

Week 8 Data visualisation¶

Name: Cheong Win Yan

Matric number: 23005189/1

- Data visualization is a graphical representation of your information or data.
- It is similar to Excel but in Python.

Importance of Data Visualization

- Making complex data understandable
- Enhancing decision making
- Storytelling with Data

In Python, data visualization is mostly performed using 3 libraries:

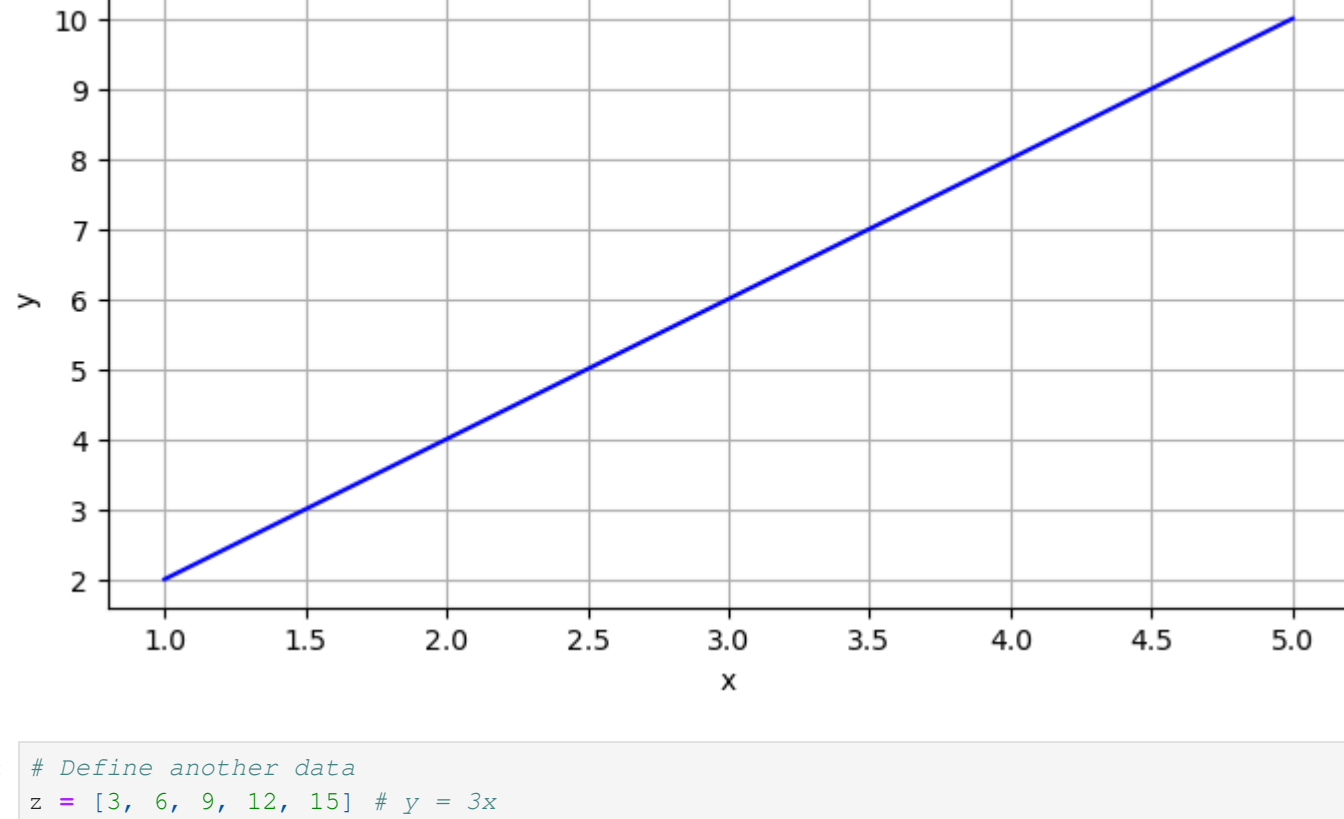
1. matplotlib
2. seaborn
3. plotly

