

# PracticalNotebook2

January 27, 2024

## 1 Practical Notebook 2

### 1.1 Pandas

In this course, we will use pandas to import the data into DataFrame objects.

Pandas is a commonly used library working with and manipulating data in various formats, such as txt, csv, excel format, and more.

You can read more about pandas [here](#), or by searching online.

```
[ ]: # The first thing we need to do is to import pandas
import pandas as pd

# We will also change how the floating point numbers are displayed
pd.set_option("display.float_format", lambda x: f"{x:.5f}")
```

#### 1.1.1 Creating our own dataset to file

We will start by creating our own data set, but later on we will import the data from a file.

```
[ ]: names = ['Alice', 'Bob', 'Charlie']
animals = ['Dog', 'Cat', None]
age = [27, 12, 43]
sex = ['Female', 'Male', 'Male']
```

We will then merge the lists together using the *zip* function.

```
[ ]: people = list(zip(names, animals, age, sex))
print(people)
```

```
[('Alice', 'Dog', 27, 'Female'), ('Bob', 'Cat', 12, 'Male'), ('Charlie', None, 43, 'Male')]
```

Now we can make our merged list into a DataFrame object by using pandas.

```
[ ]: df = pd.DataFrame(data=people, columns=['Names', 'Animals', 'Age', 'Sex'])
print(df)
```

	Names	Animals	Age	Sex
0	Alice	Dog	27	Female

1	Bob	Cat	12	Male
2	Charlie	None	43	Male

You can also export the dataframe to a csv file, where we use the function `to_csv` to export the file. You will find the file you created in the folder you are in. (In colab you will find the folder to the left.) The index parameter is set to *False*, i.e. we won't write the row names to the new file (in this case the row names are *0, 1, 2*). The header parameter is set to *True*, i.e. we will write the column names to the file (in this case the column names are *Names, Animals, Age, Sex*). You can change these parameters yourself to see the difference.

```
[ ]: df.to_csv('test_people.csv', index=False, header=True)
```

### 1.1.2 Read a dataset from file

To read the data from a csv file we will use the function `read_csv`.

```
[ ]: df = pd.read_csv('test_people.csv')
print(df)
```

	Names	Animals	Age	Sex
0	Alice	Dog	27	Female
1	Bob	Cat	12	Male
2	Charlie	NaN	43	Male

We can inspect the numerical values in the data using the function `describe`.

```
[ ]: print(df.describe())
```

```

      Age
count  3.00000
mean   27.33333
std    15.50269
min    12.00000
25%    19.50000
50%    27.00000
75%    35.00000
max    43.00000
```

And look at one specific column by using the names of the header.

```
[ ]: print(f"Here you will see the names: \n{df['Names']}")
print(f"\nHere you will see the animals: \n{df['Animals']}")
print(f"\nHere you will see the ages: \n{df['Age']}")
print(f"\nHere you will see the sex: \n{df['Sex']}")
```

```

Here you will see the names:
0      Alice
1        Bob
2     Charlie
Name: Names, dtype: object
```

Here you will see the animals:

```
0    Dog
1    Cat
2    NaN
Name: Animals, dtype: object
```

Here you will see the ages:

```
0    27
1    12
2    43
Name: Age, dtype: int64
```

Here you will see the sex:

```
0    Female
1     Male
2     Male
Name: Sex, dtype: object
```

You can also divide the groups into females and males.

```
[ ]: male, female = df['Sex'].value_counts()
      print(f"Here we have {male} male(s) and {female} female(s).")
```

Here we have 2 male(s) and 1 female(s).

By looking only at one column, as we did before, we can find some interesting data about it as well.

```
[ ]: # finding the mean value of the ages (with 2 decimals)
      print(f"mean: {df['Age'].mean():.2f}")
      # and the standard deviation (with 2 decimals)
      print(f"std: {df['Age'].std():.2f}")
```

mean: 27.33

std: 15.50

### 1.1.3 Titanic

Now we will download and use a larger dataset, to get a better understanding about the pandas library. The dataset contains passenger data from Titanic, and later on we will predict “what sort of people were most likely to survive?”. The passenger data has 7 features: Name, Sex, Socio-economic class, Siblings/Spouses Aboard, Parents/Children Aboard and Fare and a binary response variable “survived”.

```
[ ]: # Downloading the titanic dataset
      !wget https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/stuff/titanic.
      ↪ csv
```

--2024-01-27 20:19:44--

<https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/stuff/titanic.csv>

```
Resolving web.stanford.edu (web.stanford.edu)... 171.67.215.200,
2607:f6d0:0:925a::ab43:d7c8
Connecting to web.stanford.edu (web.stanford.edu)|171.67.215.200|:443...
connected.
HTTP request sent, awaiting response... 200 OK
Length: 44225 (43K) [text/csv]
Saving to: 'titanic.csv'
```

```
titanic.csv          100%[=====>]  43.19K   269KB/s   in 0.2s
```

```
2024-01-27 20:19:44 (269 KB/s) - 'titanic.csv' saved [44225/44225]
```

### Assignment a)

```
[ ]: # ASSIGNMENT:
      # Load the data and get familiar with it
      # Use the .describe() method to inspect numerical values

df = pd.read_csv("titanic.csv")
print(df.head(2), "\n\n")
print(df.describe())
```

	Survived	Pclass	Name \
0	0	3	Mr. Owen Harris Braund
1	1	1	Mrs. John Bradley (Florence Briggs Thayer) Cum...

	Sex	Age	Siblings/Spouses Aboard	Parents/Children Aboard	Fare
0	male	22.00000	1	0	7.25000
1	female	38.00000	1	0	71.28330

	Survived	Pclass	Age	Siblings/Spouses Aboard \
count	887.00000	887.00000	887.00000	887.00000
mean	0.38557	2.30552	29.47144	0.52537
std	0.48700	0.83666	14.12191	1.10467
min	0.00000	1.00000	0.42000	0.00000
25%	0.00000	2.00000	20.25000	0.00000
50%	0.00000	3.00000	28.00000	0.00000
75%	1.00000	3.00000	38.00000	1.00000
max	1.00000	3.00000	80.00000	8.00000

	Parents/Children Aboard	Fare
count	887.00000	887.00000
mean	0.38331	32.30542
std	0.80747	49.78204
min	0.00000	0.00000
25%	0.00000	7.92500

50%	0.00000	14.45420
75%	0.00000	31.13750
max	6.00000	512.32920

#### Assignment b)

```
[ ]: # ASSIGNMENT:
      # Count the number of males and females

male, female = df['Sex'].value_counts()
print(f"Number of males: {male}, Number of females: {female}")
```

