Sample Data:

```
Date,Product,Quantity,Price,City
2025-01-01,Product A,5,20,New York
2025-01-01,Product B,3,15,Los Angeles
2025-01-02,Product A,7,20,New York
2025-01-02,Product C,4,30,Chicago
2025-01-03,Product B,2,15,Chicago
2025-01-03,Product A,8,20,Los Angeles
2025-01-04,Product C,6,30,New York
2025-01-04,Product B,5,15,Los Angeles
2025-01-05,Product A,3,20,Chicago
2025-01-05,Product C,10,30,Los Angeles
```

Note:
The data cannot be displayed in the file. You can refer to the sample data provided for insights.

Sample Test Cases

monthFor...sales_dat...

```
8
9 # write the code..
10 city_sales = df.groupby("City")["Quantity"].sum()
11
12 # Find the city that sold the most products
13 best_city = city_sales.idxmax()
14 highest_quantity = city_sales.max()
15
16 # Display the result
17 print(f"City sold the most products: {best_city}")
18
```

Average time0.025 sMaximum time0.045 s
25.33 ms45.00 ms

1 out of 1 shown test case(s) passed
2 out of 2 hidden test case(s) passed

Test case 140 ms
Expected output
sales_data.csv
City sold the most products: Los Angeles
Actual output
sales_data.csv
City sold the most products: Los Angeles

TerminalTest cases

Upload files · Harshada02 · create · Course

mitaoe.codetantra.com/secure/course.jsp?euclid=6773e3f2f19c5320ca6bcd85#/contents/6773e455f19c5320ca6bd19/6773e4e4f19c5320ca6bdce/67861370a50fe56b8dd9828

CODETANTRA Home202401120055@mitaoe.ac.inSupportLogout

4.2.2. Best Selling Product14/34

Write a Python program that takes the file name of a CSV file as input, reads the data, and performs the following operations:

- The CSV file contains the columns: Date, Product, Quantity, Price, and City.
- Find the product that sold the most in terms of quantity sold.
- Display the product that sold the most and the total quantity sold for that product.

Sample Data:

Date	Product	Quantity	Price	City
2025-01-01	Product A	5	20	New York
2025-01-01	Product B	3	15	Los Angeles
2025-01-02	Product A	7	20	New York
2025-01-02	Product C	4	30	Chicago
2025-01-03	Product B	2	15	Chicago
2025-01-03	Product A	8	20	Los Angeles
2025-01-04	Product C	6	30	New York
2025-01-04	Product B	5	15	Los Angeles
2025-01-05	Product A	3	20	Chicago
2025-01-05	Product C	10	30	Los Angeles

Note:

The data cannot be displayed in the file. You can refer to the sample data provided for insights.

Sample Test Cases

monthFor...sales_data...

```
1 # Load the data
2 df = pd.read_csv(file_name)
3
4
5
6
7
8
9 product_sales = df.groupby("Product")["Quantity"].sum()
10 best_product = product_sales.idxmax()
11 highest_quantity = product_sales.max()
12
13 # Display the result
14 print(f"Best selling product: {best_product}")
15 print(f"Total quantity sold: {highest_quantity}")
```

Average time: 0.022 s22.00 msMaximum time: 0.040 s40.00 ms1 out of 1 shown test case(s) passed2 out of 2 hidden test case(s) passed

Test case 140 msExpected output: sales_data.csvActual output: sales_data.csvBest selling product: Product ABest selling product: Product ATotal quantity sold: 23Total quantity sold: 23

TerminalTest cases

Upload files · Harshada02 · create · Course

mitaoe.codetantra.com/secure/course.jsp?euclid=6773e3f2f19c5320ca6bcd85#/contents/6773e455f19c5320ca6bd19/6773e4e4f19c5320ca6bdce/67860994a50fe56b8dd8747

CODETANTRA Home202401120055@mitaoe.ac.inSupportLogout

4.2.1. Month with the Highest Total Sales03/40

Write a Python program that takes the file name of a CSV file as input, reads the data, and performs the following operations:

- The CSV file contains the columns: Date, Product, Quantity, Price, and City.
- Group the data by Month and calculate the total sales for each month.
- Find the month with the highest total sales and display it.
- Also, display the total sales for the best month.

Sample Data:

Date	Product	Quantity	Price	City
2025-01-01	Product A	5	20	New York
2025-01-01	Product B	3	15	Los Angeles
2025-01-02	Product A	7	20	New York
2025-01-02	Product C	4	30	Chicago
2025-01-03	Product B	2	15	Chicago
2025-01-03	Product A	8	20	Los Angeles
2025-01-04	Product C	6	30	New York
2025-01-04	Product B	5	15	Los Angeles
2025-01-05	Product A	3	20	Chicago
2025-01-05	Product C	10	30	Los Angeles

Note:

Sample Test Cases

monthFor...sales_data...

```
8 df["Date"] = pd.to_datetime(df["Date"])
9
10 df["Month"] = df["Date"].dt.to_period("M")
11
12 df["Total Sales"] = df["Quantity"] * df["Price"]
13
14 monthly_sales = df.groupby("Month")["Total Sales"].sum()
15
16
17 best_month = monthly_sales.idxmax()
18 highest_sales = monthly_sales.max()
```

Average time: 0.031 s30.67 msMaximum time: 0.060 s60.00 ms1 out of 1 shown test case(s) passed2 out of 2 hidden test case(s) passed

Test case 160 msExpected output: sales_data.csvActual output: sales_data.csvBest month: 2025-01Best month: 2025-01Total sales: \$1210.00Total sales: \$1210.00

TerminalTest cases

Upload files · Harshada02 · create · Course

mitaoe.codetantra.com/secure/course.jsp?eucdd=6773e3f2f1f9c5320ca6bcd85#/contents/6773e455f1f9c5320ca6bcd19/6773e4d5f1f9c5320ca6bcd2/667cf60f2365ba61af5a2ce7

CODETANTRA Home202401120055@mitaoe.ac.inSupportLogout

4.1.3. Student Information11/23

Write a program to read a text file containing student information (name, age, and grade) using Pandas. Perform the following tasks:

- Display the first five rows of the data frame.
- Calculate the average age of the students(limit the average age up to 2 decimal places).
- Filter out the students who have a grade above a certain threshold(consider the threshold grade is 'B').

