

Name: Albo, Russel Zen D.

Course and Section: CPE32S9

Date of Submission: February 9, 2024

Instructor: Engr. Roman Richard

Objectives

- Part 1: Import the Libraries and Data
- Part 2: Plot the Data
- Part 3: Perform Simple Linear Regression on the SURVIVAL feature column (you can check the internet on how you can perform simple linear regression)

Train

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

train = "/content/titanic_train.csv"
TitanicTrain = pd.read_csv(train)

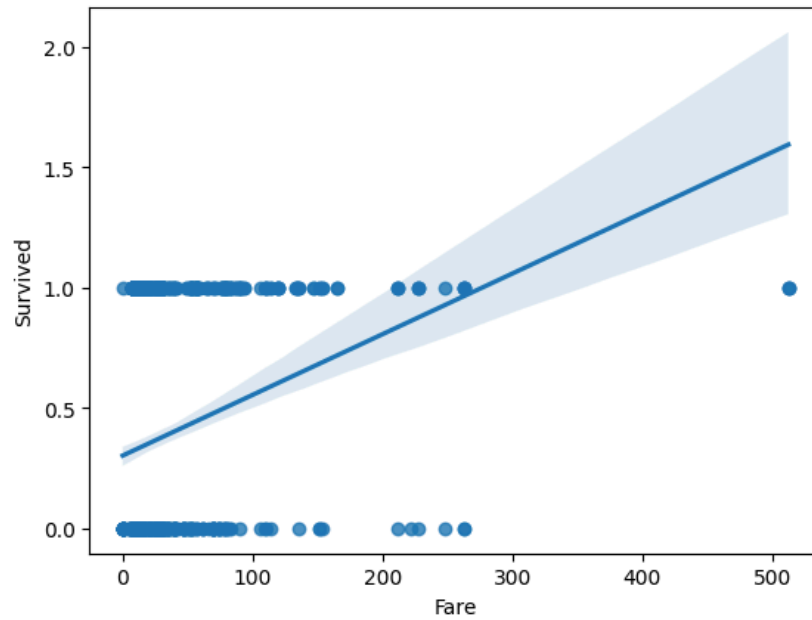
TitanicTrain.describe()
```

| | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|-------|-------------|------------|------------|------------|------------|------------|------------|
| count | 891.000000 | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000000 |
| mean | 446.000000 | 0.383838 | 2.308642 | 29.699118 | 0.523008 | 0.381594 | 32.204208 |
| std | 257.353842 | 0.486592 | 0.836071 | 14.526497 | 1.102743 | 0.806057 | 49.693429 |
| min | 1.000000 | 0.000000 | 1.000000 | 0.420000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 223.500000 | 0.000000 | 2.000000 | 20.125000 | 0.000000 | 0.000000 | 7.910400 |
| 50% | 446.000000 | 0.000000 | 3.000000 | 28.000000 | 0.000000 | 0.000000 | 14.454200 |
| 75% | 668.500000 | 1.000000 | 3.000000 | 38.000000 | 1.000000 | 0.000000 | 31.000000 |
| max | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | 6.000000 | 512.329200 |

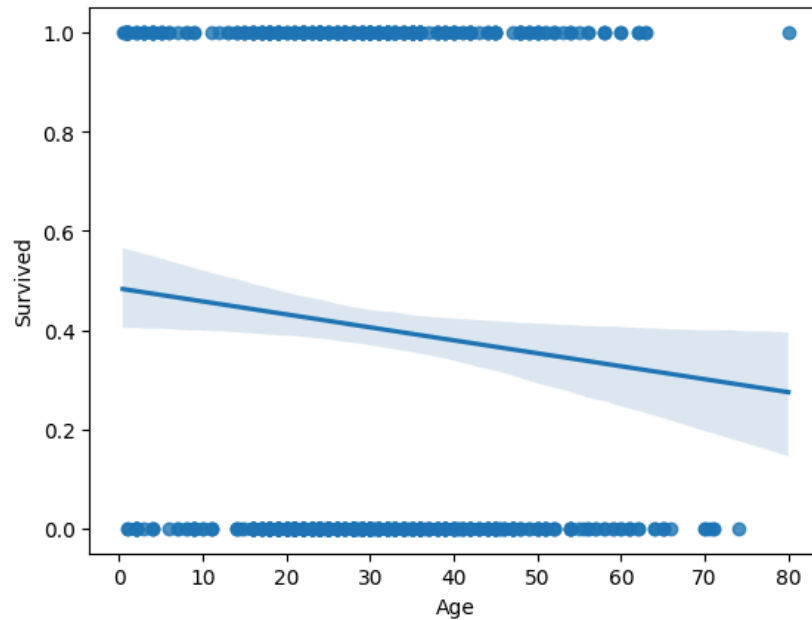
```
TitanicTrain.head()
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|-------------|----------|--------|----------------------------------------------------|--------|------|-------|-------|------------------|---------|-------|----------|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th...) | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |

```
sns.regplot(x="Fare", y="Survived", data=TitanicTrain);
```




```
sns.regplot(x="Age", y="Survived", data=TitanicTrain);
```



