Report on

# "Machine Learning"

B.E. [Computer Engineering]

Submitted By

**Shivani Jadhav**          **41**
**Vaishnavi Garghate**   **42**
**Siddhi Shimpi**           **43**
**Priti Aher**                  **44**

Under the guidance of

## Prof.Snehal Kamlapur

Academic Year: 2022-2023

Department of Computer Engineering

# Table of Contents

# Problem Statement

Build a machine learning model that predicts the type of people who survived the Titanic shipwreck using passenger data

# Introduction

The model predicts whether a passenger would survive on the titanic taking into account and comparing and finding relations amongst various features. A tragedy like the sinking of the RMS Titanic in 1912, four days into the maiden voyage of the world's largest ship, can be analyzed from many angles: the historical significance, the geopolitical consequences. The Titanic dataset provided by Kaggle is split into train and test files. The training file contains a variable called Survived (representing the number of survivors), which is our target. After downloading the dataset, you can perform an automatic Exploratory Data Analysis (EDA) to get a taste of the available variables. access to two similar datasets that include passenger information like name, age, gender, socio-economic class, etc. One dataset is titled train.csv and the other is titled test.csv.Train.csv will contain the details of a subset of the passengers on board (891 to be exact) and importantly, will reveal whether they survived or not, also known as the "ground truth". The test.csv dataset contains similar information but does not disclose the "ground truth" for each passenger. It's your job to predict these outcomes. Using the patterns you find in the train.csv data, predict whether the other 418 passengers on board (found in test.csv) survived.

## Importing Dependencies

We will be using: NumPy, pandas, matplotlib, seaborn, sklearn,

As we move ahead, you will get to know the use of each of these modules.

Now, we need to upload the downloaded dataset, into this program, so that our code can read the data and perform the necessary actions using it.
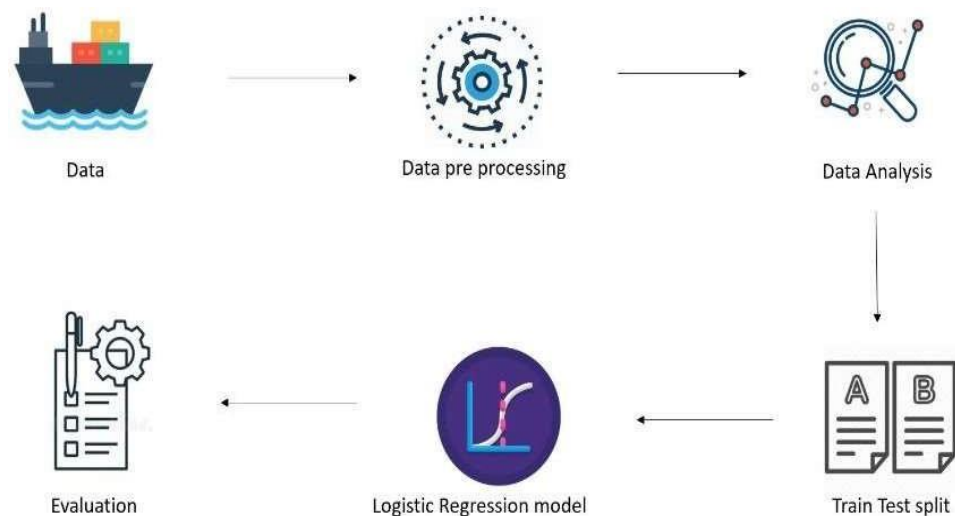
As we have downloaded a CSV file, we shall be using Pandas to store that data in a variable.

Our dataset is now stored in the variable named titanic_data.

To get a brief idea about how the data is loaded, we use the command "variable_name.head()" to get a glimpse of the dataset in the form of a table.