1. Line Plot

```
In [28]: # Import required libraries
import matplotlib.pyplot as plt
```

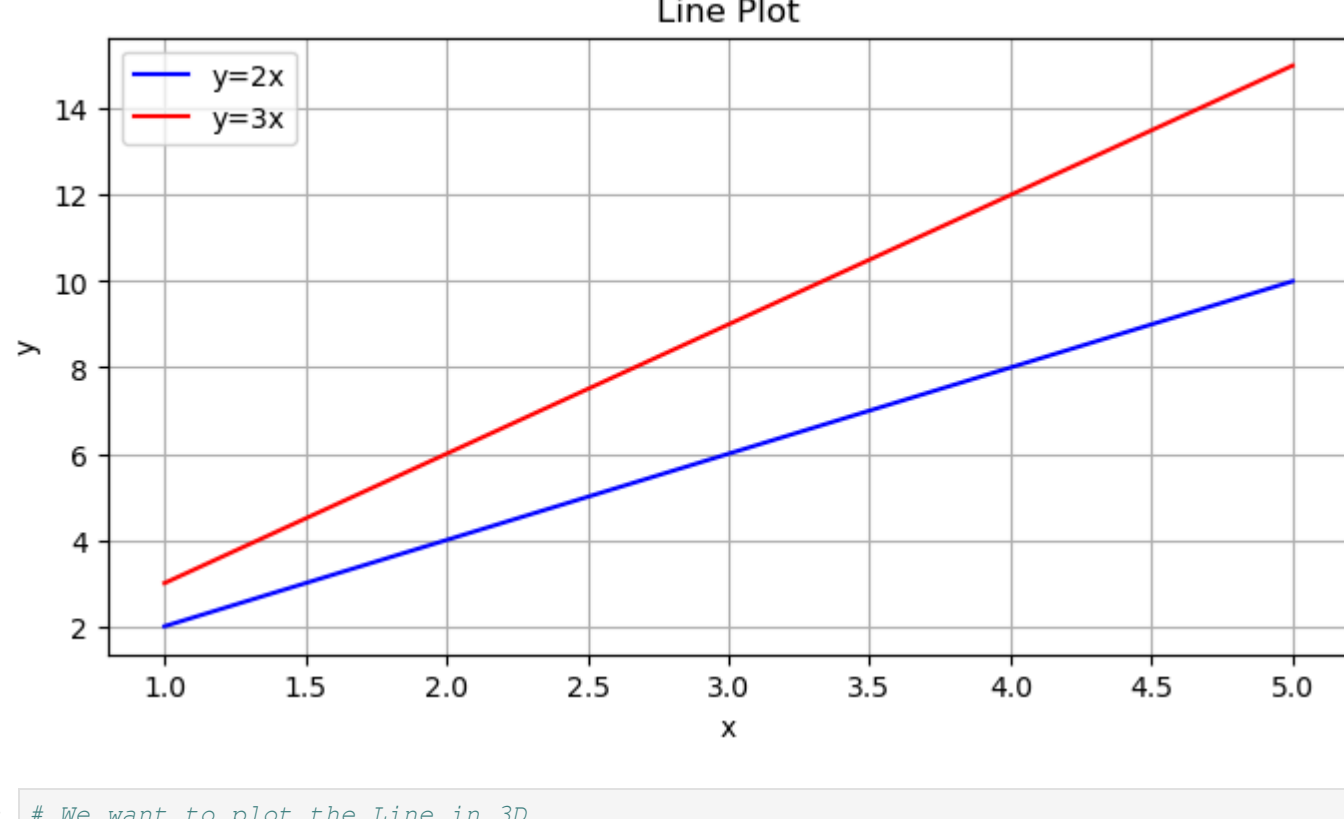
```
In [29]: # Define some data to plot
x = [1, 2, 3, 4, 5]
y = [2, 4, 6, 8, 10] # y = 2x
```

```
In [30]: # Plot x and y
plt.figure(figsize=(8, 4))
plt.plot(x, y, color='blue')
plt.title('Graph of x vs y')
plt.xlabel('x')
plt.ylabel('y')
plt.grid(True)
plt.show()
```



```
In [31]: # Define another data
z = [3, 6, 9, 12, 15] # y = 3x
```