Number of males: 573, Number of females: 314

#### Assignment c)

```
[ ]: # ASSIGNMENT:
      # Find the mean fare and display with 2 floating point numbers

mean_fare = df["Fare"].mean()
print(f"mean: {mean_fare:.2f}")

# Find the standard deviation of the fare and display with 2 floating point
↪ numbers

std_fare = df["Fare"].std()
print(f"std: {std_fare:.2f}")
```

mean: 32.31

std: 49.78

#### Assignment d)

```
[ ]: # ASSIGNMENT:
      # Count how many survived (1) and how many died (0)

      # YOUR CODE HERE

died, survived = df["Survived"].value_counts()
print(f"Number of people who died: {died}, Number of people who survived:
↪ {survived}")
```

Number of people who died: 545, Number of people who survived: 342

#### Assignment e)

```
[ ]: # ASSIGNMENT:
      # count and display the number of women who survived
      # and the number of men who survived
```

```
# YOUR CODE HERE
```

```
female_survived, male_survived = df[df["Survived"] == 1]["Sex"].value_counts()
print(f"Number of female survivors: {female_survived}, Number of male survivors:
↪ {male_survived}")
```

Number of female survivors: 233, Number of male survivors: 109

### Assignment f)

```
[ ]: # ASSIGNMENT:
# Separate the dataset from Titanic into X and y,
# where y is the column Survived, and X is the rest.
# Inspect the data. Look at for instance the function "describe" in pandas
```

```
# YOUR CODE HERE
```

```
X = df["Survived"].copy()
y = df.drop("Survived", axis=1)

x_describe = X.describe()
y_describe = y.describe()

print(x_describe, "\n\n", y_describe)
```

```
count    887.00000
mean      0.38557
std       0.48700
min       0.00000
25%       0.00000
50%       0.00000
75%       1.00000
max       1.00000
```

Name: Survived, dtype: float64

	Pclass	Age	Siblings/Spouses Aboard	Parents/Children Aboard	\
count	887.00000	887.00000	887.00000	887.00000	
mean	2.30552	29.47144	0.52537	0.38331	
std	0.83666	14.12191	1.10467	0.80747	
min	1.00000	0.42000	0.00000	0.00000	
25%	2.00000	20.25000	0.00000	0.00000	
50%	3.00000	28.00000	0.00000	0.00000	
75%	3.00000	38.00000	1.00000	0.00000	
max	3.00000	80.00000	8.00000	6.00000	

```
Fare
count 887.00000
```

```

mean    32.30542
std     49.78204
min     0.00000
25%     7.92500
50%    14.45420
75%    31.13750
max    512.32920

```

### Assignment g)

```

[ ]: # ASSIGNMENT:
# Standardize the data by subtracting the mean and dividing by the standard
    ↪ deviation.
# Inspect the data again to see that the mean is (close to) zero and the
    ↪ standard deviation is one.

# YOUR CODE HERE

X_new = (X - X.mean())/X.std()
y_new = y.copy()
y_new["Pclass"] -= y_new["Pclass"].mean()
y_new["Pclass"] /= y_new["Pclass"].std()
y_new["Age"] -= y_new["Age"].mean()
y_new["Age"] /= y_new["Age"].std()
y_new["Siblings/Spouses Aboard"] -= y_new["Siblings/Spouses Aboard"].mean()
y_new["Siblings/Spouses Aboard"] /= y_new["Siblings/Spouses Aboard"].std()
y_new["Parents/Children Aboard"] -= y_new["Parents/Children Aboard"].mean()
y_new["Parents/Children Aboard"] /= y_new["Parents/Children Aboard"].std()
y_new["Fare"] -= y_new["Fare"].mean()
y_new["Fare"] /= y_new["Fare"].std()

# Inspecting the data again:
X_new_describe = X_new.describe()
y_new_describe = y_new.describe()

print(X_new_describe, y_new_describe)

```

```

count    887.00000
mean      0.00000
std       1.00000
min      -0.79172
25%      -0.79172
50%      -0.79172
75%       1.26165
max       1.26165
Name: Survived, dtype: float64      Pclass      Age  Siblings/Spouses
Aboard  Parents/Children Aboard  \

```

count	887.00000	887.00000	887.00000	887.00000
mean	-0.00000	0.00000	-0.00000	-0.00000
std	1.00000	1.00000	1.00000	1.00000
min	-1.56040	-2.05719	-0.47559	-0.47471
25%	-0.36517	-0.65299	-0.47559	-0.47471
50%	0.83006	-0.10420	-0.47559	-0.47471
75%	0.83006	0.60392	0.42966	-0.47471
max	0.83006	3.57803	6.76640	6.95594

	Fare
count	887.00000
mean	0.00000
std	1.00000
min	-0.64894
25%	-0.48974
50%	-0.35859
75%	-0.02346
max	9.64251

## 1.2 Matplotlib

Matplotlib is a commonly used library for visualizing data in Python. Other visualization libraries exist for Python, such as seaborn, plotly, and more. Beyond the first practical notebook, we do not enforce any particular plotting library, but strongly encourage the use of Matplotlib. Below we will use the plotting functions inside of *matplotlib.pyplot*. You can read more about matplotlib [here](#) and pyplot [here](#).