## WorkFlow Of Titanic Machine Learning Model:



Data → Data pre processing → Data Analysis → Train Test split → Logistic Regression model → Evaluation
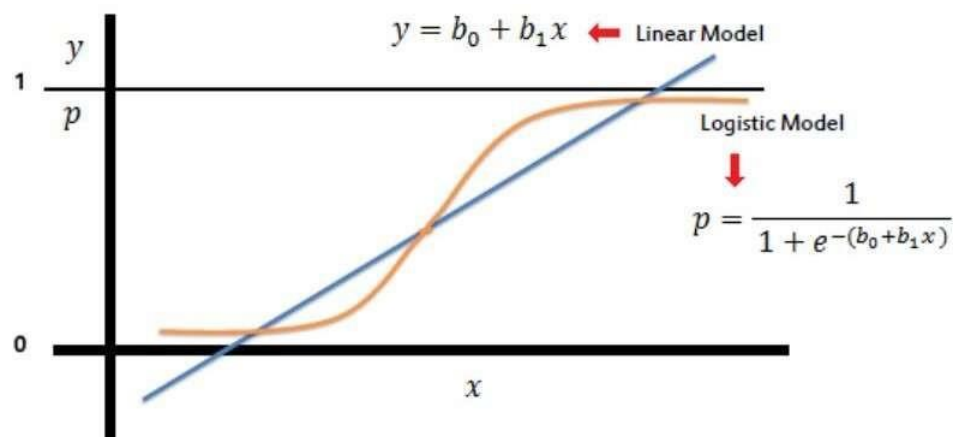
## Logistic Regression:

A simple yet crisp description of Logistic Description would be, "it is a supervised learning classification algorithm used to predict the probability of a target variable. The nature of target or dependent variable is dichotomous, which means there would be only two possible classes." Logistic Regression was used in the biological sciences in early twentieth century. It was then used in many social science applications. Logistic Regression is used when the dependent variable(target) is categorical.

For example,

To predict whether an email is spam (1) or (0)

Whether the tumor is malignant (1) or not (0)

The graph of logistic regression is as shown below:

# Implementation

Importing the Dependencies

```
[1] import numpy as np
    import pandas as pd
    import matplotlib.pyplot as plt
    import seaborn as sns
    from sklearn.model_selection import train_test_split
    from sklearn.linear_model import LogisticRegression
    from sklearn.metrics import accuracy_score
```

Data Collection & Processing

```
[2] # load the data from csv file to Pandas DataFrame
    titanic_data = pd.read_csv('/content/train.csv')
```

```
[3] # printing the first 5 rows of the dataframe
    titanic_data.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 |

| [3] | 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 |

```
[4] # number of rows and Columns
    titanic_data.shape

    (891, 12)
```

```
⏵  # getting some informations about the data
    titanic_data.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
```

```
⏵  # check the number of missing values in each column
    titanic_data.isnull().sum()
```

```
😊 PassengerId    0
    Survived       0
    Pclass         0
    Name           0
    Sex            0
    Age            177
    SibSp          0
    Parch          0
    Ticket         0
    Fare           0
    Cabin          687
    Embarked       2
    dtype: int64
```

Handling the Missing values

```
[ ] # drop the "Cabin" column from the dataframe
    titanic_data = titanic_data.drop(columns='Cabin', axis=1)
```

```
# finding the mode value of "Embarked" column
print(titanic_data['Embarked'].mode())
```

```
0    S
dtype: object
```

```
print(titanic_data['Embarked'].mode()[0])
```

```
S
```

```
# replacing the missing values in "Embarked" column with mode value
titanic_data['Embarked'].fillna(titanic_data['Embarked'].mode()[0], inplace=True)
```

```
# check the number of missing values in each column
titanic_data.isnull().sum()
```

```
PassengerId    0
Survived       0
Pclass         0
Name           0
Sex            0
Age            0
SibSp          0
Parch          0
```

```
# getting some statistical measures about the data
titanic_data.describe()
```

|       | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|-------|-------------|----------|--------|-----|-------|-------|------|
| count | 891.000000 | 891.000000 | 891.000000 | 891.000000 | 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 | 13.002015 | 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 | 22.000000 | 0.000000 | 0.000000 | 7.910400 |
| 50%   | 446.000000 | 0.000000 | 3.000000 | 29.699118 | 0.000000 | 0.000000 | 14.454200 |
| 75%   | 668.500000 | 1.000000 | 3.000000 | 35.000000 | 1.000000 | 0.000000 | 31.000000 |
| max   | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | 6.000000 | 512.329200 |

```
# finding the number of people survived and not survived
titanic_data['Survived'].value_counts()
```
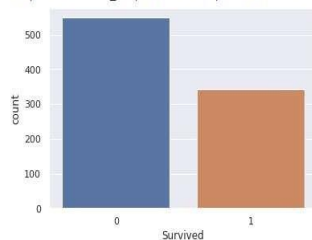
```
0    549
1    342
Name: Survived, dtype: int64
```

## Data Visualization

