

Practical Notebook 2

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.

```
In [3]: # 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}")
```

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.

```
In [4]: 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.

```
In [5]: 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.

```
In [6]: 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.

```
In [7]: df.to_csv('test_people.csv', index=False, header=True)
```

Read a dataset from file

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

```
In [8]: 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`.

```
In [9]: print(df.describe())
```

	Age
count	3.000000
mean	27.333333
std	15.50269
min	12.000000
25%	19.500000
50%	27.000000
75%	35.000000
max	43.000000

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

```
In [10]: 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.

```
In [11]: 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.

```
In [12]: # 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

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".

```
In [13]: # Downloading the titanic dataset
import wget
wget.download('https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/stuff/tit
```

Out[13]: 'titanic (1).csv'

Assignment a)

```
In [14]: # 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.describe())
```

	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)

```
In [15]: # ASSIGNMENT:
# Count the number of males and females
```

```
male, female = df['Sex'].value_counts()
print(male, female)
```

573 314

Assignment c)

```
In [16]: # 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)

```
In [17]: # ASSIGNMENT:  
# Count how many survived (1) and how many died (0)  
  
# YOUR CODE HERE  
  
died, survived = df['Survived'].value_counts()  
print(died, survived)
```

```
545 342
```

Assignment e)

```
In [18]: # 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(female_survived, male_survived)
```

```
233 109
```

Assignment f)

```
In [19]: # 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.drop('Survived', axis=1)  
y = df['Survived']  
  
x_describe = X.describe()  
y_describe = y.describe()  
  
print(x_describe, y_describe)
```

	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

Name: Survived, dtype: float64

Assignment g)

```
In [20]: # 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 1

# YOUR CODE HERE

X_new = (X-X.mean())/X.std()
y_new = (y-y.mean())/y.std()

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

print(X_new_describe, y_new_describe)
```

	Age	Fare	Parents/Children Aboard	Pclass	\
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	-2.05719	-0.64894	-0.47471	-1.56040	
25%	-0.65299	-0.48974	-0.47471	-0.36517	
50%	-0.10420	-0.35859	-0.47471	0.83006	
75%	0.60392	-0.02346	-0.47471	0.83006	
max	3.57803	9.64251	6.95594	0.83006	

	Siblings/Spouses Aboard	
count	887.00000	
mean	-0.00000	
std	1.00000	
min	-0.47559	
25%	-0.47559	
50%	-0.47559	
75%	0.42966	
max	6.76640	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

C:\Users\albin\AppData\Local\Temp\ipykernel_928\3779499128.py:7: FutureWarning: The default value of numeric_only in DataFrame.mean is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.

```
X_new = (X-X.mean())/X.std()
```

C:\Users\albin\AppData\Local\Temp\ipykernel_928\3779499128.py:7: FutureWarning: The default value of numeric_only in DataFrame.std is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.

```
X_new = (X-X.mean())/X.std()
```

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](#).

Examples

```
In [21]: # import the relevant libraries
import matplotlib.pyplot as plt
```

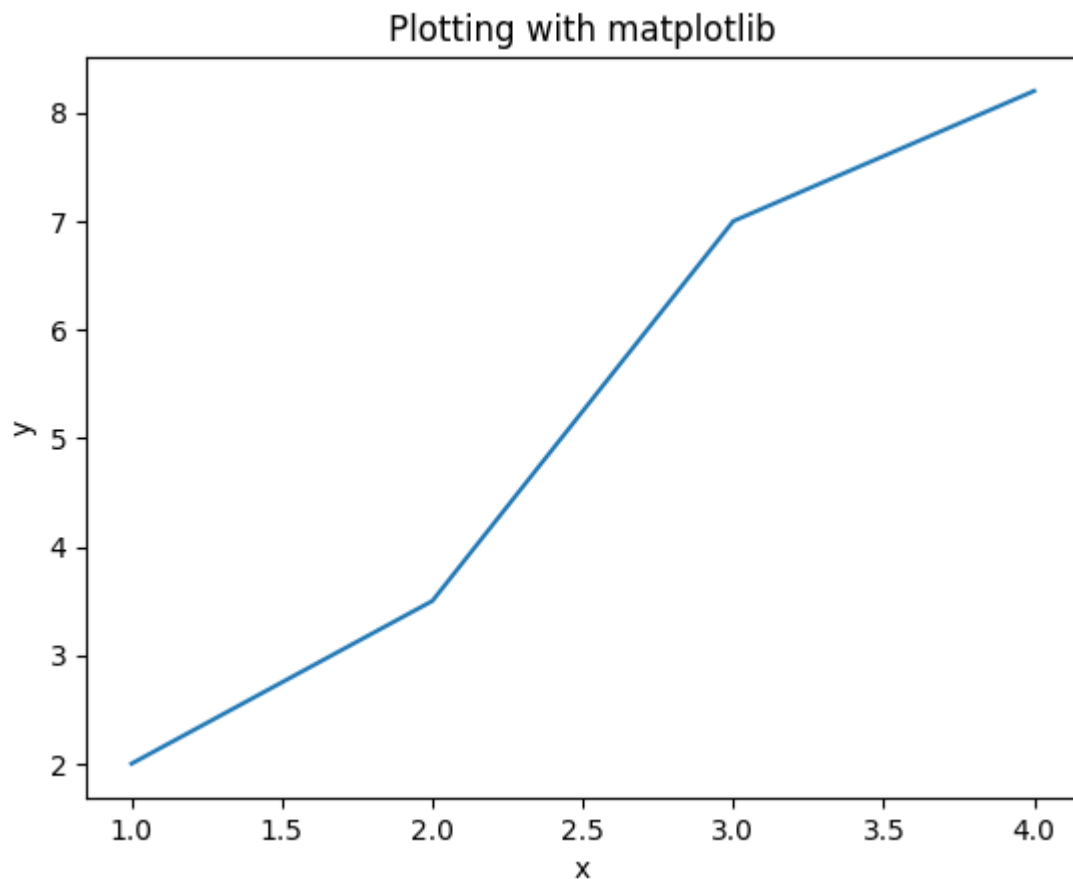
```
import numpy as np
```

We will start by looking at some small lists.

```
In [22]: # 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](#).