▼ Test

```
test = "/content/titanic_test.csv"  
TitanicTest = pd.read_csv(test)
```

```
TitanicTest.describe()
```

| | PassengerId | Pclass | Age | SibSp | Parch | Fare |  |
|-------|-------------|------------|------------|------------|------------|------------|-------------------------------------------------------------------------------------|
| count | 418.000000 | 418.000000 | 332.000000 | 418.000000 | 418.000000 | 417.000000 |  |
| mean | 1100.500000 | 2.265550 | 30.272590 | 0.447368 | 0.392344 | 35.627188 | |
| std | 120.810458 | 0.841838 | 14.181209 | 0.896760 | 0.981429 | 55.907576 | |
| min | 892.000000 | 1.000000 | 0.170000 | 0.000000 | 0.000000 | 0.000000 | |
| 25% | 996.250000 | 1.000000 | 21.000000 | 0.000000 | 0.000000 | 7.895800 | |
| 50% | 1100.500000 | 3.000000 | 27.000000 | 0.000000 | 0.000000 | 14.454200 | |
| 75% | 1204.750000 | 3.000000 | 39.000000 | 1.000000 | 0.000000 | 31.500000 | |
| max | 1309.000000 | 3.000000 | 76.000000 | 8.000000 | 9.000000 | 512.329200 | |

TitanicTest.head()

| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|-------------|--------|----------------------------------------------|--------|------|-------|-------|---------|---------|-------|----------|
| 0 | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | Q |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | S |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | Q |
| 3 | 895 | 3 | Wirz, Mr. Albert | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | S |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN | S |

Part 2

```
#Code cell 1
#import pandas
import pandas as pd
path = "/content/titanic_train.csv"
training = pd.read_csv(path)
```

```
#Code cell 2
#verify the contents of the training dataframe using the pandas info() method.
#training.?
training.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 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          714 non-null    float64
6   SibSp        891 non-null    int64
7   Parch        891 non-null    int64
8   Ticket       891 non-null    object
9   Fare         891 non-null    float64
10  Cabin        204 non-null    object
11  Embarked     889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

Are there missing values in the data set?

- Yes, there are missing values in the dataset

```
#Code cell 3
#view the first few rows of the data
#
training.head(5)
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare |
|---|-------------|----------|--------|---------------------------------------------------|--------|------|-------|-------|-----------|---------|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 |



Next steps: [View recommended plots](#)

```
#code cell 4
training["Sex"] = training["Sex"].apply(lambda toLabel: 0 if toLabel == 'male' else 1)
```

```
#code cell 5
#view the first few rows of the data again
training.head(5)
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare |
|---|-------------|----------|--------|---------------------------------------------------|-----|------|-------|-------|-----------|---------|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | 0 | 22.0 | 1 | 0 | A/5 21171 | 7.2500 |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | 1 | 38.0 | 1 | 0 | PC 17599 | 71.2833 |



Next steps: [View recommended plots](#)

```
# Code cell 6
training["Age"].fillna(training["Age"].mean(), inplace=True)

#code cell 7
#verify that the missing values for the age variable have been eliminated.
training["Age"].head(10)
```