### 1.2.1 Examples

```
[ ]: # import the relevant libraries
import matplotlib.pyplot as plt
import numpy as np
```

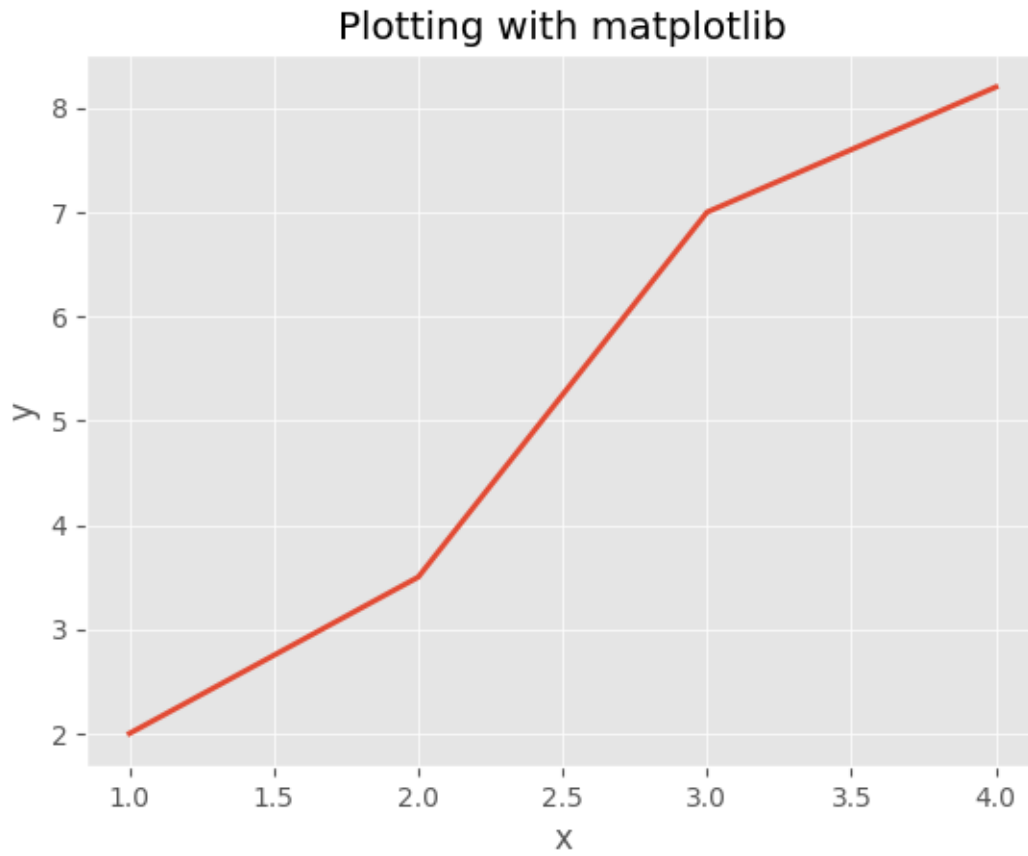
We will start by looking at some small lists.

```
[ ]: # examples of some datapoint
x = [1,2,3,4]
y = [2,3.5,7,8.2]

# plotting the data using matplotlib.pyplot.plot
plt.plot(x, y)

# It is important to add labels for the axes and a title
plt.xlabel("x")
plt.ylabel("y")
plt.title("Plotting with matplotlib")
# and always end with show(), which will show you the plot.
plt.show()
```





Plots can also be below each other, or side by side by using [subplot](#).

```
[ ]: # Vertical subplot

plt.style.use('bmh')

t = np.arange(0.0, 1.0, 0.01)
sin = np.sin(2*np.pi*t)
cos = np.cos(2*np.pi*t)

fig = plt.figure()
fig.suptitle("Sine and cosine for different t", fontsize=18)

ax1 = fig.add_subplot(2,1,1)
ax1.plot(t, sin, color='red', lw=2)
ax1.set_ylabel('Amplitude')
ax1.set_xlabel('Time')
ax1.set_title('Sine wave')

ax2 = fig.add_subplot(2,1,2)
```

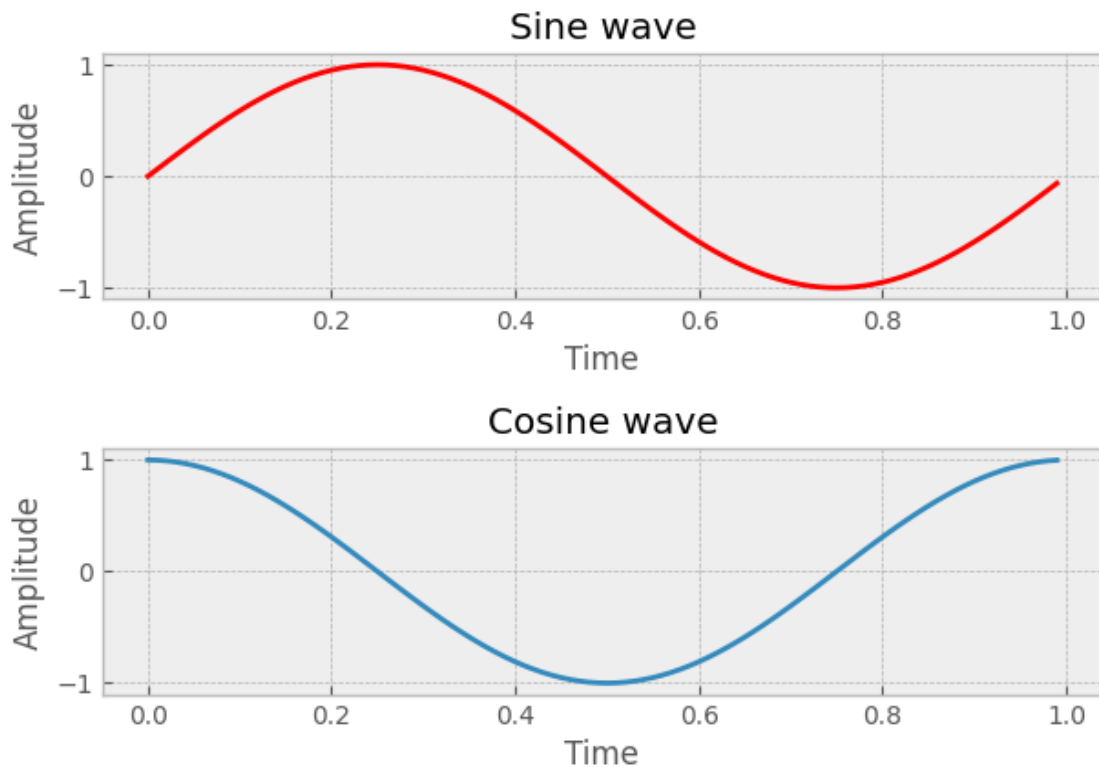
```

ax2.plot(t, cos)
ax2.set_ylabel('Amplitude')
ax2.set_xlabel('Time')
ax2.set_title('Cosine wave')

fig.tight_layout() # comment out this line to see the difference
fig.subplots_adjust(top=0.85)
plt.show()

```

## Sine and cosine for different t



```

[ ]: # Horizontal subplot

plt.style.use('bmh')

t = np.arange(0.0, 1.0, 0.01)
sin = np.sin(2*np.pi*t)
cos = np.cos(2*np.pi*t)

fig = plt.figure()
fig.suptitle("Sine and cosine for different t", fontsize=18)