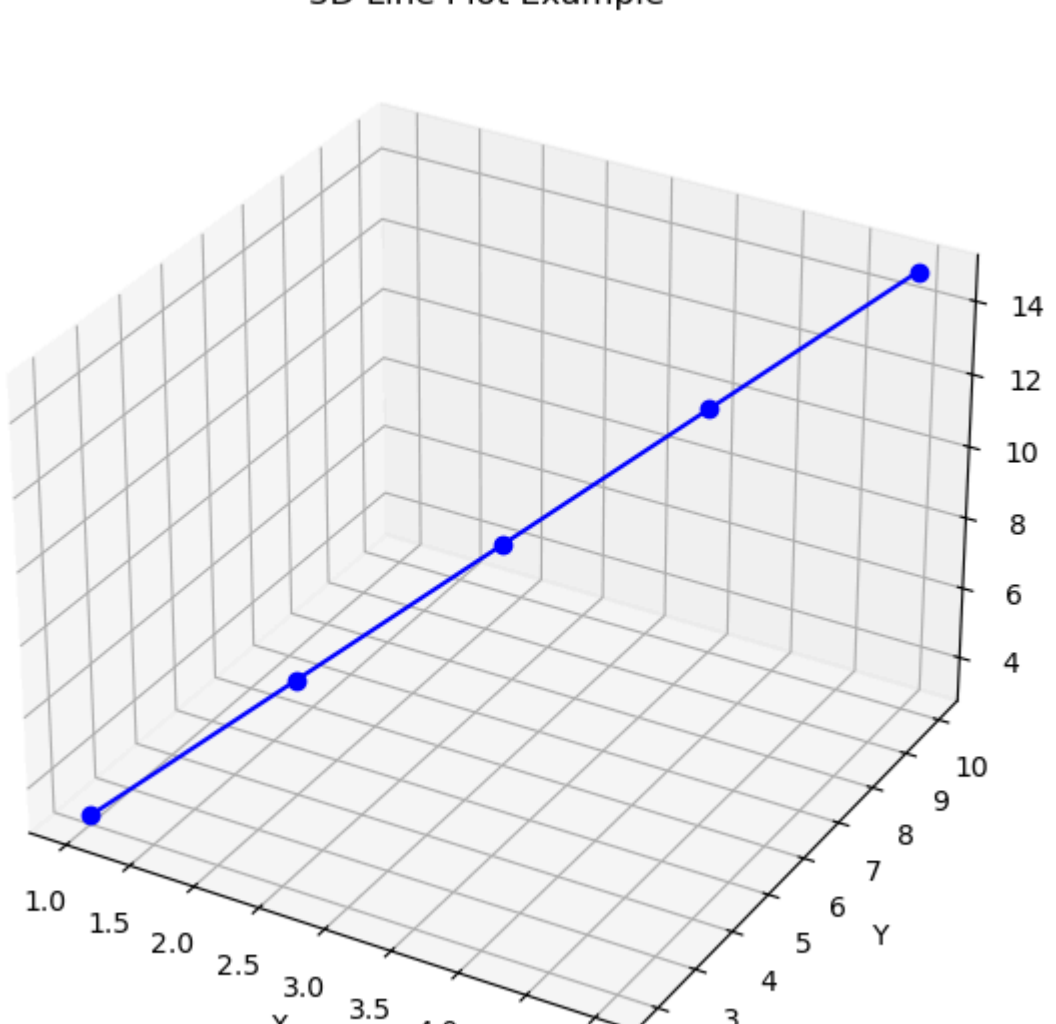
```
In [32]: # Plot two lines on the same graph
plt.figure(figsize=(8, 4))
plt.plot(x, y, color='blue', label='y=2x')
plt.plot(z, z, color='red', label='y=3x')
plt.title('Line Plot')
plt.xlabel('x')
plt.ylabel('y')
plt.grid(True)
plt.legend()
plt.show()
```



```
In [33]: # We want to plot the line in 3D
from mpl_toolkits.mplot3d import Axes3D
```

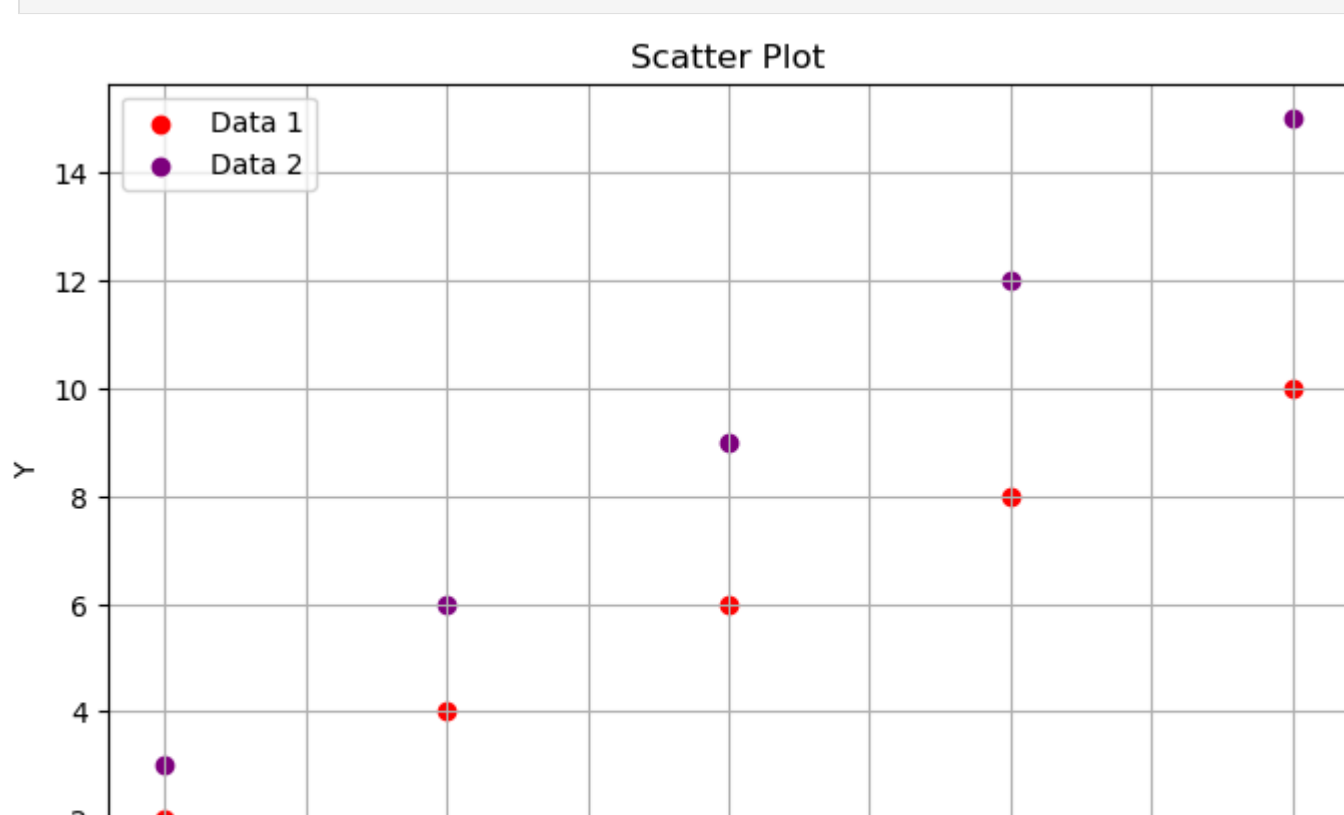
```
In [34]: fig = plt.figure(figsize=(10, 7))
ax = fig.add_subplot(111, projection='3d')
ax.plot(x, y, z, color='blue', marker='o')
ax.set_title('3D Line Plot Example')
ax.set_xlabel('x')
ax.set_ylabel('y')
ax.set_zlabel('z')
plt.show()
```

