

Topic 1. Exploratory data analysis with Pandas

Practice. Analyzing "Titanic" passengers

Fill in the missing code ("You code here").

```
In [1]: import numpy as np
import pandas as pd
from matplotlib import pyplot as plt

# Graphics in SVG format are more sharp and legible
%config InlineBackend.figure_format = 'svg'
pd.set_option("display.precision", 2)
```

Read data into a Pandas DataFrame

```
In [2]: data = pd.read_csv("titanic_train.csv", index_col="PassengerId")
```

First 5 rows

```
In [3]: data.head(5)
```

```
Out[3]:
```

	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
PassengerId									
1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.25
2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.28
3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.92
4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.10
5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.05

```
In [4]: data.describe()
```

Out[4]:

	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.00	891.00	714.00	891.00	891.00	891.00
mean	0.38	2.31	29.70	0.52	0.38	32.20
std	0.49	0.84	14.53	1.10	0.81	49.69
min	0.00	1.00	0.42	0.00	0.00	0.00
25%	0.00	2.00	20.12	0.00	0.00	7.91
50%	0.00	3.00	28.00	0.00	0.00	14.45
75%	1.00	3.00	38.00	1.00	0.00	31.00
max	1.00	3.00	80.00	8.00	6.00	512.33

Let's select those passengers who embarked in Cherbourg (Embarked=C) and paid > 200 pounds for their ticket (fare > 200).

Make sure you understand how actually this construction works.

```
In [5]: data[(data["Embarked"] == "C") & (data.Fare > 200)].head()
```

Out[5]:

	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fa
PassengerId									
119	0	1	Baxter, Mr. Quigg Edmond	male	24.0	0	1	PC 17558	247.0
259	1	1	Ward, Miss. Anna	female	35.0	0	0	PC 17755	512.0
300	1	1	Baxter, Mrs. James (Helene DeLaunieri Chaput)	female	50.0	0	1	PC 17558	247.0
312	1	1	Ryerson, Miss. Emily Borie	female	18.0	2	2	PC 17608	262.0
378	0	1	Widener, Mr. Harry Elkins	male	27.0	0	2	113503	211.0

We can sort these people by Fare in descending order.

```
In [6]: data[(data["Embarked"] == "C") & (data["Fare"] > 200)].sort_values(
        by="Fare", ascending=False
    ).head()
```

Out[6]:

	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	C
PassengerId										
259	1	1	Ward, Miss. Anna	female	35.0	0	0	PC 17755	512.33	
680	1	1	Cardeza, Mr. Thomas Drake Martinez	male	36.0	0	1	PC 17755	512.33	
738	1	1	Lesurer, Mr. Gustave J	male	35.0	0	0	PC 17755	512.33	
312	1	1	Ryerson, Miss. Emily Borie	female	18.0	2	2	PC 17608	262.38	
743	1	1	Ryerson, Miss. Susan Parker "Suzette"	female	21.0	2	2	PC 17608	262.38	

Let's create a new feature.

```
In [7]: def age_category(age):
        """
        < 30 -> 1
        >= 30, <55 -> 2
        >= 55 -> 3
        """
        if age < 30:
            return 1
        elif age < 55:
            return 2
        elif age >= 55:
            return 3
```

```
In [8]: age_categories = [age_category(age) for age in data.Age]
        data["Age_category"] = age_categories
```

Another way is to do it with `apply`.

```
In [9]: data["Age_category"] = data["Age"].apply(age_category)
```

1. How many men/women were there onboard?

- 412 men and 479 women
- 314 men and 577 women
- 479 men and 412 women
- **577 men and 314 women**

```
In [10]: print((data["Sex"] == "male").sum(), "men", "and", (data["Sex"] == "female").sum(), "women")
```

577 men and 314 women

2. Print the distribution of the `Pclass` feature. Then the same, but for men and women separately. How many men from second class were there onboard?

- 104
- **108**
- 112
- 125

```
In [11]: print("All:\n", data["Pclass"].value_counts())
```

```
All:
Pclass
3      491
1      216
2      184
Name: count, dtype: int64
```

```
In [12]: print("Men:\n", data[data["Sex"] == "male"]["Pclass"].value_counts())
```

```
Men:
Pclass
3      347
1      122
2      108
Name: count, dtype: int64
```

```
In [13]: print("Women:\n", data[data["Sex"] == "female"]["Pclass"].value_counts())
```

```
Women:
Pclass
3      144
1       94
2       76
Name: count, dtype: int64
```

```
In [14]: print("\nMen from second class:", ((data["Sex"] == "male") & (data["Pclass"] == 2)).sum())
```

Men from second class: 108

3. What are median and standard deviation of `Fare` ?. Round to two decimals.

- **median is 14.45, standard deviation is 49.69**
- median is 15.1, standard deviation is 12.15
- median is 13.15, standard deviation is 35.3
- median is 17.43, standard deviation is 39.1

```
In [15]: print(f"median is {round(data['Fare'].median(), 2)}, standard deviation is {round(data['Fare'].std(), 2)}")
```

median is 14.45, standard deviation is 49.69

4. Is that true that the mean age of survived people is higher than that of passengers who eventually died?

- Yes

- No

```
In [16]: print(data[data["Survived"] == 1]["Age"].mean() > data[data["Survived"] == 0]["Age"].mean())
```

False

5. Is that true that passengers younger than 30 y.o. survived more frequently than those older than 60 y.o.? What are shares of survived people among young and old people?

- 22.7% among young and 40.6% among old
- **40.6% among young and 22.7% among old**
- 35.3% among young and 27.4% among old
- 27.4% among young and 35.3% among old

```
In [17]: print(f"{round(data[data['Age'] < 30]['Survived'].mean() * 100, 1)}% among young and {round(data[data['Age'] >= 60]['Survived'].mean() * 100, 1)}% among old")
```

40.6% among young and 22.7% among old

6. Is that true that women survived more frequently than men? What are shares of survived people among men and women?

- 30.2% among men and 46.2% among women
- 35.7% among men and 74.2% among women
- 21.1% among men and 46.2% among women
- **18.9% among men and 74.2% among women**

```
In [18]: print(f"{round(data[data['Sex'] == 'male']['Survived'].mean() * 100, 1)}% among men and {round(data[data['Sex'] == 'female']['Survived'].mean() * 100, 1)}% among women")
```

18.9% among men and 74.2% among women

7. What's the most popular first name among male passengers?

- Charles
- Thomas
- **William**
- John

```
In [19]: def get_first_name(name: str) -> str:
    first_name = name.split(",")[1].split()[1]
    return first_name

first_names = [get_first_name(name) for name in data.Name]
data["First name"] = first_names

print(data["First name"].value_counts().idxmax())
```

William

8. How is average age for men/women dependent on **Pclass** ? Choose all correct statements:

- **On average, men of 1 class are older than 40**
- On average, women of 1 class are older than 40
- **Men of all classes are on average older than women of the same class**

- On average, passengers of the first class are older than those of the 2nd class who are older than passengers of the 3rd class

```
In [20]: print(data[(data["Sex"] == "male") & (data["Pclass"] == 1)]["Age"].mean())
print(data[(data["Sex"] == "female") & (data["Pclass"] == 1)]["Age"].mean())
print(all([
    data[(data["Sex"] == "male") & (data["Pclass"] == class_num)]["Age"].
    data[(data["Sex"] == "female") & (data["Pclass"] == class_num)]["Age"]
    for class_num in data["Pclass"].unique()
]))
print(data[data["Pclass"] == 1]["Age"].mean() > data[data["Pclass"] == 2]
```

True

False

True

True