```

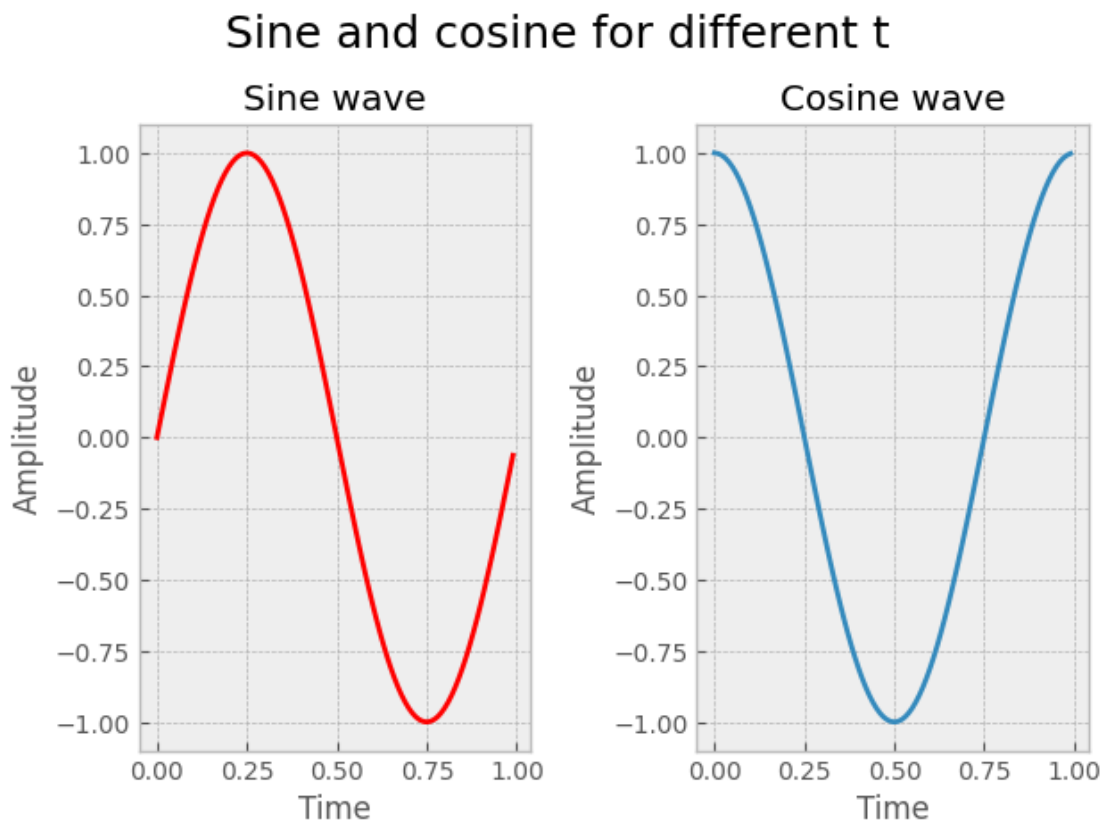
```

ax1 = fig.add_subplot(1,2,1) # we have changed (2,1,1) to (1,2,1)
ax1.plot(t, sin, color='red', lw=2)
ax1.set_ylabel('Amplitude')
ax1.set_xlabel('Time')
ax1.set_title('Sine wave')

ax2 = fig.add_subplot(1,2,2) # we have changed (2,1,2) to (1,2,2)
ax2.plot(t, cos)
ax2.set_ylabel('Amplitude')
ax2.set_xlabel('Time')
ax2.set_title('Cosine wave')

fig.tight_layout() # comment out this line to see the difference
fig.subplots_adjust(top=0.85)
plt.show()

```



And with different stylings

```

[ ]: # Here are all the different "pre-configured" styles matplotlib lib supports
# https://matplotlib.org/tutorials/intermediate/artists.
    ↪html#sphx-glr-tutorials-intermediate-artists-py

```

```
plt.style.available
```

```
[ ]: ['Solarize_Light2',
      '_classic_test_patch',
      '_mpl-gallery',
      '_mpl-gallery-nogrid',
      'bmh',
      'classic',
      'dark_background',
      'fast',
      'fivethirtyeight',
      'ggplot',
      'grayscale',
      'seaborn-v0_8',
      'seaborn-v0_8-bright',
      'seaborn-v0_8-colorblind',
      'seaborn-v0_8-dark',
      'seaborn-v0_8-dark-palette',
      'seaborn-v0_8-darkgrid',
      'seaborn-v0_8-deep',
      'seaborn-v0_8-muted',
      'seaborn-v0_8-notebook',
      'seaborn-v0_8-paper',
      'seaborn-v0_8-pastel',
      'seaborn-v0_8-poster',
      'seaborn-v0_8-talk',
      'seaborn-v0_8-ticks',
      'seaborn-v0_8-white',
      'seaborn-v0_8-whitegrid',
      'tableau-colorblind10']
```

The plots can also be both below each other and side by side at the same time (as a matrix) as you can see below. Here we have also plotted two graphs together in every figure, and added a color and a label for each one of them.

```
[ ]: # Matrix subplot

fig = plt.figure()
fig.suptitle("Sine and cosine for different t", fontsize=18)

i = 1
for freq in [1, 2, 3]:
    for t_max in [1, 2]:
        t = np.arange(0.0, t_max, 0.01)
        sin = np.sin(2*freq*np.pi*t)
        cos = np.cos(2*freq*np.pi*t)
```

```

ax = fig.add_subplot(3,2,i)
ax.plot(t, sin, color='red', lw=2, label='sine')
ax.plot(t, cos, color='blue', lw=2, label='cosine')
ax.set_ylabel('Amplitude')
ax.set_xlabel('Time')
ax.legend(fontsize=6)
ax.set_title(f'freq = {freq}', fontsize=10)
i += 1

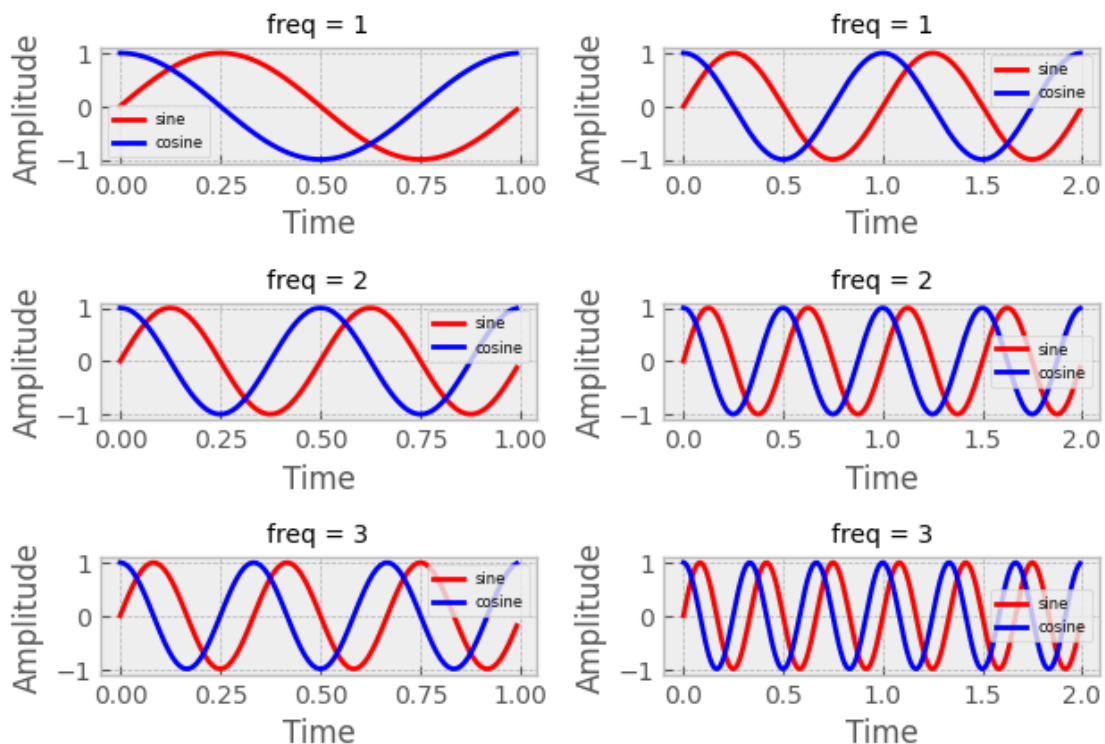
```

```

fig.tight_layout() # comment out this line to see the difference
fig.subplots_adjust(top=0.85)
plt.show()

