# </talentlabs>

# Assignment 6 - Exploratory Data Analysis and Plotting Systems in Python (21 points)

## Instructions

1. Answer the below question in the boxes if needed.
2. For coding exercises, code in a single google colab notebook and zip all your code before submission.
3. Please submit the assignment through TalentLabs Learning System.

## Part 1: Concept Questions

**Question 1 (1 point)**
Which of these are graph plotting systems in Python? Select all that are correct.

1. Scikit - Learn
2. Pandas
3. Numpy
4. ggplot
5. Matplotlib
6. Tidyverse
7. Seaborn
8. Tableau

Matplotlib, Seaborn

</talentlabs>

**Question 2 (2 points)**
Mark the steps that are part of a Exploratory Data Analysis project.

1.  Build a model
2.  Plot a histogram and boxplot to answer a question
3.  Fetch data from a website
4.  Make a dashboard for your stakeholders.
5.  Removing Missing Values
6.  Look at outliers.
7.  Create tables and write data into a database

---

The steps of Exploratory Data Analysis project are as below:

1. Fetch data from a website.
2. Create tables and write data into a database.
3. Removing Missing Values.
4. Plot a histogram and boxplot to answer a question.
5. Look at outliers.
6. Build a model
7. Make a dashboard for your stakeholders.

---

**Question 3 (2 points)**
What inconsistencies do you spot in the `pakistan_intellectual_capital.csv` dataset?
We are looking for inconsistencies of the type:
-   data entry errors (could be related to different ways of looking at value or data type related)
-   missing values
-   duplicates

Tell us what inconsistencies do you spot in:
1)  Department
2)  Designation
3)  Year
4)  Country

Example:
Field: Terminal Degree
Inconsistencies: Duplicates
Explanation: This column has duplicates such as phd and PhD, MS and MSCS ( can be seen as data entry errors).

Tips

</talentlabs>

Your Answer:

1. Field: Department
Inconsistencies: Data entry error
Explanation: The column has inconsistency in data entry error, such as "Computer Sciences" and "Computer Science" shall be in same category.

2. Field: Designation
Inconsistencies: Missing values
Explanation: This column has missing values such as nan.

3. Field: Year
Inconsistencies: Missing values
Explanation: This column has missing values such as nan.

4. Field: Country
Inconsistencies: Data entry error
Explanation: The column has inconsistency in data entry error, such as "German" and "german" shall be in same category.

**Question 4 (2 points)**
Match the examples below to where these types of analytics are (Descriptive or Predictive)?

| Data Analytics Example | Descriptive / Predictive |
|---|---|
| Early Detection of Allergic Reactions | **Predictive** |

</talentlabs>

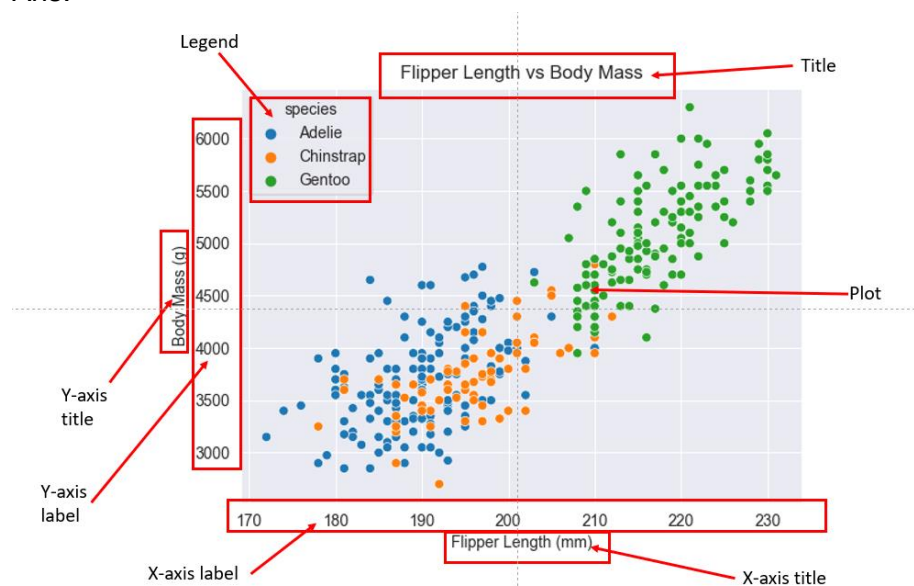| What genres and TV shows interest their subscribers most | **Descriptive** |
|---|---|
| Change in Year over Year customer behavior | **Descriptive** |
| Forecasting Future Cash Flow for a company | **Predictive** |

**Question 5 (2 points)**
Identify and label at least 5 elements of this graph. Annotate by editing the image.
Hint: Look at "elements of a graph" slides.