3D Line Plot Example



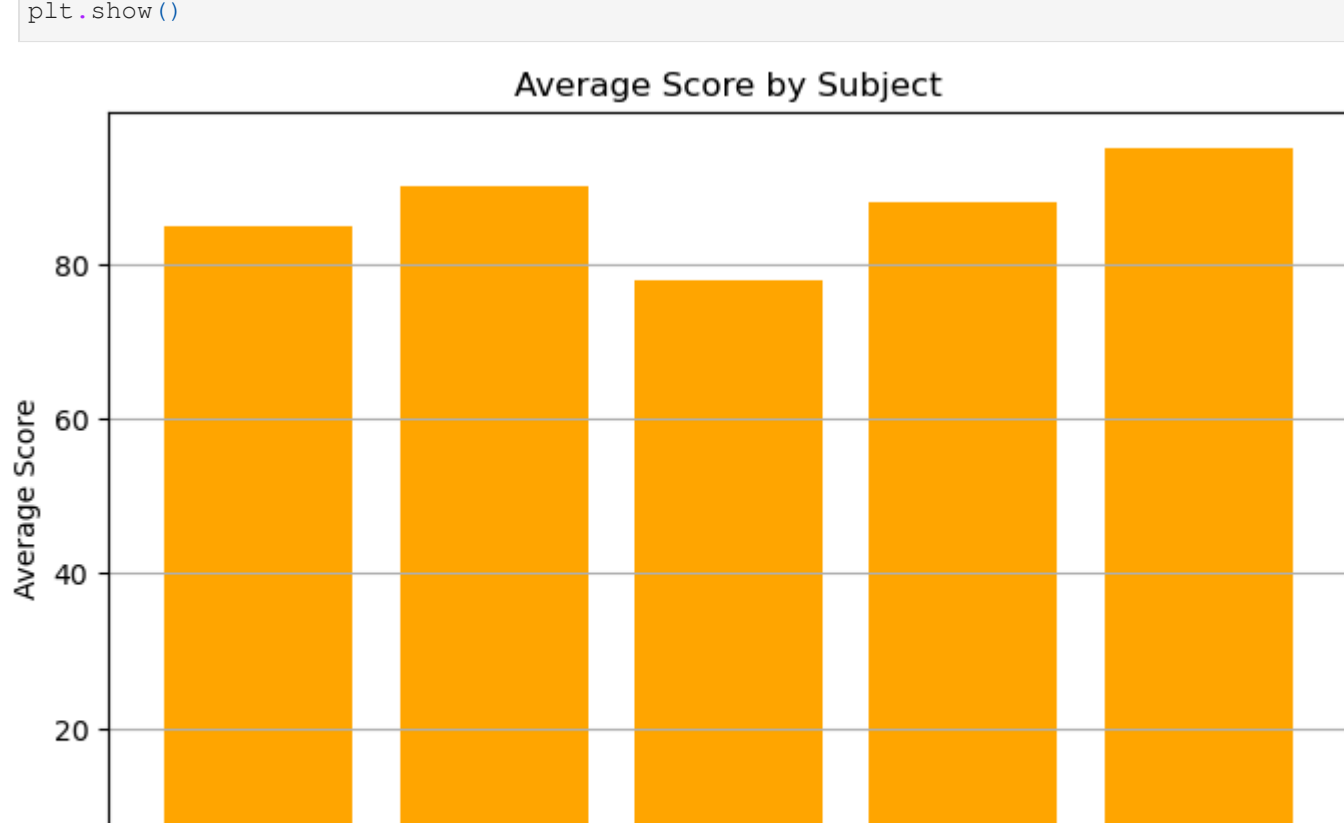
2. Scatter Plot

```
In [37]: plt.figure(figsize=(8, 5))
plt.scatter(x, y, label='Data 1', color='red')
plt.scatter(z, z, label='Data 2', color='purple')
plt.title('Scatter Plot')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.grid(True)
plt.show()
```



