

4.2.8. Titanic Dataset Analysis and Data Cleaning - 4

00:39

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,

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

Sample Test Cases

+

titanicDat...

Submit

```
1 import pandas as pd
2 import numpy as np
3
4 # Load the Titanic dataset
5 data = pd.read_csv('Titanic-Dataset.csv')
6 data = pd.get_dummies(data, columns=['Embarked'], drop_first=True)
7
8
9 survivors_by_gender = data[data['Survived'] == 1]['Sex'].value_counts()
10 print(survivors_by_gender)
11
12 non_survivors_by_gender = data[data['Survived'] == 0]
13 ['Sex'].value_counts()
14 print(non_survivors_by_gender)
15
16 #3. Get the number of survivors by embarked location (Embarked_S)
17 survivors_by_embarked_s = data[data['Survived'] == 1]
18 ['Embarked_S'].value_counts()
19 print(survivors_by_embarked_s)
20
21 #4. Get the number of non-survivors by embarked location (Embarked_S)
22 non_survivors_by_embarked_s = data[data['Survived'] == 0]
23 ['Embarked_S'].value_counts()
```

Terminal Test cases

Activate Windows
Go to Settings to activate Windows.

< Prev Reset Submit Next >

4.2.7. Titanic Dataset Analysis and Data Cleaning - 3

00:58

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,

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

Sample Test Cases

+

titanicDat...

Submit

```
1 import pandas as pd
2 import numpy as np
3
4 # Load the Titanic dataset
5 data = pd.read_csv('Titanic-Dataset.csv')
6 data['FamilySize'] = data['SibSp'] + data['Parch']
7 data['IsAlone'] = np.where(data['FamilySize'] > 0, 0, 1)
8 data = pd.get_dummies(data, columns=['Embarked'], drop_first=True)
9 # 1. Calculate the survival rate by class
10
11 print(data.groupby('Pclass')['Survived'].mean())
12 # 2. Calculate the survival rate by embarked location
13
14 print(data.groupby('Embarked_S')['Survived'].mean())
15 # 3. Calculate the survival rate by family size
16 print(data.groupby('FamilySize')['Survived'].mean())
17
18 # 4. Calculate the survival rate by being alone
19
20 print(data.groupby('IsAlone')['Survived'].mean())
21 # 5. Get the average fare by class
22 print(data.groupby('Pclass')['Fare'].mean())
23
24 # 6. Get the average age by class
25
26 print(data.groupby('Pclass')['Age'].mean())
```

Terminal

Test cases

Activate Windows
Go to Settings to activate Windows.

< Prev Reset Submit Next >

4.2.6. Titanic Dataset Analysis and Data Cleaning - 2

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

The Titanic dataset contains columns as shown below,

Pas sen ger id	Sur vive d	Pcl ass	Na me	Sex	Age	Sib Sp	Par ch	Tick et	Fare	Cab in	Em bar ked

Sample Test Cases +

```
titanicDat... Submit
1 import pandas as pd
2 import numpy as np
3
4 # Load the Titanic dataset
5 data = pd.read_csv('Titanic-Dataset.csv')
6 data['FamilySize'] = data['SibSp'] + data['Parch']
7
8 # 1. Create a new column 'IsAlone' (1 if alone, 0 otherwise)
9 data['IsAlone'] = np.where(data['FamilySize'] == 0, 1, 0)
10
11 # 2. Convert 'Sex' to numeric (male: 0, female: 1)
12 data['Sex'] = data['Sex'].map({'male': 0, 'female': 1})
13
14 # 3. One-hot encode the 'Embarked' column
15 data = pd.get_dummies(data, columns=['Embarked'])
16
17 # 4. Get the mean age of passengers
18 mean_age = data['Age'].mean()
19 print(mean_age)
20
21 # 5. Get the median fare of passengers
22 median_fare = data['Fare'].median()
23 print(median_fare)
24
25 # 6. Get the number of passengers by class
26 print(data['Pclass'].value_counts())
```

4.2.5. Titanic Dataset Analysis and Data Cleaning

00:32

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.

1. Display the first 5 rows of the dataset.
2. Display the last 5 rows of the dataset.
3. Get the shape of the dataset (number of rows and columns).
4. Get a summary of the dataset (using `.info()`).
5. Get basic statistics (mean, standard deviation, etc.) of the dataset using `.describe()`.
6. Check for missing values and display the count of missing values for each column.
7. Fill missing values in the 'Age' column with the median age.
8. Fill missing values in the 'Embarked' column with the most frequent value (mode).
9. Drop the 'Cabin' column due to many missing values.
10. 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...