```
In [23]: # 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)

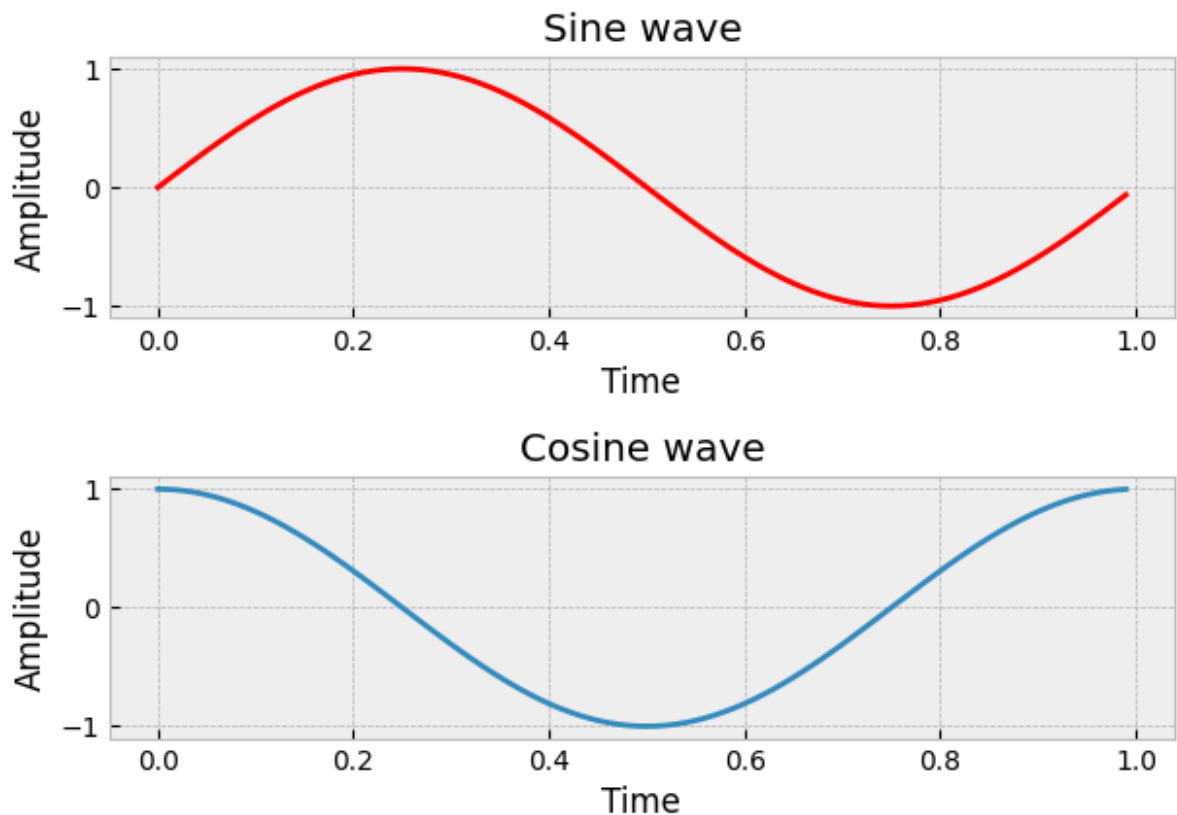
ax1 = fig.add_subplot(2,1,1)
ax1.plot(t, sin, color='red', lw=2)
ax1.set_ylabel('Amplitude')
ax1.set_xlabel('Time')
ax1.set_title('Sine wave')

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



```

In [24]: # 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()

```

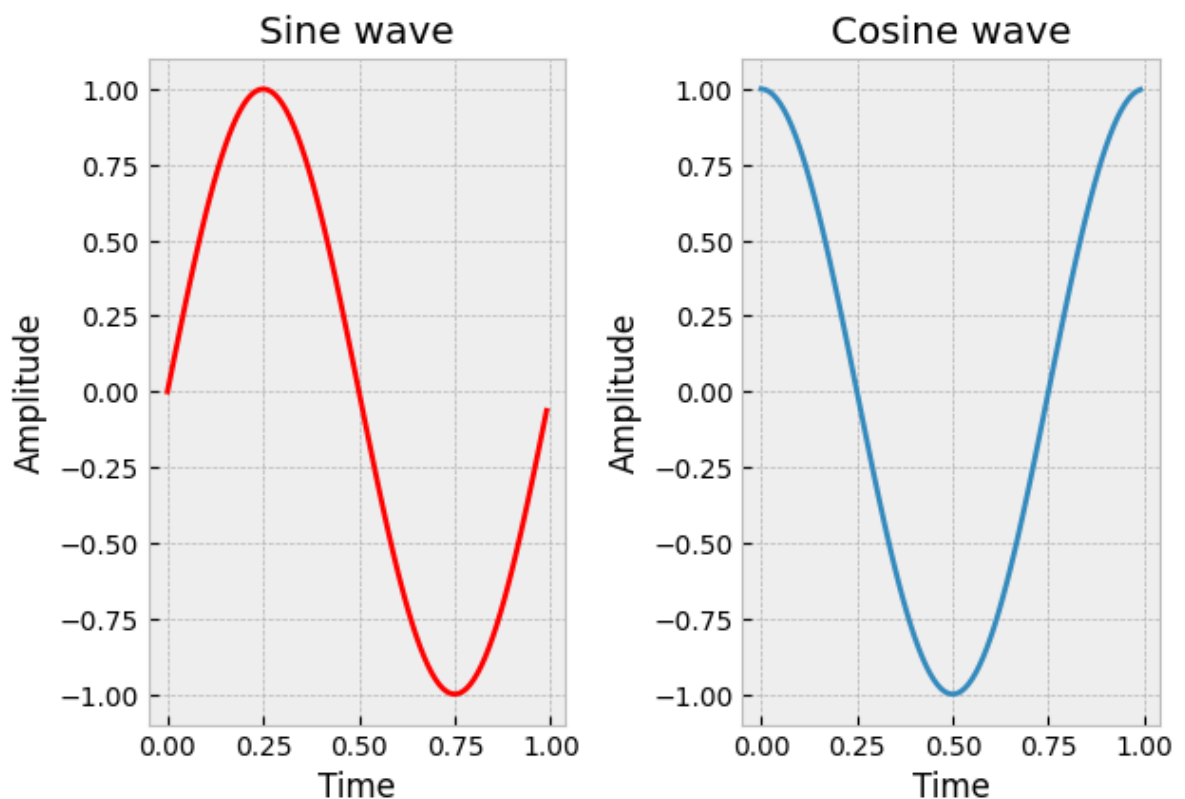
```
fig.suptitle("Sine and cosine for different t", fontsize=18)

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()
```

Sine and cosine for different t



And with different stylings

```
In [25]: # Here are all the different "pre-configured" styles matplotlib supports
# https://matplotlib.org/tutorials/intermediate/artists.html#sphx-glr-tutorials-int
plt.style.available
```

```
Out[25]: ['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.