3. Bar Plot

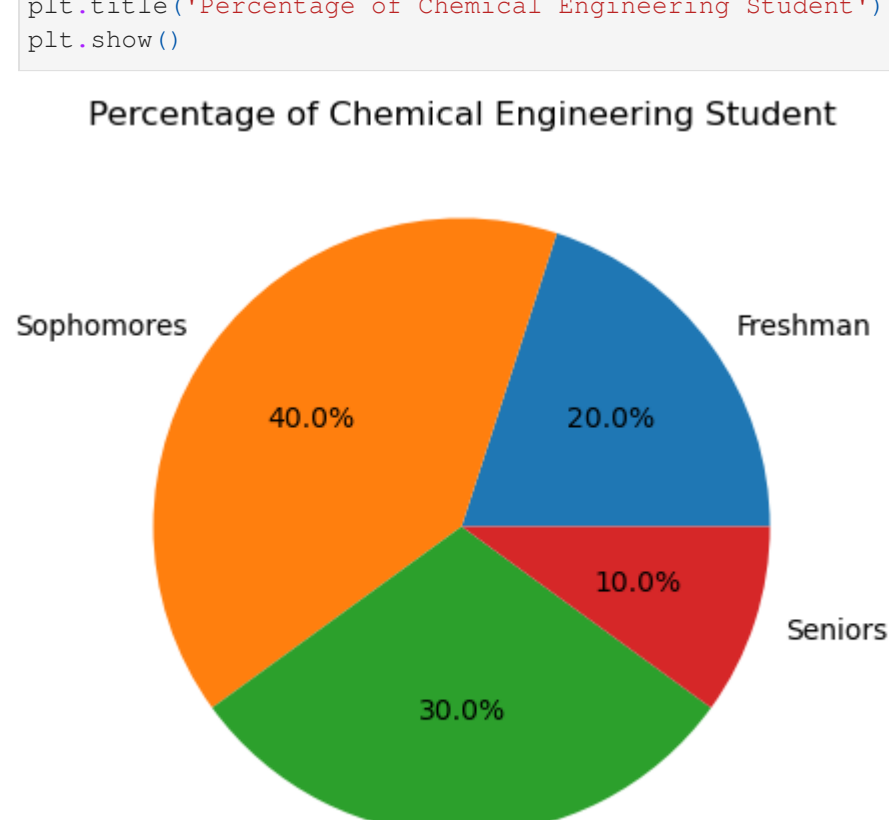
```
In [38]: #Generate some data
subjects = ('Math', 'Science', 'English', 'History', 'Art')
average_score = [85, 90, 78, 88, 92]
plt.figure(figsize=(8, 5))
plt.bar(subjects, average_score, color='orange')
plt.title('Average Score by Subject')
plt.xlabel('Subjects')
plt.ylabel('Average Score')
plt.grid(axis='y')
plt.show()
```



4. Pie Plot

```
In [41]: student_level = ['Freshman', 'Sophomores', 'Juniors', 'Seniors']
num_students=[100, 200, 150, 50]
plt.figure(figsize=(8, 5))
plt.pie(num_students, labels=student_level, autopct='%1.1f%%')
plt.title('Percentage of Chemical Engineering Student')
plt.show()
```

Percentage of Chemical Engineering Student



Exercise: Data Visualization on Titanic Dataset

Objective: Using the Titanic dataset, explore the data and create visualizations to gain insights. At the end, please answer the following questions:

1. Did majority of the passengers survive?
2. Is there any relationship between fare and survival?

Instructions:

1. Load the Dataset

- Load the titanic.csv into your Jupyter Notebook environment
- Display first few rows using head()

2. Visualize Data

- Bar Plot: To show the number of passengers in each class (Pclass).
- Pie Chart: Show the proportion of passengers who survived vs. did not survive (Survived column).
- Line Plot: Plot the age distribution of passengers. Use Age on the x-axis and the count on the y-axis.
- Scatter Plot: Visualize the relationship between Fare and Age, using different colors for survivors and non-survivors.