```

## Sine and cosine for different t



### 1.2.2 Plotting data from Pandas

Now we will plot some of the datapoints from the titanic dataset to visualize it.

```

[ ]: # Downloading the titanic dataset
!wget https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/stuff/titanic.
     ↪ csv

```

```
--2024-01-27 20:19:46--
https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/stuff/titanic.csv
Resolving web.stanford.edu (web.stanford.edu)... 171.67.215.200,
2607:f6d0:0:925a::ab43:d7c8
Connecting to web.stanford.edu (web.stanford.edu)|171.67.215.200|:443...
connected.
HTTP request sent, awaiting response... 200 OK
Length: 44225 (43K) [text/csv]
Saving to: 'titanic.csv.1'

titanic.csv.1      100%[=====>]  43.19K   270KB/s   in 0.2s

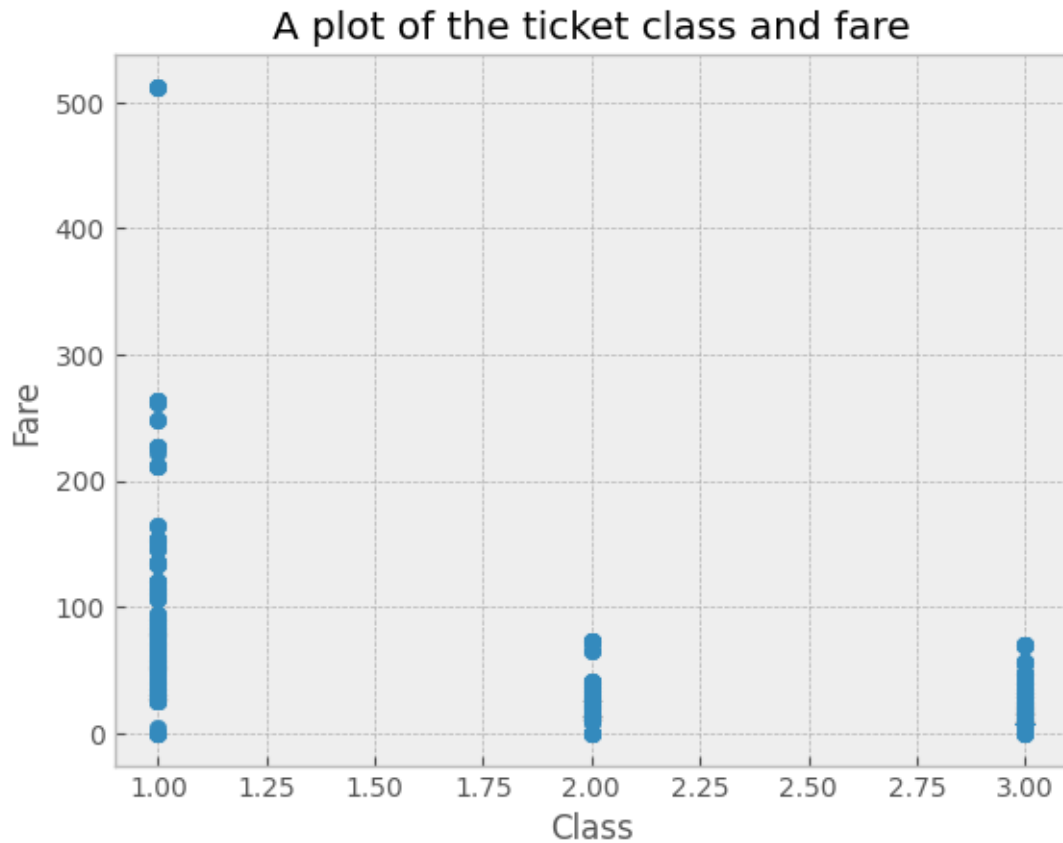
2024-01-27 20:19:47 (270 KB/s) - 'titanic.csv.1' saved [44225/44225]
```

```
[ ]: # Load the titanic dataset for plotting
import pandas as pd
df = pd.read_csv('titanic.csv')
```