```
0    22.000000
1    38.000000
2    26.000000
3    35.000000
4    35.000000
5    29.699118
6    54.000000
7     2.000000
8    27.000000
9    14.000000
Name: Age, dtype: float64
```

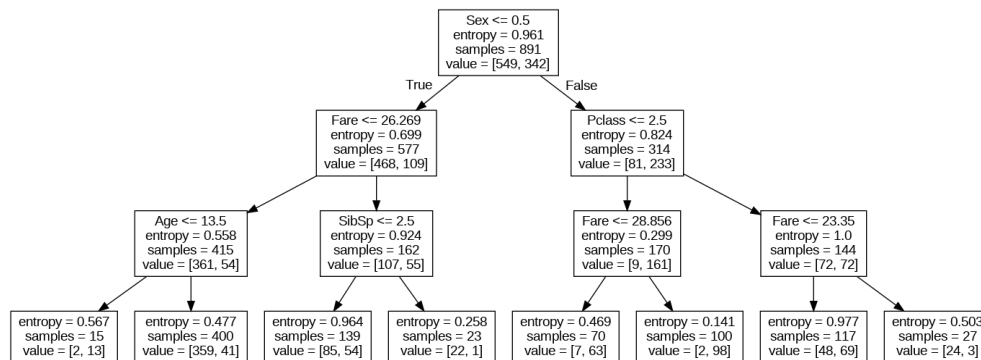
```
#code cell 8
#create the array for the target values
y_target = training["Survived"].values
```


0.8226711560044894

```
#code cell 12
from six import StringIO
with open("/content/titanic.dot", 'w') as f:
    f = tree.export_graphviz(clf_train, out_file=f, feature_names=columns)
```

```
#code cell 13
#run the Graphviz dot command to convert the .dot file to .png
!dot -Tpng /content/titanic.dot -o /content/titanic.png
```

```
#code cell 14
#import the Image module from the Ipython.display library
from IPython.display import Image
#display the decision tree graphic
Image("/content/titanic.png")
```



What describes the group that had the most deaths by number? Which group had the most survivors?

- The group that had the most deaths by numbers are Male while the most survivors are Female

```
#code cell 15
#import the file into the 'testing' dataframe.
path = "/content/titanic_test.csv"
testing = pd.read_csv(path)
testing.info(5)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 11 columns):
#   Column      Non-Null Count  Dtype
---  -
0   PassengerId  418 non-null    int64
1   Pclass       418 non-null    int64
2   Name         418 non-null    object
3   Sex          418 non-null    object
4   Age         332 non-null    float64
5   SibSp        418 non-null    int64
6   Parch        418 non-null    int64
7   Ticket       418 non-null    object
8   Fare         417 non-null    float64
9   Cabin        91 non-null     object
10  Embarked     418 non-null    object
dtypes: float64(2), int64(4), object(5)
memory usage: 36.0+ KB
```

How many records are in the data set?

- there are 418 records in the dataset

Which important variables(s) are missing values and how many are missing?

- The missing values of age are 86 and fare had 1 missing values

```
#code cell 16
#replace the Gender labels in the testing dataframe
# Hint: look at code cell 4
testing["Sex"] = testing["Sex"].apply(lambda toLabel: 0 if toLabel == 'male' else 1)
training.head(5)
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare |
|---|-------------|----------|--------|------------------------------------------------------------------|-----|------|-------|-------|-----------|---------|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | 0 | 22.0 | 1 | 0 | A/5 21171 | 7.2500 |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | 1 | 38.0 | 1 | 0 | PC 17599 | 71.2833 |

Next steps:

[View recommended plots](#)

```
#code cell 17
#Use the fillna method of the testing dataframe column "Age"
#to replace missing values with the mean of the age values.
testing["Age"].fillna(testing["Age"].mean(), inplace=True)
testing["Fare"].fillna(testing["Fare"].mean(), inplace=True)

#code cell 18
#verify the data preparation steps. Enter and run both the info and head
#methods from here, by entering and running one and then the other.
testing.info()
testing.head(10)
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 418 entries, 0 to 417
```

```
Data columns (total 11 columns):
```

```
#   Column      Non-Null Count  Dtype
---  -
0   PassengerId  418 non-null      int64
1   Pclass       418 non-null      int64
2   Name         418 non-null      object
3   Sex          418 non-null      int64
4   Age          418 non-null      float64
5   SibSp        418 non-null      int64
6   Parch        418 non-null      int64
7   Ticket       418 non-null      object
8   Fare         418 non-null      float64
9   Cabin        91 non-null       object
10  Embarked     418 non-null      object
```

```
dtypes: float64(2), int64(5), object(4)
```

```
memory usage: 36.0+ KB
```

| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|-------------|--------|----------------------------------|-----|------|-------|-------|--------|---------|-------|----------|
| 0 | 892 | 3 | Kelly, Mr. James | 0 | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | 1 | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | 0 | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | |
| 3 | 895 | 3 | Wirz, Mr. Albert | 0 | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander | 1 | 62.0 | 1 | 1 | 310406 | 10.0000 | NaN | |

Next steps:

[View recommended plots](#)

```
#code cell 19
```

```
#create the variable X_input to hold the features that the classifier will use
```

```
X_input = testing[list(columns)].values
```

```
#code cell 20
```

```
#apply the model to the testing data and store the result in a pandas dataframe.
```

```
#Use X_input as the argument for the predict() method of the clf_train classifier object
```

```
target_labels = clf_train.predict(X_input)
```

```
#convert the target array into a pandas dataframe using the pd.DataFrame() method and target as argument
```

```
target_labels = pd.DataFrame({'Est_Survival':target_labels, 'Name':testing['Name']})
```

```
testing.head(20)
```


| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cab |
|----|-------------|--------|----------------------------------------------|-----|----------|-------|-------|-----------|---------|-----|
| 0 | 892 | 3 | Kelly, Mr. James | 0 | 34.50000 | 0 | 0 | 330911 | 7.8292 | N |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | 1 | 47.00000 | 1 | 0 | 363272 | 7.0000 | N |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | 0 | 62.00000 | 0 | 0 | 240276 | 9.6875 | N |
| 3 | 895 | 3 | Wirz, Mr. Albert | 0 | 27.00000 | 0 | 0 | 315154 | 8.6625 | N |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | 1 | 22.00000 | 1 | 1 | 3101298 | 12.2875 | N |
| 5 | 897 | 3 | Svensson, Mr. Johan Cervin | 0 | 14.00000 | 0 | 0 | 7538 | 9.2250 | N |
| 6 | 898 | 3 | Connolly, Miss. Kate | 1 | 30.00000 | 0 | 0 | 330972 | 7.6292 | N |
| 7 | 899 | 2 | Caldwell, Mr. Albert Francis | 0 | 26.00000 | 1 | 1 | 248738 | 29.0000 | N |
| 8 | 900 | 3 | Abraham, Mrs. Joseph (Sophie Halaut Easu) | 1 | 18.00000 | 0 | 0 | 2657 | 7.2292 | N |
| 9 | 901 | 3 | Davies, Mr. John Samuel | 0 | 21.00000 | 2 | 0 | A/4 48871 | 24.1500 | N |
| 10 | 902 | 3 | Ilieff, Mr. Ylio | 0 | 30.27259 | 0 | 0 | 349220 | 7.8958 | N |

Next steps:

[View recommended plots](#)

```
#code cell 21
#import the numpy library as np
import numpy as np
# Load data for all passengers in the variable all_data
path = "/content/titanic_train.csv"
all_data = pd.read_csv(path)
# Merging using the field Name as key, selects only the rows of the two datasets that refer to the same passenger
testing_results = pd.merge(target_labels, all_data[['Name','Survived']], on=['Name'])
# Compute the accuracy as a ratio of matching observations to total osbervations. Store this in in the variable acc.
acc = np.sum(testing_results['Est_Survival'] == testing_results['Survived']) /float(len(testing_results))
# Print the result
print(acc)
```

1.0

```
#code cell 22
#import the titanic_all.csv file into a dataframe called all_data. Specify the list of columns to import.
path = "/content/titanic_train.csv"
all_data = pd.read_csv(path, usecols=['Survived','Pclass','Sex','Age','SibSp','Fare'])
#View info for the new dataframe
all_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
```

```
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Survived    891 non-null    int64
1   Pclass      891 non-null    int64
2   Sex         891 non-null    object
3   Age         714 non-null    float64
4   SibSp       891 non-null    int64
5   Fare        891 non-null    float64
dtypes: float64(2), int64(3), object(1)
memory usage: 41.9+ KB
```