3. Summarize the Results

- Answer the 2 questions stated in the Objective

```
In [60]: import pandas as pd
df = pd.DataFrame()
```

```
In [66]: # Load Dataset
titanic = pd.read_csv('titanic.csv')
```

```
In [67]: # Display the first 5 rows
titanic.head()
```

```
Out[67]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	7.2500
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	71.2833
2	3	1	3	Heikinen, Miss. Laina	female	26.0	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	53.1000
4	5	0	3	Allan, Mr. William Henry	male	35.0	8.0500

```
In [68]: titanic.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 7 columns):
 #   Column             Non-Null Count  Dtype
---  --
 0   PassengerId         891 non-null    int64
 1   Survived            891 non-null    int64
 2   Pclass              891 non-null    int64
 3   Name                891 non-null    object
 4   Sex                 891 non-null    object
 5   Age                 891 non-null    float64
 6   Fare                891 non-null    float64
dtypes: float64(2), int64(3), object(2)
memory usage: 48.9+ KB
```

```
In [69]: titanic.describe()
```

```
Out[69]:
```

	PassengerId	Survived	Pclass	Age	Fare
count	891.000000	891.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	32.204208
std	257.353842	0.486582	0.836071	13.002015	49.693429
min	1.000000	0.000000	1.000000	0.420000	0.000000
25%	223.500000	0.000000	2.000000	22.000000	7.915400
50%	446.000000	0.000000	3.000000	29.699118	14.454200
75%	668.500000	1.000000	3.000000	35.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	512.329200

The dataset contains:

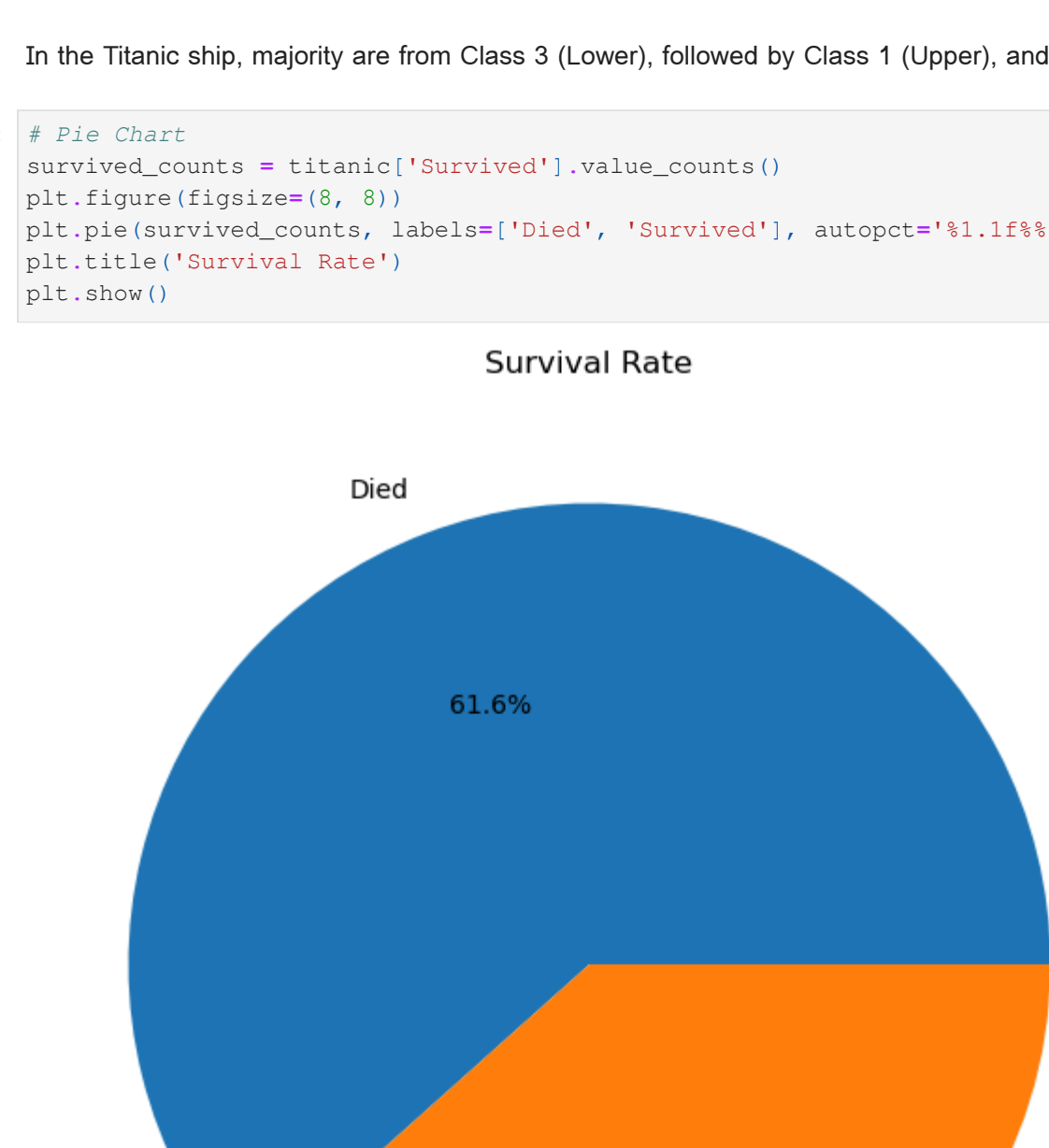
- PassengerId: A number system to indicate a passenger
- Survived: 0 = Died, 1 = Survived
- Pclass: Indicator of the passenger's social status (1 = Upper, 2 = Middle, 3 = Lower)
- Name: Passenger's name
- Sex: Passenger's gender
- Age: Passenger's age
- Fare: Ticket price

```
In [72]: # Bar plot
plt.figure(figsize=(8, 5))
titanic['Pclass'].value_counts().sort_index().plot(kind='bar')
plt.title('Number of Passengers in Each Class')
plt.xlabel('Passenger Class')
plt.ylabel('Count')
plt.grid(axis='y')
plt.show()
```



In the Titanic ship, majority are from Class 3 (Lower), followed by Class 1 (Upper), and Class 2 (Middle).