Note:
Refer to the displayed test cases for better understanding.

Sample Test Cases

studentin...studentdat...

```
1 import pandas as pd
2
3 # Prompt the user for the file name
4 file_name = input()
5
6 # Load the data into a Pandas DataFrame
7 df = pd.read_csv(file_name, sep="\s+", names=["Name", "Age", "Grade"])
8
9 # Display the first five rows
10 print("First five rows:")
11 print(df.head())
```

Average time0.046 s46.90 msMaximum time0.062 s62.00 ms

1 out of 1 shown test case(s) passed1 out of 1 hidden test case(s) passed

Test case 162 ms

Expected outputActual output

studentdata.txtstudentdata.txt

First five rows:First five rows:

...Name Age Grade...Name Age Grade

0 John 25 A0 John 25 A

1 Allin 22 B1 Allin 22 B

2 Emma 24 A2 Emma 24 A

TerminalTest cases

Upload files · Harshada02 · create · Course

mitaoe.codetantra.com/secure/course.jsp?eucdd=6773e3f2f1f9c5320ca6bcd85#/contents/6773e455f1f9c5320ca6bcd19/6773e4d5f1f9c5320ca6bcd2/667aa30728220463f649bb62

CODETANTRA Home202401120055@mitaoe.ac.inSupportLogout

4.1.2. Dictionary to dataframe13/26

A dictionary of lists has been provided to you in the editor. Create a DataFrame from the dictionary of lists and perform the listed operations, then display the DataFrame before and after each manipulation.

Create the DataFrame:

- Convert the dictionary to a Pandas DataFrame.

Add a new row:

- Take inputs from the user for the new row data (name, age).
- Add the new row to the DataFrame.
- Display the DataFrame after adding the new row.

Modify a row:

- Modify a specific row by changing the age. Take the row index and new age value from the user.
- Display the DataFrame after modifying the row.

Delete a row:

- Take the row index to be deleted from the user.
- Remove the specified row.
- Display the DataFrame after deleting the row.

Add a new column:

Sample Test Cases

datafram...

```
16 # Adding a new row
17 new_name = input("New name: ")
18 new_age = int(input("New age: "))
19 df.loc[len(df)] = [new_name, new_age]
20
21 # Display the DataFrame after adding a new row
22 print("After adding a row:\n",df)
23
24 # Modifying a row
25 row_to_modify = int(input("Index of row to modify: "))
26 new_age = int(input("New age: "))
```

Average time0.102 s102.80 msMaximum time0.114 s114.00 ms

1 out of 1 shown test case(s) passed1 out of 1 hidden test case(s) passed

Test case 1114 ms

Expected outputActual output

Original DataFrame:Original DataFrame:

...Name Age...Name Age

0 Alice 250 Alice 25

1 Bob 301 Bob 30

2 Charlie 352 Charlie 35

New name: SusanNew name: Susan

TerminalTest cases

Upload files · Harshada02 · create · Course

mitaoe.codetantra.com/secure/course.jsp?euclid=6773e3f2f1f9c5320ca6bc85#/contents/6773e455f1f9c5320ca6bd19/6773e4d5f1f9c5320ca6bd2/667a7e0428220463f6497a3c

CODETANTRAHome202401120055@mitaoe.ac.inSupportLogout

4.1.1.1. Pandas - series creation and manipulation0.026s

Write a Python program that takes a list of numbers from the user, creates a Pandas series from it, and then calculates the mean of even and odd numbers separately using the `groupby` and `mean()` operations.

Input Format:

- The user should enter a list of numbers separated by space when prompted.

Output Format:

- The program should display the mean of even and odd numbers separately.
- Each mean value should be displayed with a label indicating whether it corresponds to even or odd numbers.

Sample Test Cases

seriesMa...

1import pandas as pd
2
3# Take inputs from the user to create a list of numbers
4numbers = list(map(int, input().split()))
5
6# Create a Pandas series from the list of numbers
7series = pd.Series(numbers)
8# Grouping by even and odd numbers and calculating the mean
9grouped = series.groupby(series % 2 == 0).mean()
10
11# Display the mean of even and odd numbers with labels

Average time0.013 s13.33 msMaximum time0.027 s27.00 ms

3 out of 3 shown test case(s) passed
3 out of 3 hidden test case(s) passed

Test case 10.027 ms
Expected output1 2 3 4 5 6 7 8 9 10
Mean of even and odd numbers: 5.0
Odd: 5.0
Even: 6.0
dtype: float64
Actual output1 2 3 4 5 6 7 8 9 10
Mean of even and odd numbers: 5.0
Odd: 5.0
Even: 6.0
dtype: float64

TerminalTest cases

PrevResetSubmitNext