#### Assignment h)

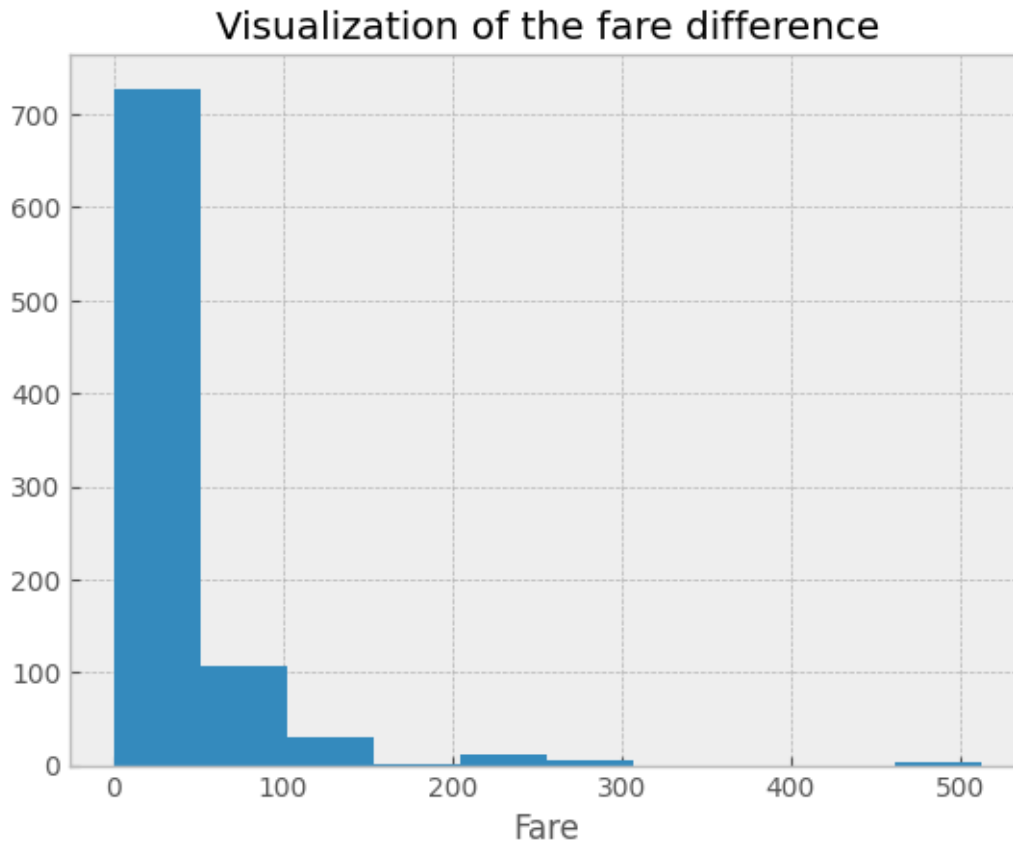
```
[ ]: # ASSIGNMENT:
# make a scatterplot of the class of ticket in the x-axis
# and the fare on the y-axis
# label the plot and the axes appropriately

plt.scatter(df["Pclass"], df["Fare"])
plt.xlabel("Class")
plt.ylabel("Fare")
plt.title("A plot of the ticket class and fare")
plt.show()
```



**Assignment i)** It might also be a good idea to plot a histogram over the data, to get a better understanding of how the data looks. This can be done using the function *hist* from matplotlib.

```
[ ]: fare = df["Fare"]
plt.hist(fare)
plt.xlabel("Fare")
plt.title("Visualization of the fare difference")
plt.show()
```

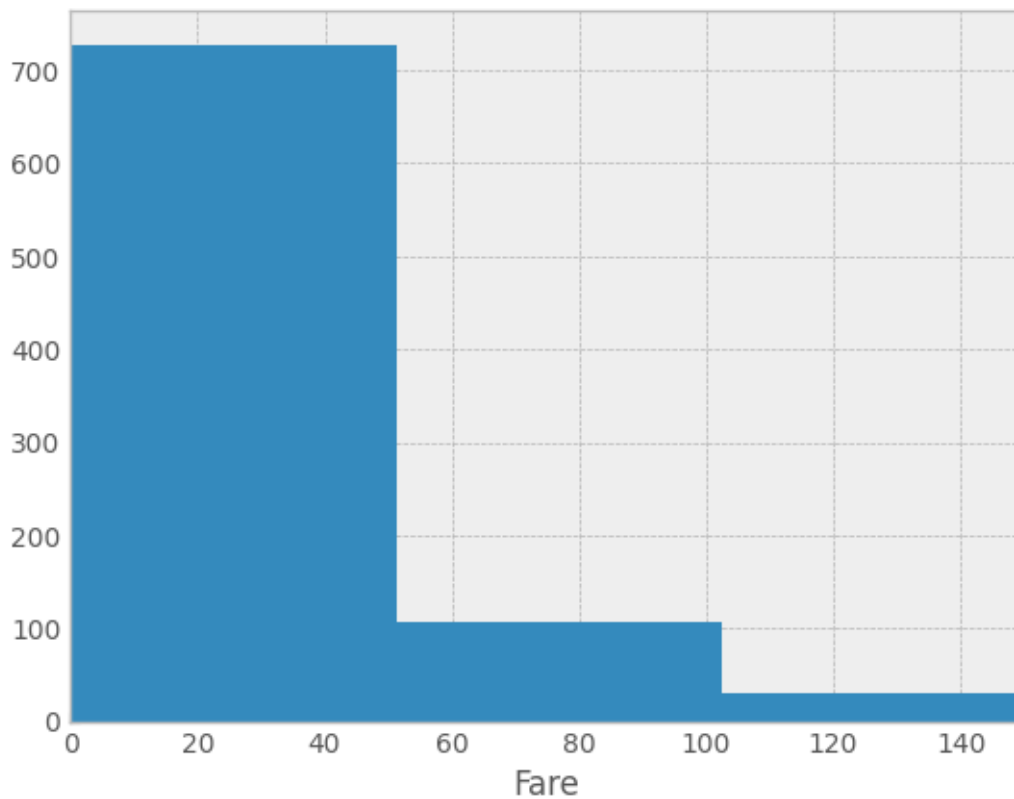


As you can see, most of the people paid less than 150 for the ticket.

```
[ ]: # ASSIGNMENT:  
# Plot a histogram over the people who paid less than, or equal to, 150.  
# label the plot and the axes appropriately  
  
fare = df["Fare"]  
plt.hist(fare)  
plt.xlabel("Fare")  
plt.xlim(0, 150)  
plt.title("Visualization of the fare difference between 0 and 150")  
plt.show()
```



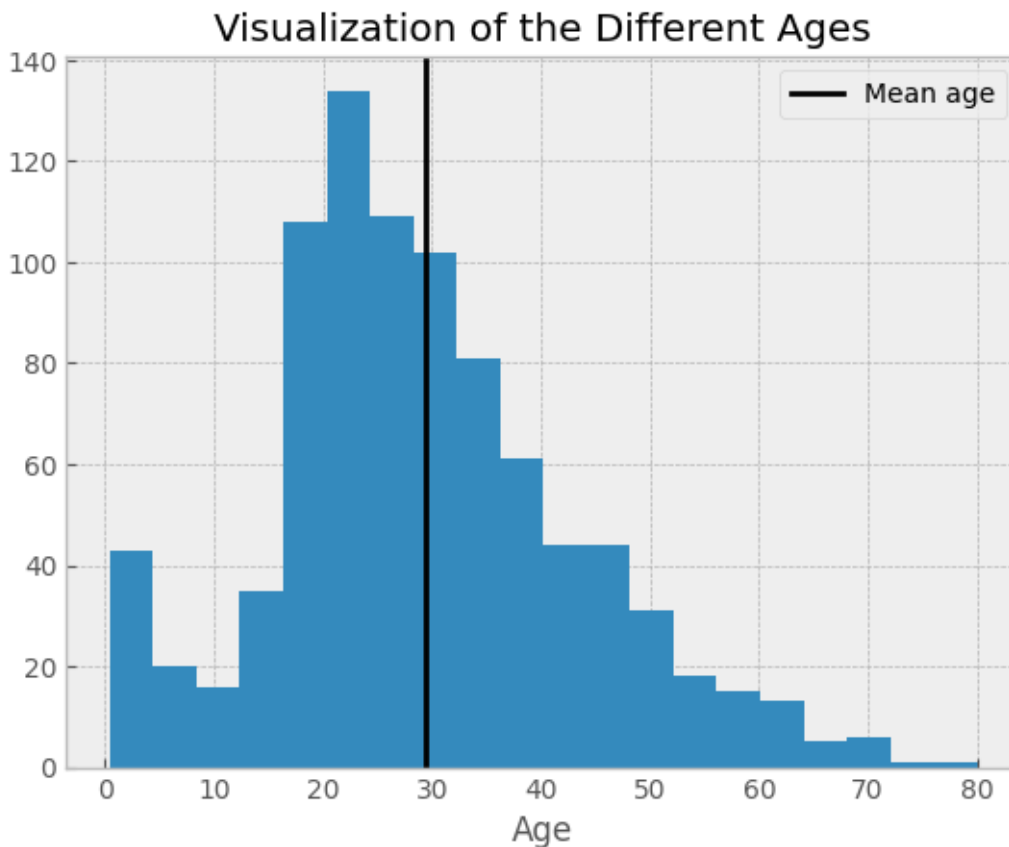
Visualization of the fare difference between 0 and 150