```
[ ]  sns.set()
```

```
# making a count plot for "Survived" column
sns.countplot('Survived', data=titanic_data)
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only vali
  FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7fd6c77f16d0>
```
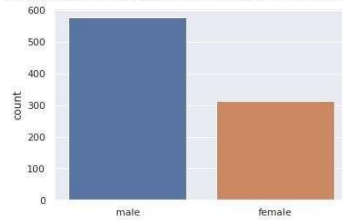
```
titanic_data['Sex'].value_counts()
```

```
male      577
female    314
Name: Sex, dtype: int64
```

```
# making a count plot for "Sex" column
sns.countplot('Sex', data=titanic_data)
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only vali
  FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7fd6cbeb1d90>
```
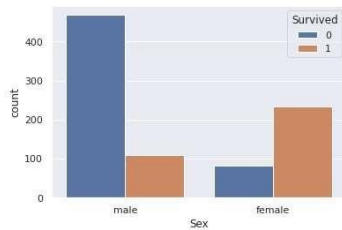


```
# number of survivors Gender wise
sns.countplot('Sex', hue='Survived', data=titanic_data)
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only vali
  FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7fd6c77d0dd0>
```
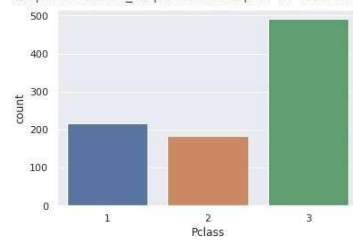


```
# making a count plot for "Pclass" column
sns.countplot('Pclass', data=titanic_data)
```

```
# making a count plot for "Pclass" column
sns.countplot('Pclass', data=titanic_data)
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only val
  FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7fd6c5f7bfd0>
```



```
sns.countplot('Pclass', hue='Survived', data=titanic_data)
```

```
titanic_data['Sex'].value_counts()
```

```
male      577
female    314
Name: Sex, dtype: int64
```

```
titanic_data['Embarked'].value_counts()
```

```
S    646
C    168
Q     77
Name: Embarked, dtype: int64
```

```
# converting categorical Columns
```

```
titanic_data.replace({'Sex':{'male':0,'female':1}, 'Embarked':{'S':0,'C':1,'Q':2}}, inplace=True)
```

```
titanic_data.head()
```

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

```
X = titanic_data.drop(columns = ['PassengerId','Name','Ticket','Survived'],axis=1)
Y = titanic_data['Survived']
```

```
print(X)
```

```
     Pclass  Sex        Age  SibSp  Parch      Fare  Embarked
0         3    0  22.000000      1      0    7.2500         0
1         1    1  38.000000      1      0   71.2833         1
2         3    1  26.000000      0      0    7.9250         0
3         1    1  35.000000      1      0   53.1000         0
4         3    0  35.000000      0      0    8.0500         0
..      ...  ...        ...    ...    ...       ...       ...
886       2    0  27.000000      0      0   13.0000         0
887       1    1  19.000000      0      0   30.0000         0
888       3    1  29.699118      1      2   23.4500         0
889       1    0  26.000000      0      0   30.0000         1
890       3    0  32.000000      0      0    7.7500         2

[891 rows x 7 columns]
```

```
print(Y)
```

```
0    0
1    1
2    1
```

Splitting the data into training data & Test data

```
[ ]  X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size=0.2, random_state=2)
```

```
[ ]  print(X.shape, X_train.shape, X_test.shape)
```

```
     (891, 7) (712, 7) (179, 7)
```

Model Training

Logistic Regression

```
⏵  model = LogisticRegression()
```

```
                                                   + Code      + Text
```

```
[ ]  # training the Logistic Regression model with training data
     model.fit(X_train, Y_train)
```