Submit

```
1 import pandas as pd
2 import numpy as np
3
4 # Load the Titanic dataset
5 data = pd.read_csv('Titanic-Dataset.csv')
6
7 # 1. Display the first 5 rows of the dataset
8 print(data.head())
9
10 # 2. Display the last 5 rows of the dataset
11 print(data.tail())
12
13 # 3. Get the shape of the dataset
14
15 print(data.shape)
16 # 4. Get a summary of the dataset (info)
17
18 print(data.info())
19 # 5. Get basic statistics of the dataset
20 print(data.describe())
21
22 # 6. Check for missing values
23 print(data.isnull().sum())
24
25 # 7. Fill missing values in the 'Age' column with the median age
26 data['Age'].fillna(data['Age'].median(), inplace=True)
```

Terminal Test cases

Activate Windows
Go to Settings to activate Windows.

< Prev Reset Submit Next >

4.2.4. Most Frequently Sold Product Pairs

22:31

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:

• 2025-01-01: Product A, Product B

Sample Test Cases

+

```
1 import pandas as pd
2 from itertools import combinations
3 from collections import Counter
4
5 # Prompt user to input the file name
6 file_name = input()
7
8 # Read data from the specified CSV file
9 df = pd.read_csv(file_name)
10
11 # write the code
12
13 grouped=df.groupby('Date')['Product'].apply(list)
14 product_combination=[]
15 for products in grouped:
16     product_combination.extend(combinations(sorted(set(products)),2))
17 combinations_count=Counter(product_combination)
18 max_count=combinations_count.most_common(1)[0][1]
19 for combo, count in combinations_count.items():
20     if count==max_count:
21         print(f"{combo[0]} and {combo[1]}: {count} times")
22
23 # Output the most frequent product pairsd
```

Terminal Test cases

Activate Windows
Go to Settings to activate Windows.

< Prev Reset Submit Next >

4.2.3. City that Sold the Most Products

03:13

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

+

Explorer

monthFor... sales_dat...

Submit

```
1 import pandas as pd
2
3 # Prompt the user for the file name
4 file_name = input()
5
6 # Load the data
7 df = pd.read_csv(file_name)
8
9 # write the code..
10 city_sales = df.groupby("City")["Quantity"].sum()
11
12 # Find the city with the highest total quantity sold
13 best_city = city_sales.idxmax()
14 # Display the result
15 print(f"City sold the most products: {best_city}")
16
```

Debugger

Terminal Test cases

Activate Windows
Go to Settings to activate Windows.

[< Prev](#) [Reset](#) [Submit](#) [Next >](#)

4.2.2. Best Selling Product

07:35

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_dat...

Submit

```
1 import pandas as pd
2
3 # Prompt the user for the file name
4 file_name = input()
5
6 # Load the data
7 df = pd.read_csv(file_name)
8
9
10 # Find the product with the highest total quantity sold
11 product_sales = df.groupby("Product")["Quantity"].sum()
12 best_product = product_sales.idxmax()
13 highest_quantity = product_sales.max()
14
15 # Display the result
16 print(f"Best selling product: {best_product}")
17 print(f"Total quantity sold: {highest_quantity}")
18
```

Terminal Test cases

Activate Windows
Go to Settings to activate Windows.

[< Prev](#) [Reset](#) [Submit](#) [Next >](#)

4.2.1. Month with the Highest Total Sales

19:47

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:

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

Sample Test Cases

+

monthFor...

sales_dat...

Submit

```
1 import pandas as pd
2
3 # Prompt the user for the file name
4 file_name = input()
5
6 # Load the data
7 df = pd.read_csv(file_name)
8 df['Total_Sales'] = df['Quantity'] * df['Price']
9 df['Date'] = pd.to_datetime(df['Date'])
10 df['Month'] = df['Date'].dt.to_period('M')
11 Monthly_sales = df.groupby('Month')['Total_Sales'].sum()
12 best_month = Monthly_sales.idxmax()
13 highest_sales = Monthly_sales.max()
14 print(f"Best month: {best_month}")
15 print(f"Total sales: ${highest_sales:.2f}")
16
```

Terminal Test cases

Activate Windows
Go to Settings to activate Windows.

< Prev Reset Submit Next >