#### Assignment j)

```
[ ]: # ASSIGNMENT:
# plot a histogram over all the ages with 20 bins. Draw a vertical line at the
# mean age.
# label the plot and the axes appropriately

plt.hist(df["Age"], 20)
plt.xlabel("Age")
plt.title("Visualization of the Different Ages")
plt.axvline(x = df["Age"].mean(), color = "k", label="Mean age")
plt.legend()
plt.show()
```



**Assignment k)** Sometimes it is better to plot the figures together in one figure instead. This can be done with subplot, as shown in the examples above.

```
[ ]: # ASSIGNMENT:
# Make a subplot over the Fare, Class, and Age
# label the plot and the axes appropriately

fig = plt.figure()
fig.suptitle("Histograms for Fare, Class, and Age", fontsize=18)

ax1 = fig.add_subplot(2, 2, 1) # we have changed (2,1,1) to (1,2,1)
ax1.hist(df["Fare"])
ax1.set_xlabel('Fare')
ax1.set_title('Visualisation of the fares')

ax2 = fig.add_subplot(2, 2, 2) # we have changed (2,1,2) to (1,2,2)
ax2.hist(df["Pclass"])
ax2.set_xlabel('Pclass')
```

```

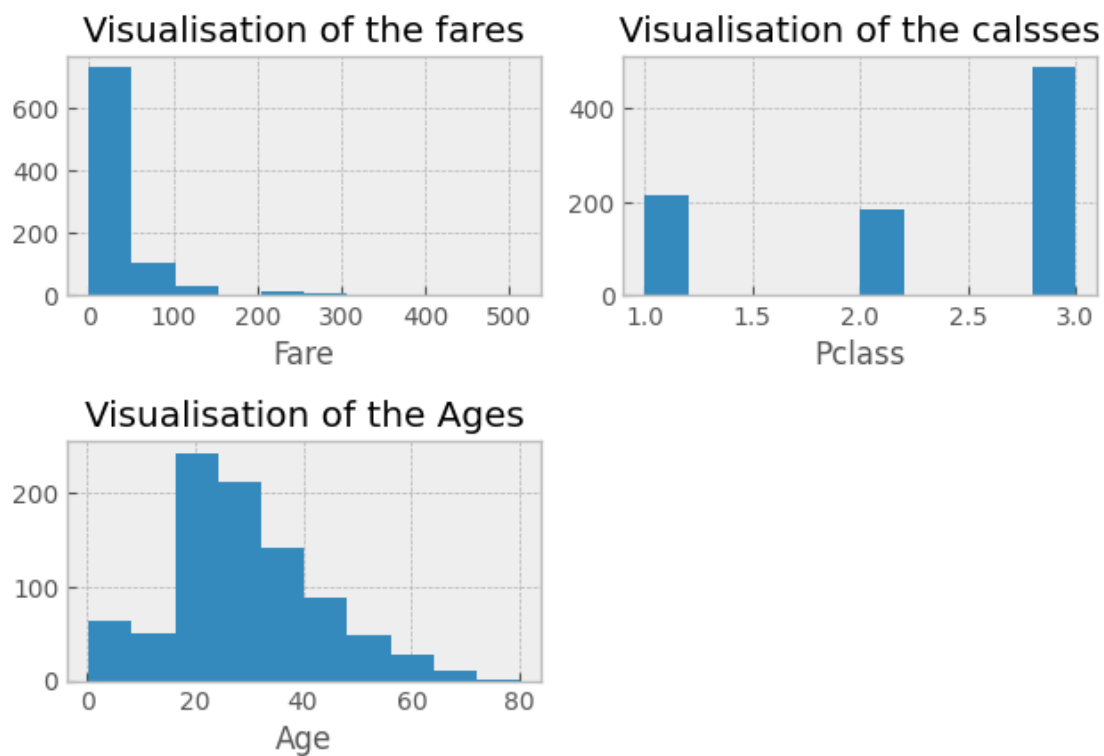
ax2.set_title('Visualisation of the calsses')

ax3 = fig.add_subplot(2, 2, 3) # we have changed (2,1,2) to (1,2,2)
ax3.hist(df["Age"])
ax3.set_xlabel('Age')
ax3.set_title('Visualisation of the Ages')

fig.tight_layout() # comment out this line to see the difference
fig.subplots_adjust(top=0.85)
plt.show()