```
     /usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:940: ConvergenceWarning: lbfgs failed to converge (status=1):
     STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

```
# accuracy on training data
X_train_prediction = model.predict(X_train)
```

```
print(X_train_prediction)
```

```
[0 1 0 0 0 0 0 1 0 0 0 1 0 0 1 0 0 1 0 1 0 0 0 0 0 1 0 0 1 0 0 1 0 1 1 0 0 1 0 1
 0 0 0 0 0 0 1 1 0 0 1 0 1 0 1 0 0 0 0 0 0 1 0 1 0 0 1 1 0 0 1 1 0 1 0 0 1
 0 0 0 0 0 0 1 0 0 0 1 0 0 0 1 0 1 0 0 1 0 0 1 1 1 0 1 0 0 0 0 0 1 0 0 0
 1 1 0 0 1 0 0 1 0 0 1 0 0 1 0 1 0 1 0 1 0 1 1 1 1 1 1 0 0 1 1 1 0 0 1 0 0
 0 0 0 1 0 1 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 1 1 0 0 1 0 1 0 1 1 1
 0 0 1 0 0 0 1 0 0 1 1 0 1 0 0 0 0 0 1 1 0 1 1 1 1 0 0 0 0 0 0 0
 0 1 0 0 1 1 1 0 0 1 0 1 1 1 0 0 1 0 0 0 0 1 0 0 0 1 0 0 0 1 0 1 0 1 0 0 0
 0 0 0 0 0 0 1 0 1 0 0 1 0 0 1 0 1 0 1 1 0 0 0 0 1 0 1 0 0 1 0 0 0 1 0 0 0
 0 1 1 0 0 0 0 0 1 0 1 0 0 0 0 0 1 1 1 0 0 0 1 0 1 0 0 0 0 0 0 0 1 1 0 1 1
 0 1 1 1 0 0 0 0 0 0 0 0 1 0 0 1 1 1 0 1 0 0 0 0 1 1 0 0 0 1 0 1 1 1 0 0
 0 0 1 0 0 0 1 1 0 0 1 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 1 0 1 1 1 0 1 1 0 0 0
 0 1 0 1 0 0 1 1 0 0 0 0 1 0 0 0 0 1 1 0 1 0 1 0 0 0 0 0 0 1 0 0 0 0 1 1 0 0
 1 0 1 0 0 1 0 0 0 0 0 0 0 0 1 0 0 1 1 0 0 0 1 1 0 1 0 0 1 0 0 0 1 1 0 1 0
 0 0 0 0 1 0 0 1 0 1 1 0 0 1 0 0 1 0 0 0 1 0 1 1 0 0 1 1 0 1 0 1 1 1 0 1 0
 0 1 0 0 1 0 0 1 0 0 0 0 1 1 0 0 1 0 1 0 1 0 0 0 0 0 0 1 1 1 0 0 1 1 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 0 0 0 0
 0 0 1 0 0 0 0 0 1 0 1 0 1 0 0 1 0 1 0 1 1 1 0 0 0 1 0 1 0 0 0 0 1 1 1 0 0 1 1
 0 0 0 1 0 1 0 0 0 0 0 1 1 0 1 1 1 0 0 0 1 0 0 0 0 1 0 0 0 1 0 0 1 0 0 1 0 0 0
```

```
training_data_accuracy = accuracy_score(Y_train, X_train_prediction)
print('Accuracy score of training data : ', training_data_accuracy)
```

```
Accuracy score of training data :  0.8075842696629213
```

```
# accuracy on test data
X_test_prediction = model.predict(X_test)
```

```
print(X_test_prediction)
```

```
[0 0 1 0 0 0 0 0 0 0 1 1 0 0 1 0 0 1 0 0 1 0 1 1 0 1 0 1 1 0 0 0 0 0 0 0 0 1 1
 0 0 0 0 0 1 0 0 1 1 0 0 1 0 0 0 0 0 0 1 0 0 0 1 0 0 0 1 0 1 0 1 0 0 0 1 0 1 0
 1 0 0 0 1 0 1 0 0 0 1 1 0 0 1 0 0 0 0 0 0 1 0 1 0 0 1 0 1 1 1 1 0 0 0 0
 0 0 0 1 1 0 1 0 0 1 0 0 0 0 0 0 1 0 0 0 0 1 1 0 0 0 0 0 0 1 1 1 1 0 1 0 0
 0 1 0 0 0 0 1 0 0 1 1 0 1 0 0 0 1 1 0 0 1 0 0 1 1 1 0 0 0 0 0]
```

```
test_data_accuracy = accuracy_score(Y_test, X_test_prediction)
print('Accuracy score of test data : ', test_data_accuracy)
```

```
Accuracy score of test data :  0.7821229050279329
```

# Conclusion

Thus, we have Successfully build a machine learning model that predicts the type of people who survived the Titanic shipwreck using passenger data.