```
In [74]: # Pie chart
survived_counts = titanic['Survived'].value_counts()
plt.figure(figsize=(8, 8))
plt.pie(survived_counts, labels=['Died', 'Survived'], autopct='%1.1f%%')
plt.title('Survival Rate')
plt.show()
```



Most of the people on the Titanic ship did NOT survive (61.6%). Only 38.4% of the passengers survived..

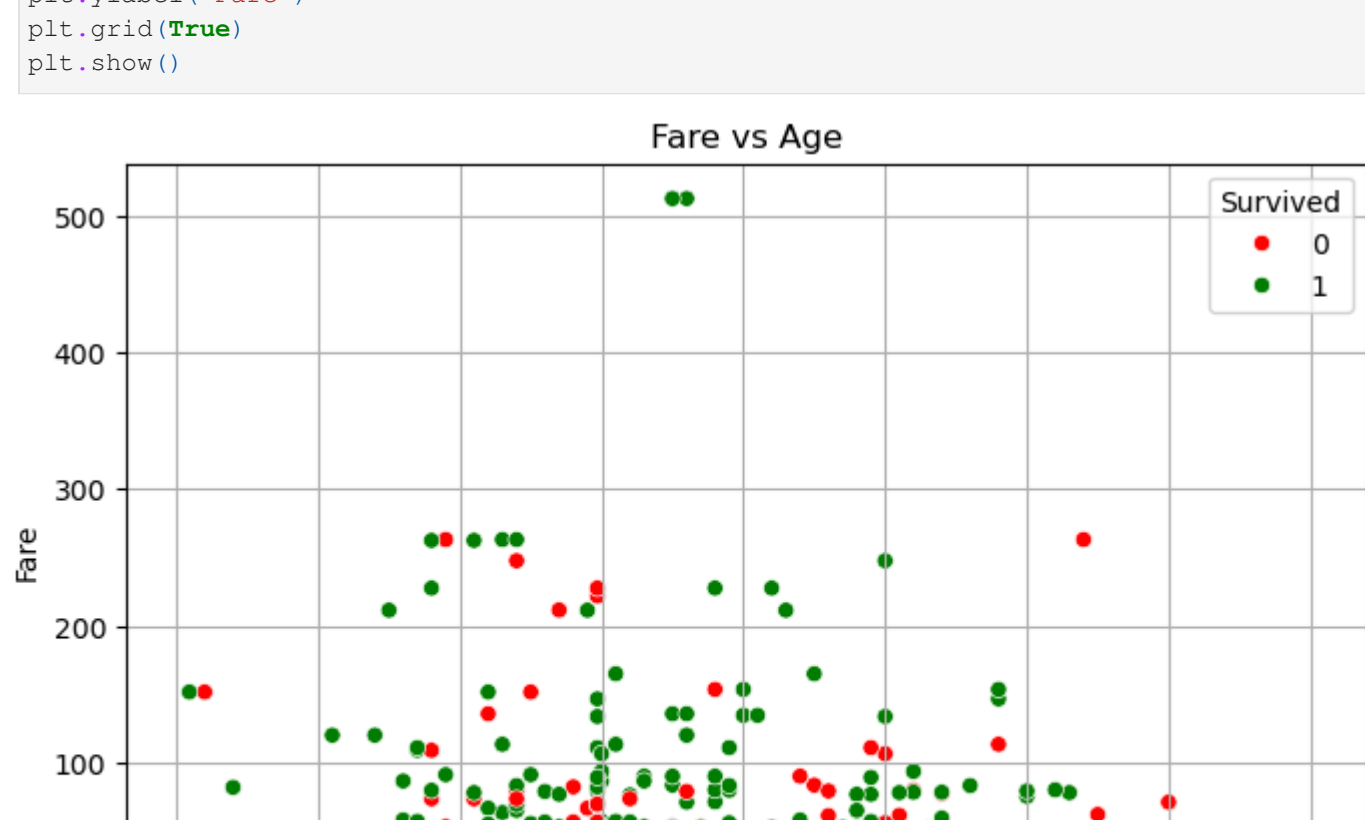
```
In [75]: # Line plot
plt.figure(figsize=(8, 5))
titanic['Age'].value_counts().sort_index().plot(kind='line')
plt.title('Age Distribution of Passengers')
plt.xlabel('Age')
plt.ylabel('Count')
plt.grid()
plt.show()
```



Most of the passengers are around 15-35 years old with majority being 29 years old. The youngest being below 1 year old and oldest being 80 years old.