```

## Histograms for Fare, Class, and Age



**Assignment 1)** Now we want to compare the fare and class, as we did before, but this time we want to divide them into two colors, depending on if they survived or not.

```

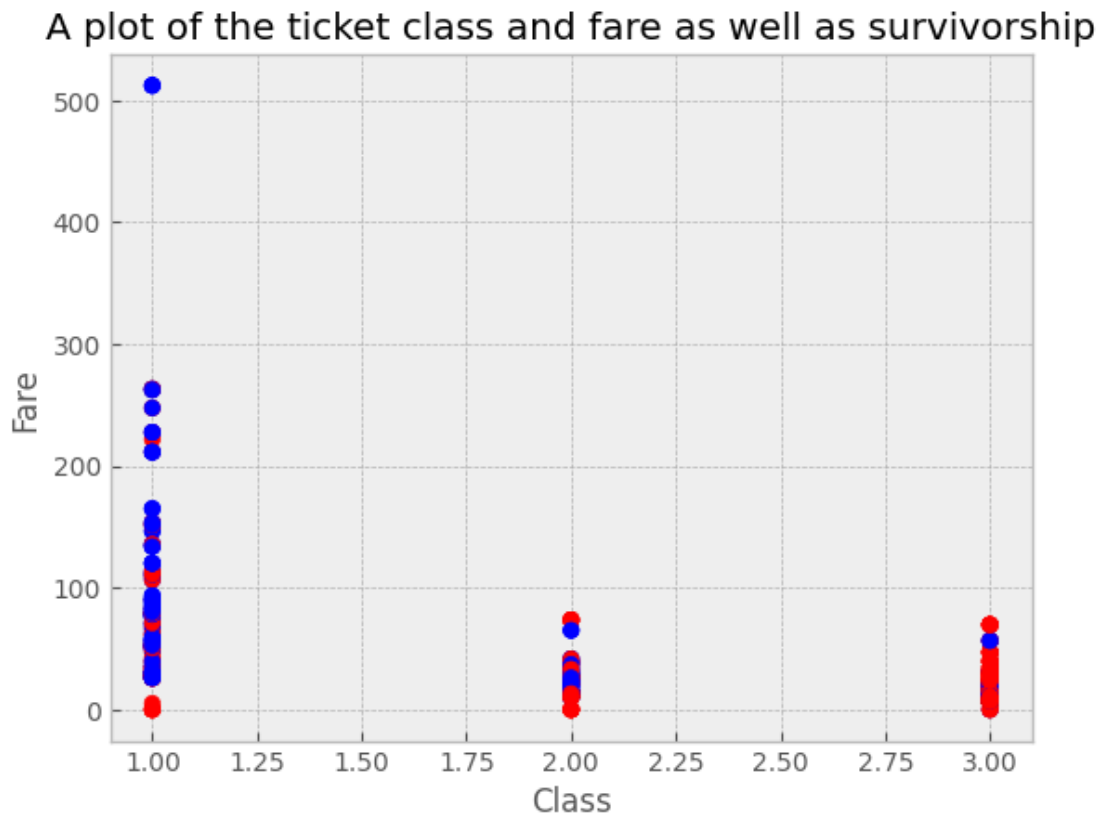
[ ]: # ASSIGNMENT:
# Make a scatter plot with fare on the y-axis
# and class on the x-axis
# using red dots for all the people who died
# and blue dots for the people who survived.
# use different markers for the survived and died points

```

```
# label the plot and the axes appropriately

c = df["Survived"].replace(1, value="b").replace(0, value="r")

plt.scatter(df["Pclass"], df["Fare"], color=c)
plt.xlabel("Class")
plt.ylabel("Fare")
plt.title("A plot of the ticket class and fare as well as survivorship")
plt.show()
```

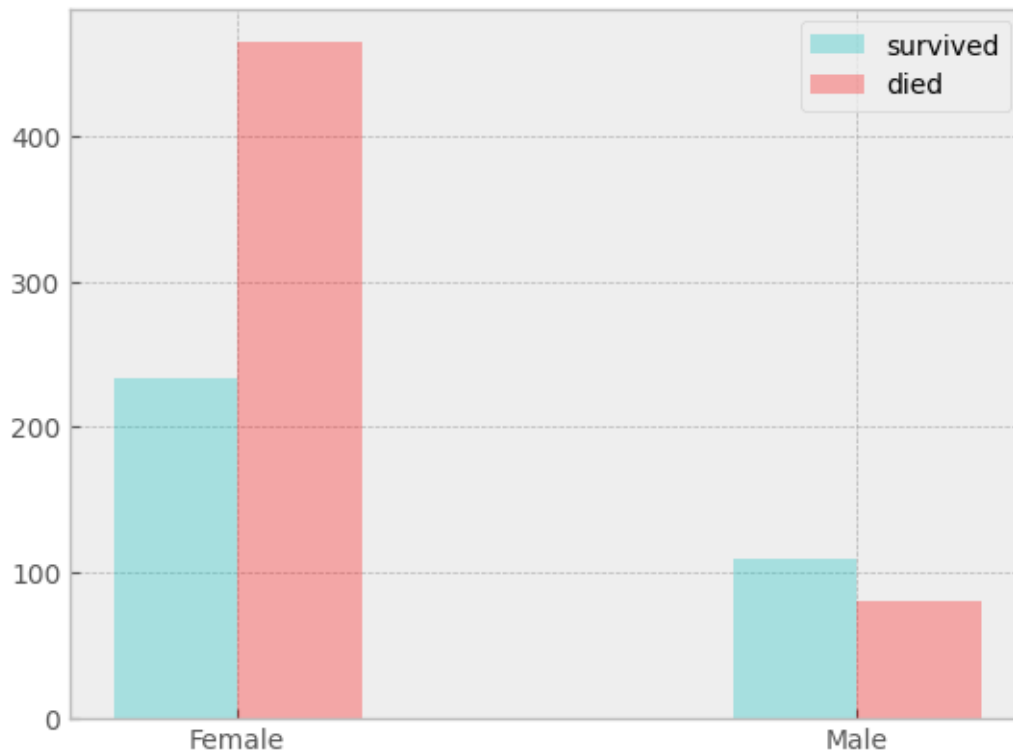


**Assignment m)** It might also be interesting to visualize how many of the men and women survived. This can be done with the bar function, which will be given to you.

```
[ ]: # ASSIGNMENT:
# Calculate how many women and men died and survived.
# label the plot and the axes appropriately

female_died, male_died = df[df["Survived"] == 0]["Sex"].value_counts()
```

```
plt.bar([0.9,1.9], [female_survived, male_survived] , color='c',
        label='survived', width=0.2, alpha=0.3)
plt.bar([1.1, 2.1], [female_died, male_died] , color='r', label='died', width=0.
        2, alpha=0.3)
plt.xticks([1,2], ['Female','Male'])
plt.legend()
plt.show()
```



[ ]: *### (Optional) Plotting a histogram of a random distribution*

OPTIONAL:

Plotting a Histogram of Random values

Your task is to generate 10000 random numbers that follows the normal distribution, with a mean,  $\mu = 1$ , and variance  $\sigma^2 = 0.25$ .

Plot the **normalized** histogram with 50 bars and a contour plot.

```
[ ]: import numpy as np
import matplotlib.pyplot as plt

plt.style.use('ggplot')
np.random.seed(42)
```

```
# OPTIONAL ASSIGNMENT:
# Draw 10000 random values from a normal distribution with:
#   mu = 1, sigma2 = 0.25
#
# Plot the histogram and cumulative distribution
# label the plot and the axes appropriately

# YOUR CODE HERE

plt.show()
```