Ans:

</talentlabs>

**Question 6 (6 points)**

Load the titanic dataset using seaborn using:

```
import seaborn as sns
df = sns.load_dataset('titanic');
```

## Data Dictionary

| Variable | Definition | Key |
|---|---|---|
| survival | Survival | 0 = No, 1 = Yes |
| pclass | Ticket class | 1 = 1st, 2 = 2nd, 3 = 3rd |
| sex | Sex | |
| Age | Age in years | |
| sibsp | # of siblings / spouses aboard the Titanic | |
| parch | # of parents / children aboard the Titanic | |
| ticket | Ticket number | |
| fare | Passenger fare | |
| cabin | Cabin number | |
| embarked | Port of Embarkation | C = Cherbourg, Q = Queenstown, S = Southampton |

1. How many columns and rows does the dataset have? (½ point)

Rows: 891
Columns: 15

2. Print the column data types and number of missing values (½ point)

| No. | Column | Data type | Number of missing values |
|---|---|---|---|
| 1 | survived | int64 | 0 |
| 2 | pclass | int64 | 0 |
| 3 | sex | object | 0 |
| 4 | age | float64 | 177 |
| 5 | sibsp | int64 | 0 |

</talentlabs>

| 6 | parch | int64 | 0 |
|---|---|---|---|
| 7 | fare | float64 | 0 |
| 8 | embarked | object | 2 |
| 9 | class | category | 0 |
| 10 | who | object | 0 |
| 11 | adult_male | bool | 0 |
| 12 | deck | category | 688 |
| 13 | embark_town | object | 2 |
| 14 | alive | object | 0 |
| 15 | alone | bool | 0 |

3. Run descriptive statistics on the dataset and report the mean and standard deviation for
   - age
   - fare, and

   And the most frequent value for
   - sex
   - embark_town.

   (2 point)

**Age:**
Mean - 29.699118
Standard Deviation - 14.526497

**Fare:**
Mean - 32.204208
Standard Deviation - 49.693429
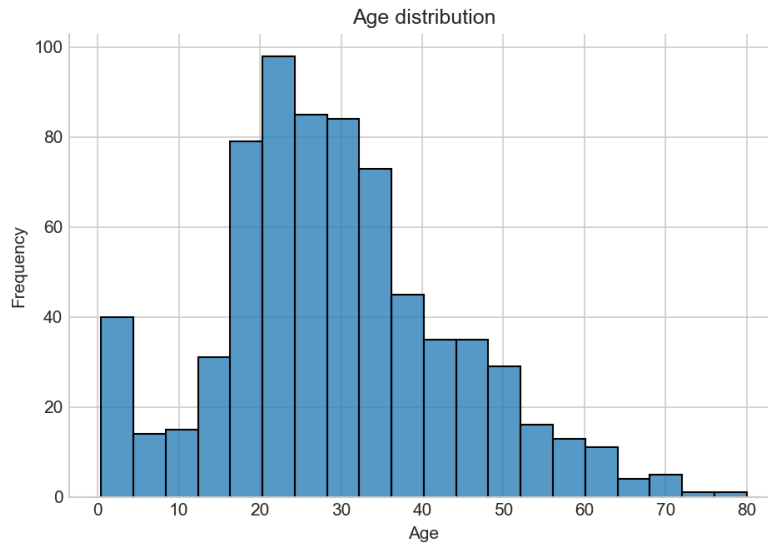
**Sex:**
Most frequent value -  "male"

**Embark_town:**
Most frequent value – "Sothampton"

4. The most convenient way to take a quick look at a univariate distribution in seaborn is the displot() function. By default, this will draw a histogram. Plot the histogram of age and add a title, x label, y label, gridlines. Count all the infants on board (age less than 3) and all the children ages 5-10. (3 points)

# </talentlabs>

Screenshot of the chart:



Age distribution

Number of infants onboard (age less than 3): 24

Number of children ages 5-10:  24

# </talentlabs>

## Part 2: Coding Exercises

For each of the exercises below, please write code in the same Google Colaboratory notebook or Jupyter Notebook, and create visualizations according to the instructions. You should also include the Google Colaboratory notebook or Jupyter Notebook in your submission.
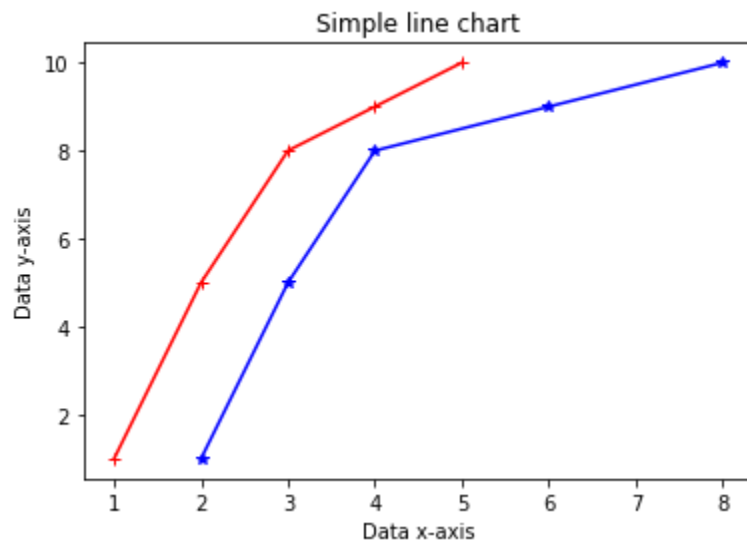
**Question 1 (2 points)**
Write a Python program to plot several lines with different format styles in one command using arrays. Also give it a title, x label and y label shown below for the chart. You could use any color for each.

The arrays are below:
```
a = np.array([1,2,3,4,5])
b = np.array([2,3,4,6,8])
c = np.array([1,5,8,9,10])
```

Use a and b on the x axis, c on the y axis.
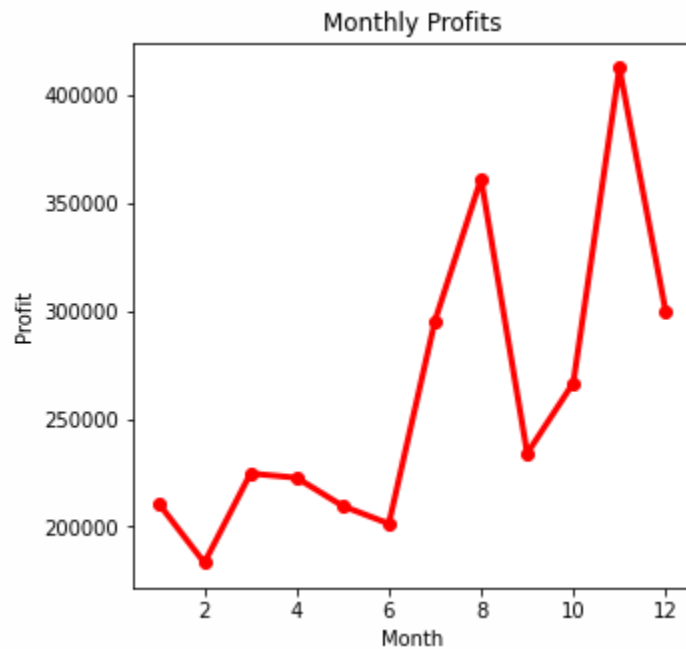
A sample output is included below:

</talentlabs>

**Question 2 (2 points)**

Ingest the company_sales_data.csv (attached in the assignment materials) and work to get total profit of all months and show line plot with the following style properties. Generated line plot must include following Style properties:

- X label name = Month
- Y label name = Profit
- Title Monthly Profits
- Add a circle marker.
- Make a line plot
- Line marker color as red
- Line width should be 3

Sample output:

</talentlabs>

**Question 3 (2 points)**

Ingest the company_sales_data.csv (attached in the assignment materials) and for each product column we see the number of units sold for various months, Read face cream and face wash product sales data and show it using the bar chart. The bar chart should display the number of units sold for facecream and face wash in the month of June. Add a separate bar for face cream and face wash in the same chart. Add title, labels mentioned below in sample. (2 points).

Sample output below:



Number of units sold for facecream and facewash in the month of June