```
In [81]: import seaborn as sns
```

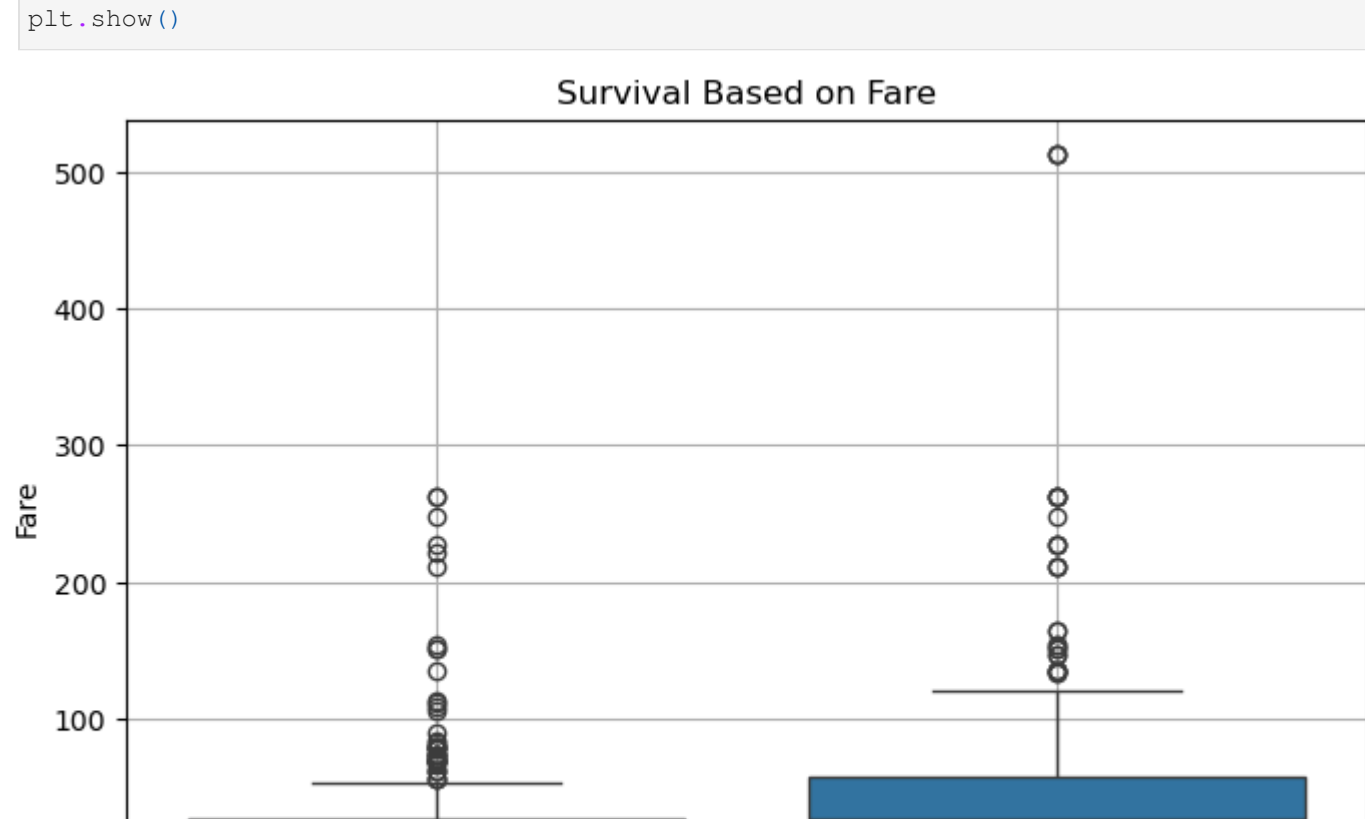
```
# Scatter plot
plt.figure(figsize=(8, 5))
sns.scatterplot(x='Age', y='Fare', hue='Survived', data=titanic, palette=['red', 'green'])
plt.title('Fare vs Age')
plt.xlabel('Age')
plt.ylabel('Fare')
plt.grid(True)
plt.show()
```



Most people would've thought that fare doesn't really contribute to survival. However, we can see that there is more green (Survived) concentrated at the upper side of y-axis (More expensive fare) while red (Died) are more concentrated at the lower side of y-axis (Cheaper fare).

```
In [82]: # Find the relationship between Fare and Survived
```

```
plt.figure(figsize=(8, 5))
sns.boxplot(x='Survived', y='Fare', data=titanic)
plt.title('Survival Based on Fare')
plt.xlabel('Survived')
plt.ylabel('Fare')
plt.grid(True)
plt.show()
```



This boxplot just proves it.

Summary of findings:

