

Chapter 6 - Exercise 3: Titanic Disaster

Vào ngày 15 tháng 4 năm 1912, trong chuyến hải trình đầu tiên của mình, tàu Titanic đã chìm sau khi va chạm với một tảng băng trôi, đã có 1502 nạn nhân mãi mãi ra đi trong tổng số 2224 hành khách và thủy thủ đoàn.

Thông tin về Titanic Disaster có thể xem tại: <https://www.kaggle.com/c/titanic/data> (<https://www.kaggle.com/c/titanic/data>).

Dựa trên tập tin *train.csv*, hãy thực hiện các yêu cầu sau:

Yêu cầu

Câu 1:

- a) Đọc dữ liệu từ tập tin *train.csv* và lưu vào biến *titanic*. Hiển thị 5 dòng dữ liệu đầu của *titanic*
- b) Thiết lập cột index cho *titanic* là *PassengerId*.

Câu 2:

- Tạo pie chart thể hiện tỷ lệ hành khách nam/nữ trên tàu.

Câu 3:

- Cho biết có bao nhiêu người còn sống sót
- Cho biết tỷ lệ bao nhiêu người còn sống sót

Câu 4:

- 4a) Vẽ biểu đồ histogram của cột vé (*Fare*), bổ sung các thông tin *xlabel*, *ylabel*, *title*. Bạn nhận xét gì về biểu đồ vừa vẽ
- 4b) Vẽ biểu đồ histogram của cột vé (*Fare*). Bạn nhận xét gì về biểu đồ vừa vẽ

Câu 5:

- a) Xem thông tin mô tả (*describe*) của biến *Fare* và vẽ biểu đồ *Boxplot*
- b) Vẽ biểu đồ *Barplot* thể hiện tổng số hành khách trong mỗi class
- c) Vẽ biểu đồ *Barplot* thể hiện tổng số hành khách tính theo tỷ lệ trong mỗi class
- d) Vẽ biểu đồ *Barplot* thể hiện tuổi trung bình theo Nam/Nữ
- e) Vẽ biểu đồ *Barplot* thể hiện tuổi trung bình theo class
- f) Vẽ biểu đồ *Barplot* thể hiện tổng số hành khách và số chết trong mỗi class
- g) vẽ biểu đồ *StackedBar* thể hiện tổng số hành khách Nam/Nữ theo mỗi class

```
In [1]: 1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

```
In [2]: 1 # Câu 1:
2 # a) Đọc dữ liệu từ tập tin train.csv và lưu vào biến titanic.
3 titanic = pd.read_csv(r'data\train.csv', sep = ',')
4 # Hiển thị 5 dòng dữ liệu đầu của titanic
5 titanic.head()
```

Out[2]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	

```
In [3]: 1 # b) Thiết lập cột index cho titanic là PassengerId.
2 # Hiển thị lại 5 dòng dữ liệu đầu của titanic lúc này.
3 titanic.set_index('PassengerId', inplace=True)
```

In [4]:

1 titanic.head()

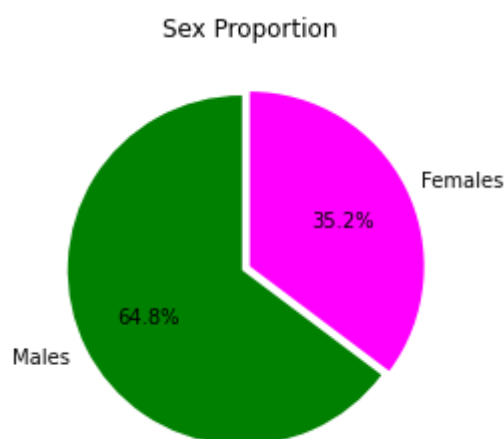
Out[4]:

	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarke
PassengerId											
1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	
2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	
3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	
4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	
5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	

```

In [5]: 1 # Câu 2: Tạo pie chart thể hiện tỷ lệ hành khách nam/nữ trên tàu.
        2 '''
        3 Gợi ý:
        4 - Tạo biến males, female là tổng nam và tổng nữ.
        5 - Tạo biến proportions là list có 2 phần tử là male và female
        6 - Vẽ biểu đồ: với dữ liệu là proportions, nhãn là ['Males', 'Female'],
        7     màu là ['green', 'magenta']
        8 - Thiết lập title là Sex Proportion
        9 '''
        10
        11 # Tạo biến males, female là tổng nam và tổng nữ
        12 males = (titanic['Sex'] == 'male').sum()
        13 females = (titanic['Sex'] == 'female').sum()
        14
        15 # Tạo biến proportions là list có 2 phần tử là male và female
        16 proportions = [males, females]
        17
        18 # Vẽ biểu đồ
        19 plt.pie(proportions, labels = ['Males', 'Females'], shadow = False,
        20         colors = ['green', 'magenta'],
        21         explode = (0.05, 0), startangle = 90, autopct = '%1.1f%%')
        22
        23 plt.axis('off')
        24
        25 # Thiết lập title là Sex Proportion
        26 plt.title("Sex Proportion")
        27
        28 # Show the plot
        29 plt.show()

```



```

In [6]: 1 # Câu 3: Cho biết có bao nhiêu người còn sống sót
        2

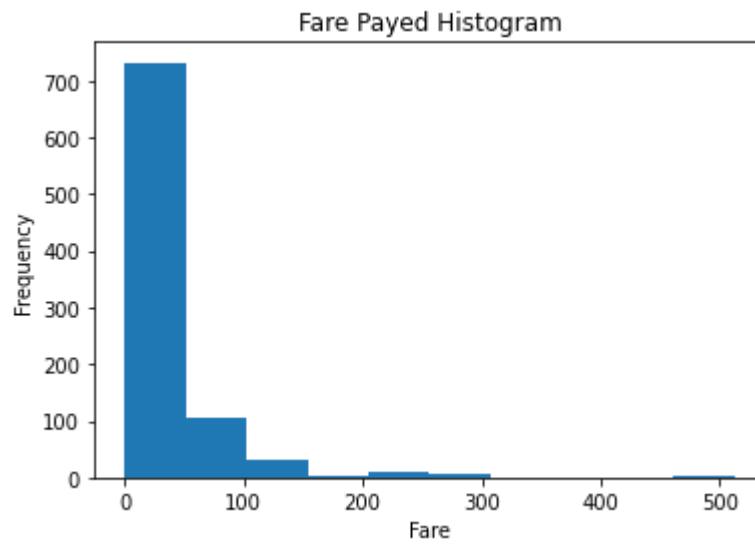
```

```

In [7]: 1 # Câu 4a: Vẽ biểu đồ histogram của cột vé (Fare),
        2 # bổ sung các thông tin xlabel, ylabel, title
        3 # Bạn nhận xét gì về biểu đồ vừa vẽ
        4

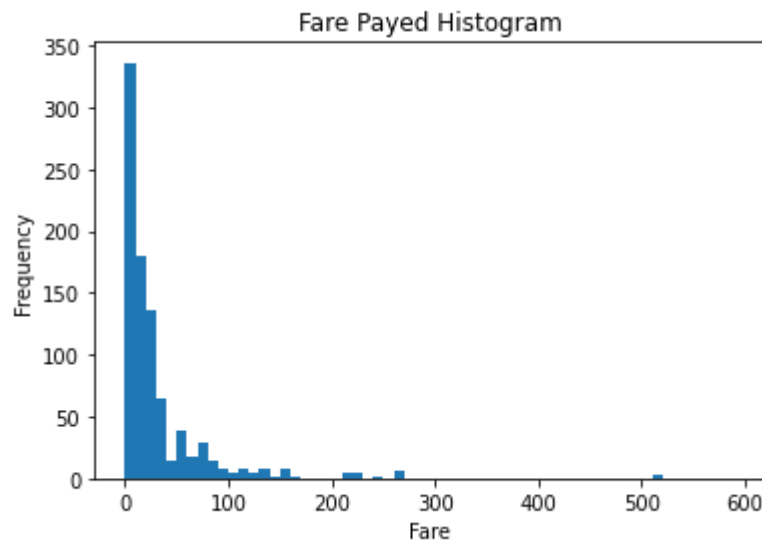
```

Nhấn vào đây để xem kết quả!



```
In [8]: 1 # Câu 4b: Vẽ biểu đồ histogram của cột vé (Fare)
2 # Bạn nhận xét gì về biểu đồ vừa vẽ
3 '''
4 Gợi ý:
5 - Tạo binsVal = np.arange(0,600,10)
6 - Vẽ histogram với dữ liệu là titanic, bins = binsVal.
7 - Bổ sung các thông tin xlabel, ylabel, title
8 '''
9 # Tạo binsVal = np.arange(0,600,10)
10 binsVal = np.arange(0,600,10)
11
```

Nhấn vào đây để xem kết quả!



```
In [9]: 1 # Câu 5a: xem thông tin mô tả (describe) của biến Fare và vẽ biểu đồ Boxplot
2
```

```
In [10]: 1 # Câu 5b: vẽ biểu đồ Barplot thể hiện tổng số hành khách trong mỗi class
2
```

```
In [11]: 1 # Câu 5c: vẽ biểu đồ Barplot
2 # thể hiện tổng số hành khách tính theo tỷ lệ trong mỗi class
3
```

In [12]:

1

Câu 5d: vẽ biểu đồ Barplot thể hiện tuổi trung bình theo Nam/Nữ

2

In [13]:

1

Câu 5e: vẽ biểu đồ Barplot thể hiện tuổi trung bình theo class

2

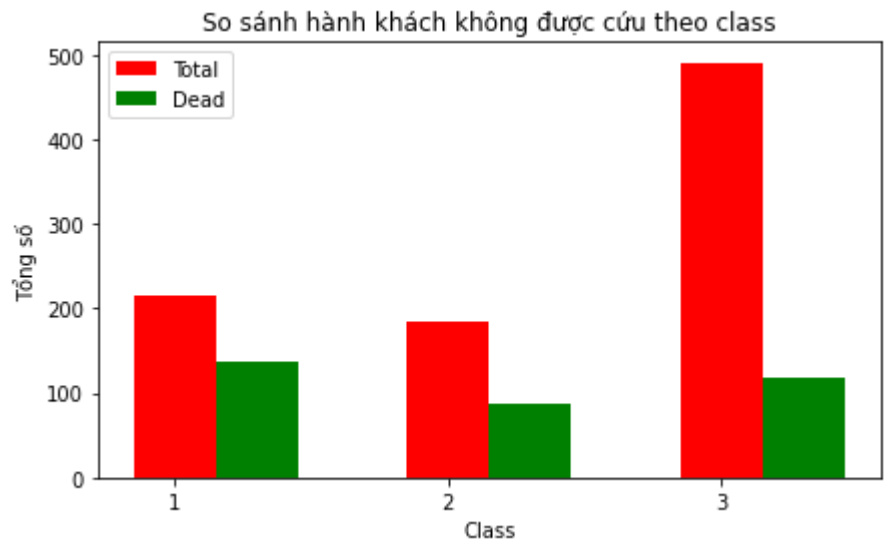
In [14]:

1

Câu 5f: vẽ biểu đồ Barplot thể hiện tổng số hành khách và số chết trong mỗi class

2

Nhấn vào đây để xem kết quả!



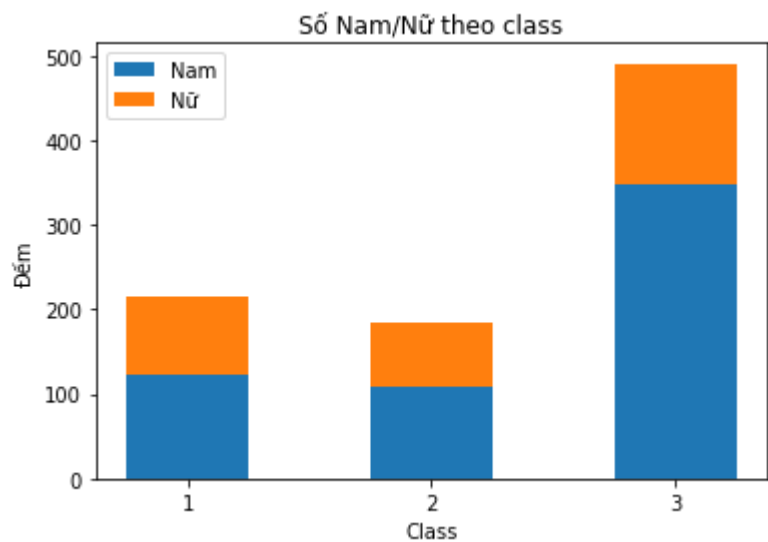
In [15]:

1

Câu 5g: vẽ biểu đồ StackedBar thể hiện tổng số hành khách Nam/Nữ theo mỗi class

2

Nhấn vào đây để xem kết quả!



In []:

1