```
In [26]: # Matrix subplot

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

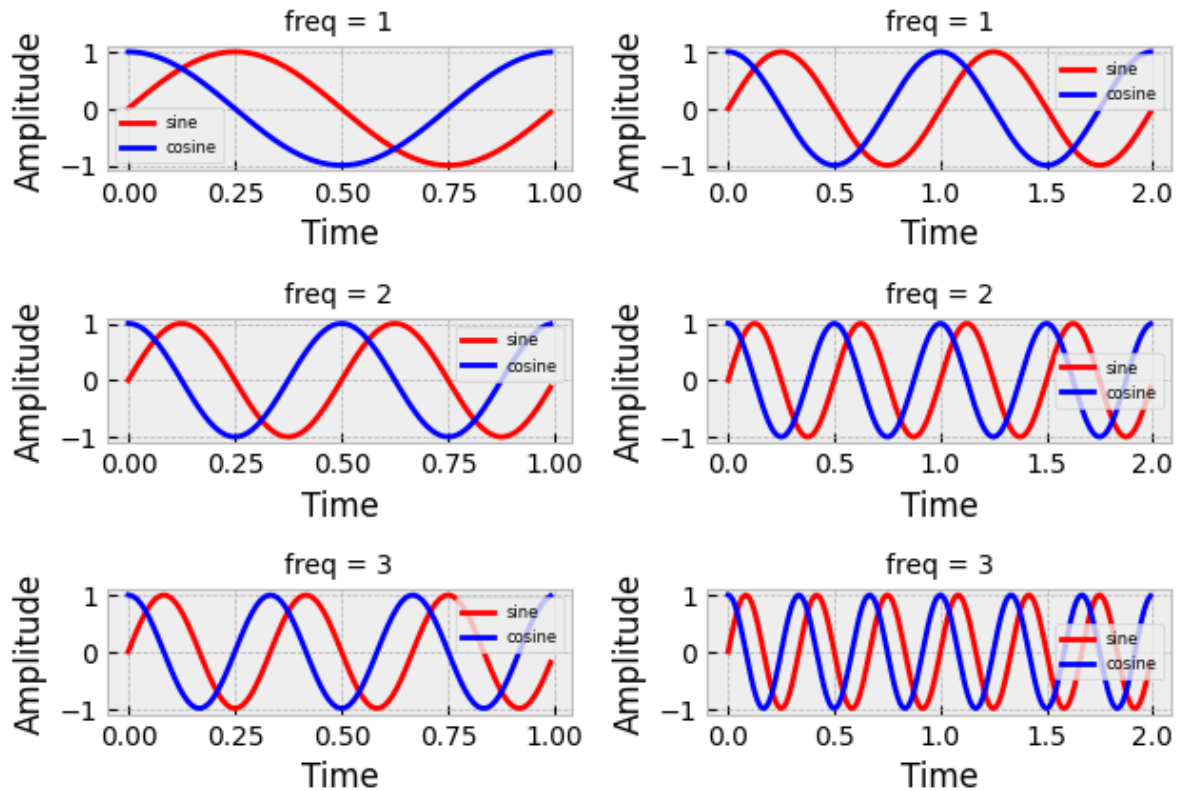
i = 1
for freq in [1, 2, 3]:
    for t_max in [1, 2]:
        t = np.arange(0.0, t_max, 0.01)
        sin = np.sin(2*freq*np.pi*t)
        cos = np.cos(2*freq*np.pi*t)

        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



Plotting data from Pandas

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

```
In [27]: # Downloading the titanic dataset
import wget
wget.download('https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/stuff/tit
```

```
Out[27]: 'titanic (2).csv'
```

```
In [28]: # Load the titanic dataset for plotting
import pandas as pd
df = pd.read_csv('titanic.csv')
```

Assignment h)

```
In [29]: # 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

# YOUR CODE HERE
ticket_class = df['Pclass']
```

```

fare = df['Fare']

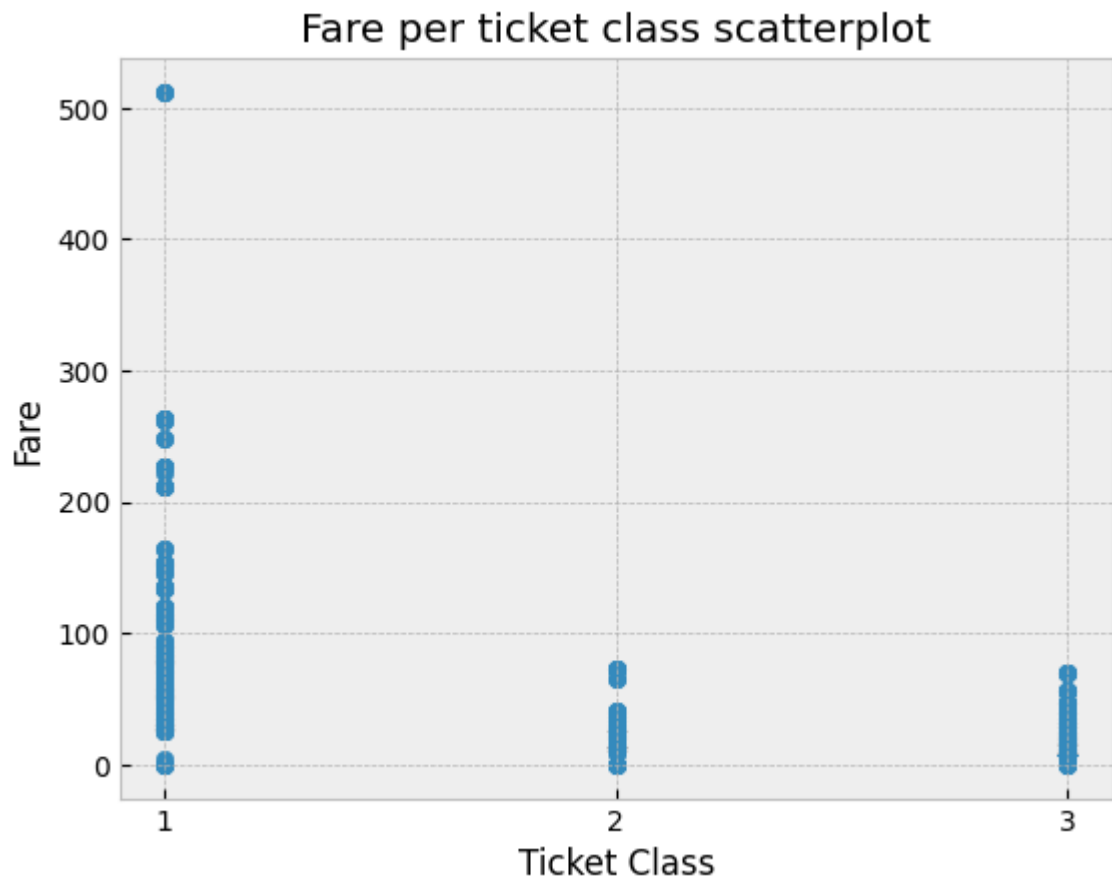
plt.title('Fare per ticket class scatterplot')

plt.xlabel('Ticket Class')
plt.xticks([1,2,3])

plt.ylabel('Fare')

plt.scatter(ticket_class, 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.

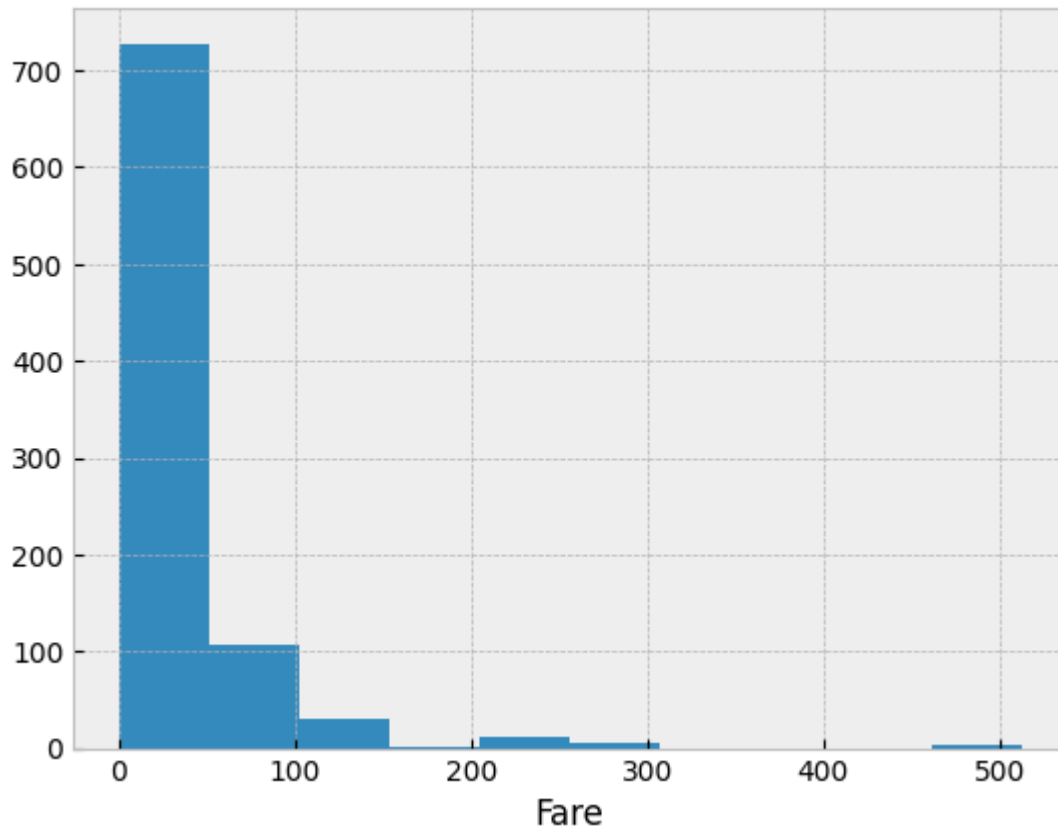
```

In [30]: fare = df["Fare"]

plt.hist(fare)
plt.xlabel("Fare")
plt.title("Visualizaion of the fare difference")
plt.show()

```

Visualizaion of the fare difference

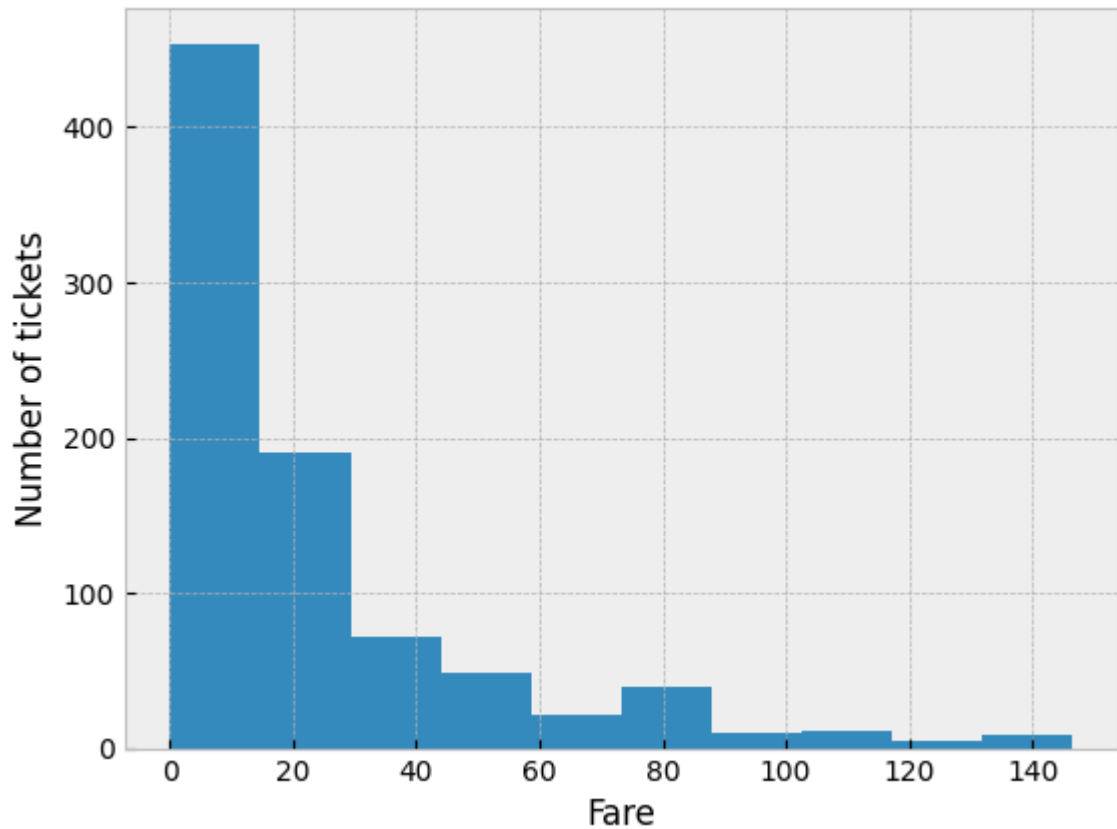


As you can see, most of the people paid less than 150 for the ticket.

```
In [31]: # ASSIGNMENT:
# Plot a histogram over the people who paid less than, or equal to, 150.
# Label the plot and the axes appropriately

# YOUR CODE HERE
fare_limit = 150
low_fare = df[df['Fare'] <= fare_limit]['Fare']
plt.title('Low fare difference visualization')
plt.xlabel('Fare')
plt.ylabel('Number of tickets')
plt.hist(low_fare)
plt.show()
```

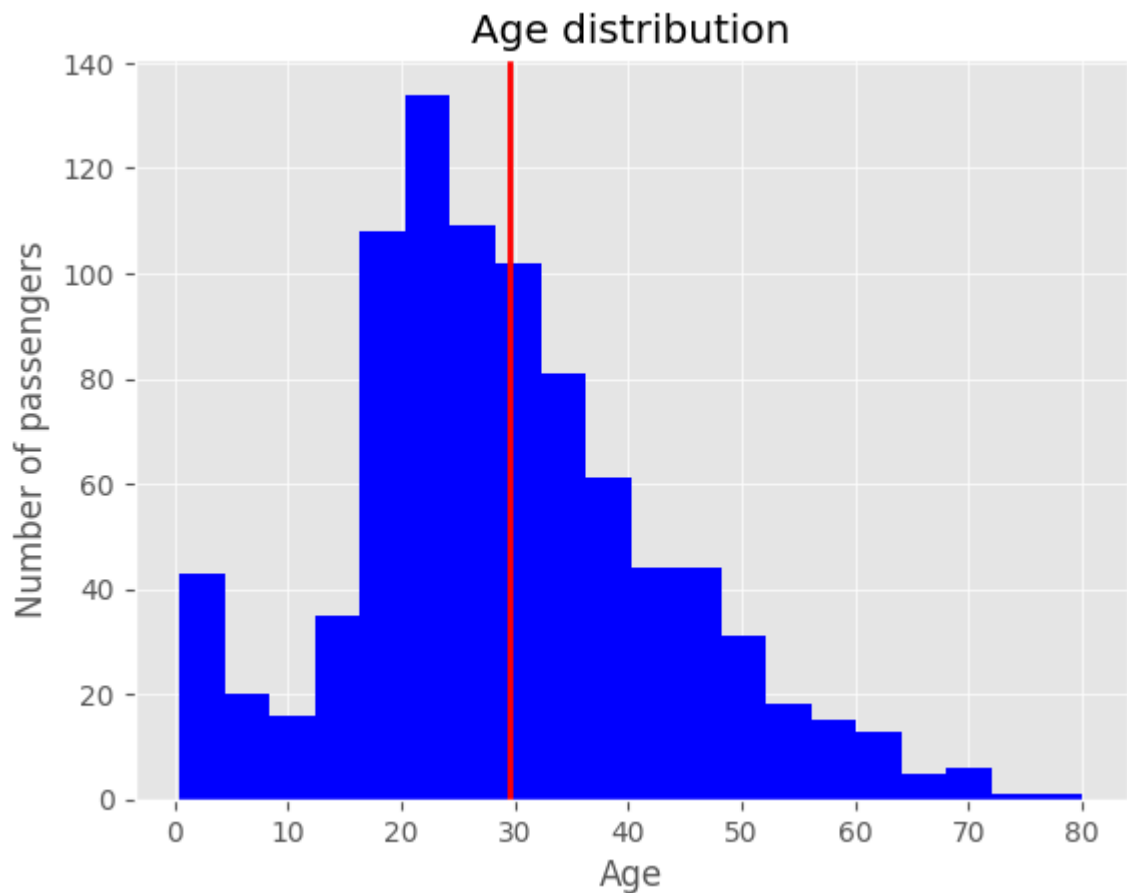
Low fare difference visualization



Assignment j)

```
In [44]: # ASSIGNMENT:
# plot a histogram over all the ages with 20 bins. Draw a vertical line at the mean
# label the plot and the axes appropriately

# YOUR CODE HERE
age = df['Age']
plt.title('Age distribution')
plt.xlabel('Age')
plt.ylabel('Number of passengers')
plt.hist(age, bins=20, color='b')
plt.axvline(age.mean(), color='r')
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.

```
In [33]: # ASSIGNMENT:
# Make a subplot over the Fare, Class, and Age
# Label the plot and the axes appropriately

# YOUR CODE HERE
fare = df['Fare']
ticket_class = df['Pclass']
age = df['Age']

fig = plt.figure()
fig.suptitle('Visualisations of passenger fare, class and age', fontsize=18)
fig.set_figheight(9)

ax = fig.add_subplot(3,1,1)
ax.set_title('Fare histogram')
ax.set_xlabel('Fare')
ax.hist(fare)

ax = fig.add_subplot(3,1,2)
ax.set_title('Fares for each ticket class')
ax.set_xlabel('Ticket class')
ax.set_ylabel('Fare')
```



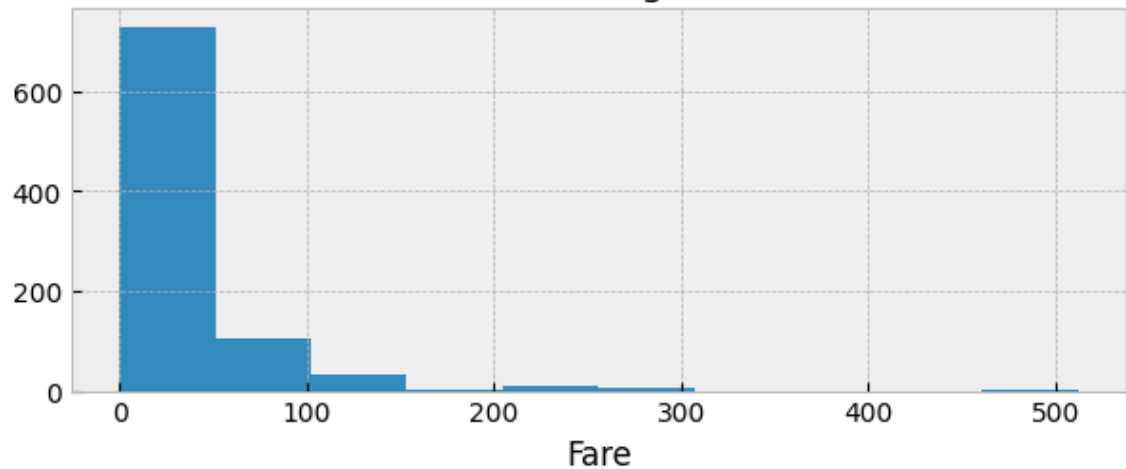
```
ax.set_xticks([1,2,3])
ax.scatter(ticket_class, fare)

ax = fig.add_subplot(3,1,3)
ax.set_title('Age histogram')
ax.set_xlabel('Age')
ax.set_ylim(0,250)
ax.hist(age)

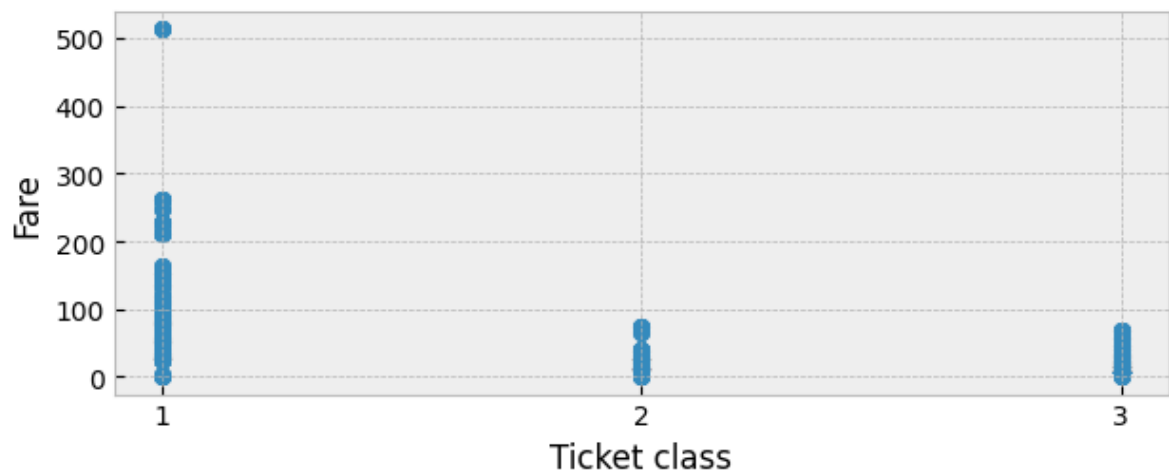
fig.subplots_adjust(top=0.85)
fig.tight_layout() # comment out this line to see the difference
plt.show()
```

Visualisations of passenger fare, class and age

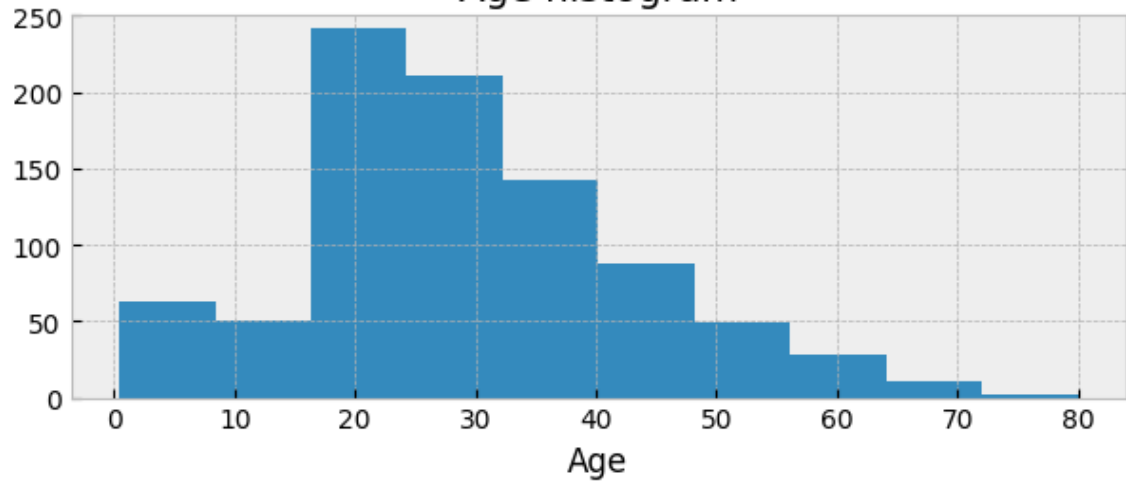
Fare histogram



Fares for each ticket class



Age histogram



Assignment I)

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.

```

In [34]: # 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

# YOUR CODE HERE
died = df[df['Survived']==0]
survived = df[df['Survived']==1]

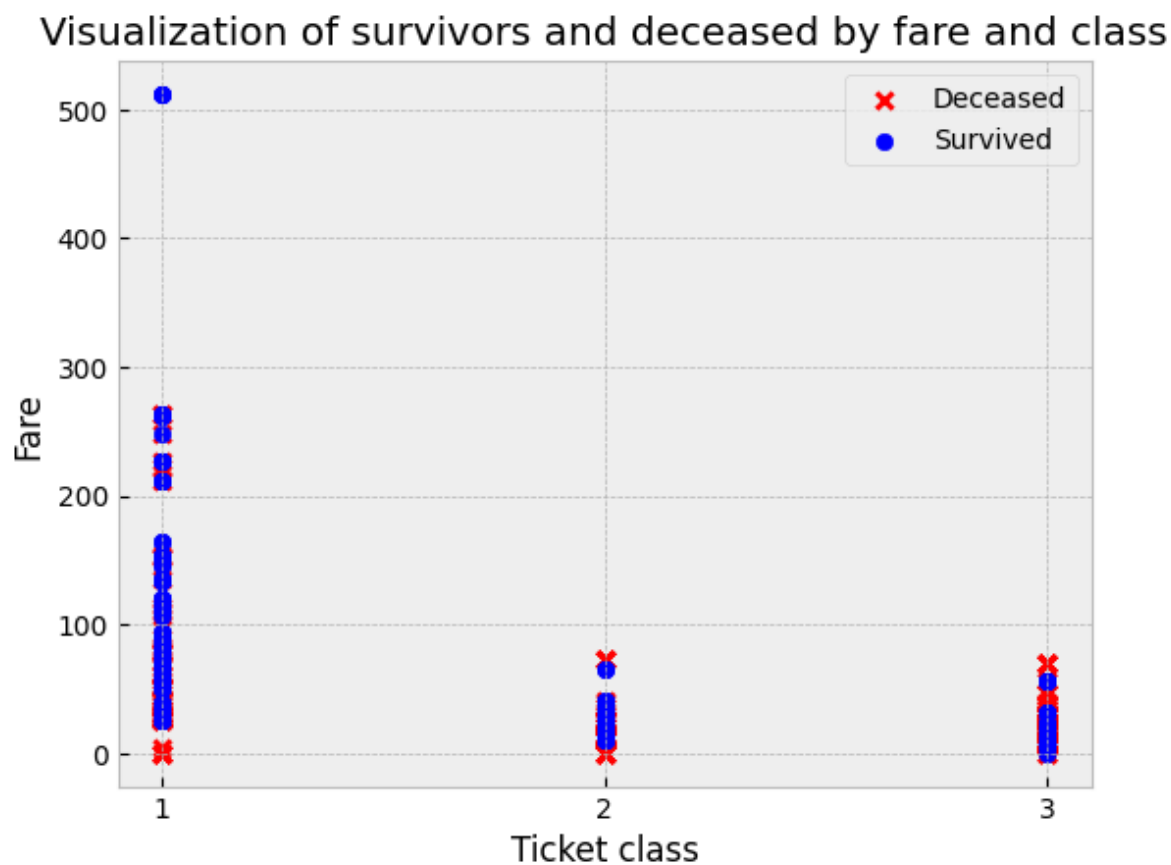
d_class = died['Pclass']
d_fare = died['Fare']

s_class = survived['Pclass']
s_fare = survived['Fare']

plt.figure()
plt.title('Visualization of survivors and deceased by fare and class')
plt.xlabel('Ticket class')
plt.xticks([1,2,3])
plt.ylabel('Fare')

plt.scatter(d_class, d_fare, color='r', marker='x')
plt.scatter(s_class, s_fare, color='b', marker='o')
plt.legend(['Deceased', 'Survived'])
plt.show()

```



Assignment m)

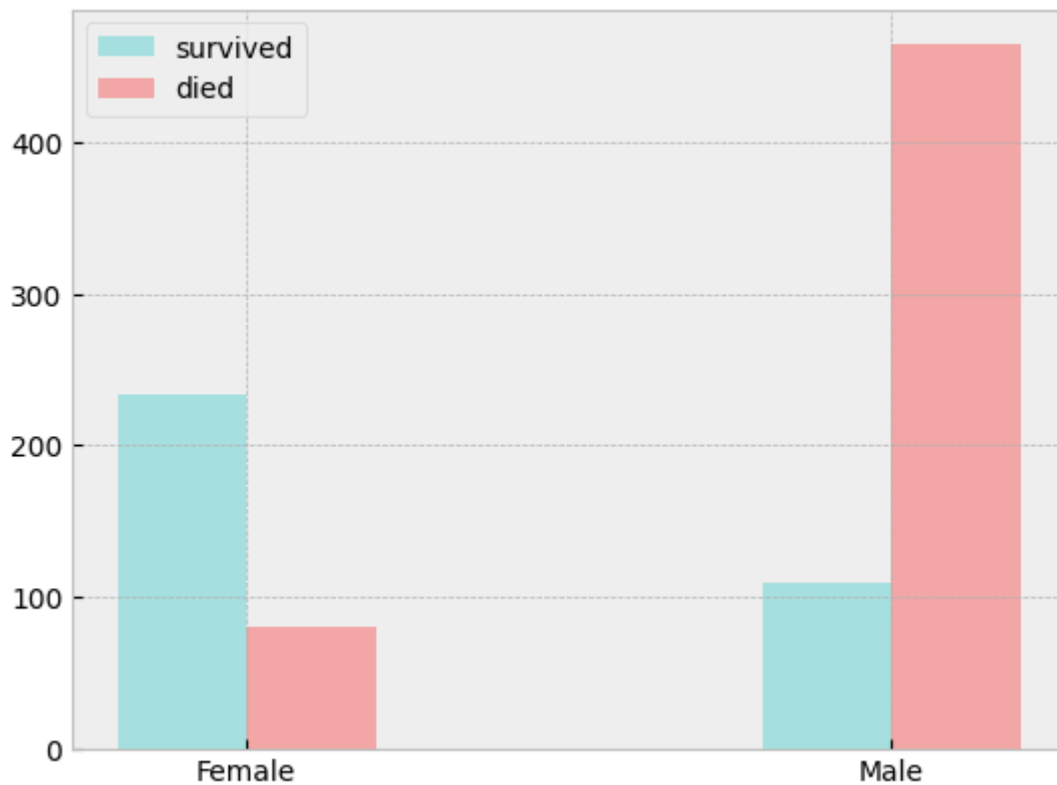
It might also be interesting to visualize how many of the men and women survived. This can be done with the bar function, which will be given to you.

```
In [35]: # ASSIGNMENT:
# Calculate how many women and men died and survived.
# Label the plot and the axes appropriately

# YOUR CODE HERE
female_survived, male_survived = survived['Sex'].value_counts()
female_died, male_died = died['Sex'].value_counts()

assert female_survived + female_died == len(df[df['Sex']=='female'].index)
assert male_survived + male_died == len(df[df['Sex']=='male'].index)

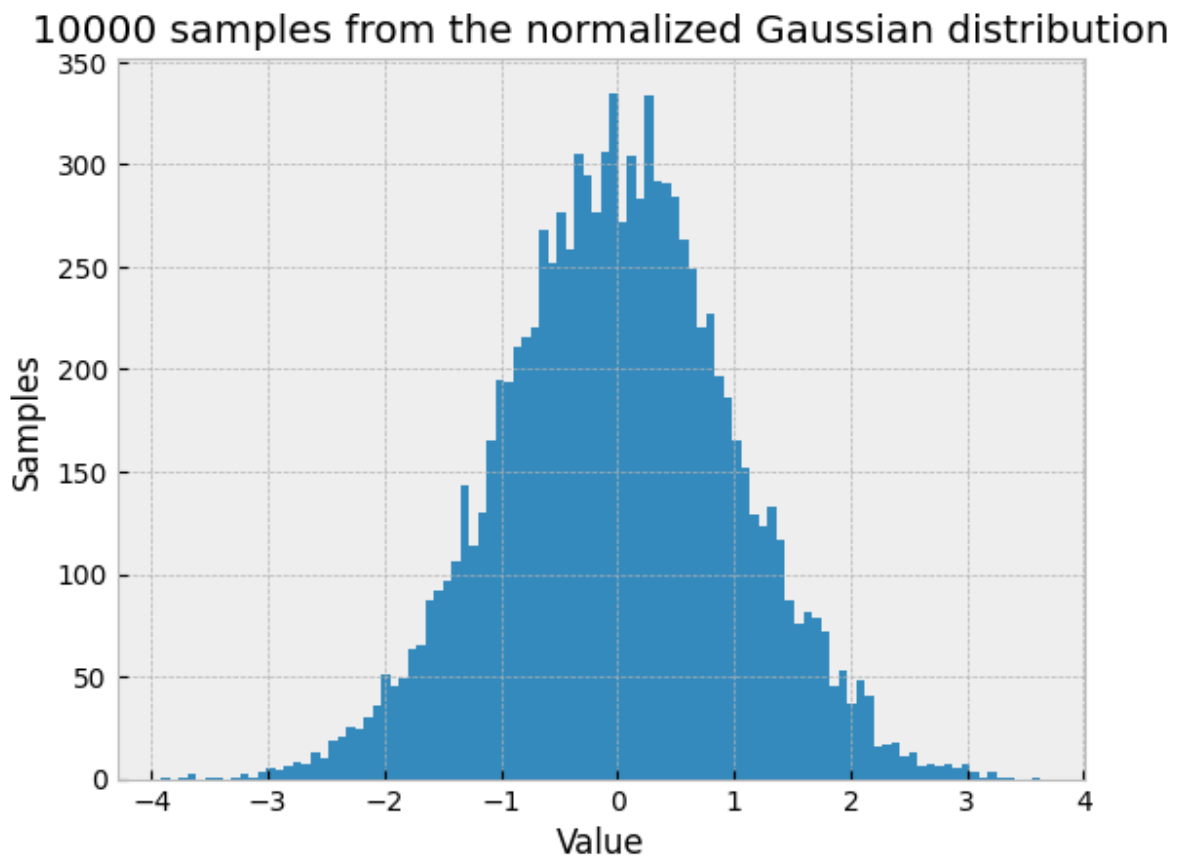
plt.bar([0.9,1.9], [female_survived, male_survived] , color='c', label='survived',
plt.bar([1.1, 2.1], [female_died, male_died] , color='r', label='died', width=0.2,
plt.xticks([1,2], ['Female', 'Male'])
plt.legend()
plt.show()
```



```
In [36]: ### (Optional) Plotting a histogram of a random distribution
import numpy as np
import matplotlib.pyplot as plt

samples = 10000
n_vec = np.random.normal(size=samples)
plt.xlabel('Value')
plt.ylabel('Samples')
```

```
plt.title(f'{samples} samples from the normalized Gaussian distribution')
plt.hist(n_vec, bins=100)
plt.show()
```



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.

```
In [1]: 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

no_samples = 10000
```

```

mu=1
sigma2=0.25
samples = np.random.normal(mu, sigma2, size=no_samples)
normalized = (samples - mu) / sigma2
fig = plt.figure()
ax = fig.add_subplot(2,1,1)
ax.set_xlabel('Value')
ax.set_ylabel('Samples')
ax.set_title(f'{no_samples} samples from the normalized Gaussian distribution')
ax.hist(normalized, bins=50)

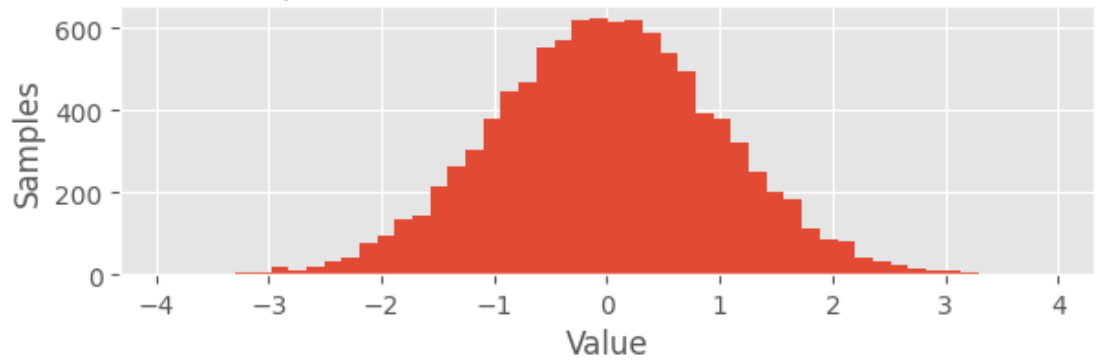
ax = fig.add_subplot(2,1,2)
ax.set_title('Cumulative distribution')
ax.set_xlabel('Value')
ax.set_ylabel('Cumulative probability')
ax.hist(normalized, cumulative=True, bins=50)

# TODO...

fig.tight_layout()
plt.show()

```

10000 samples from the normalized Gaussian distribution



Cumulative distribution

