# <u>**SmartBridge Applied Data Science**</u>

<u>**Name:**</u> Routhu Sai Praneeth

<u>**Mail:**</u>  routhusai.praneeth2020@vitstudent.ac.in

## <u>**ADS Assignment 2**</u>

**Titanic Ship Case Study:**

## **Problem Description:**

On April 15, 1912, during her maiden voyage, the Titanic sank after colliding with an iceberg, killing 1502 out of 2224 passengers and crew. Translated 32% survival rate.

- One of the reasons that the shipwreck led to such loss of life was that there were not enough lifeboats for the passengers and crew.
- Although there was some element of luck involved in surviving the sinking, some groups of people were more likely to survive than others, such as women, children, and the upper-class.
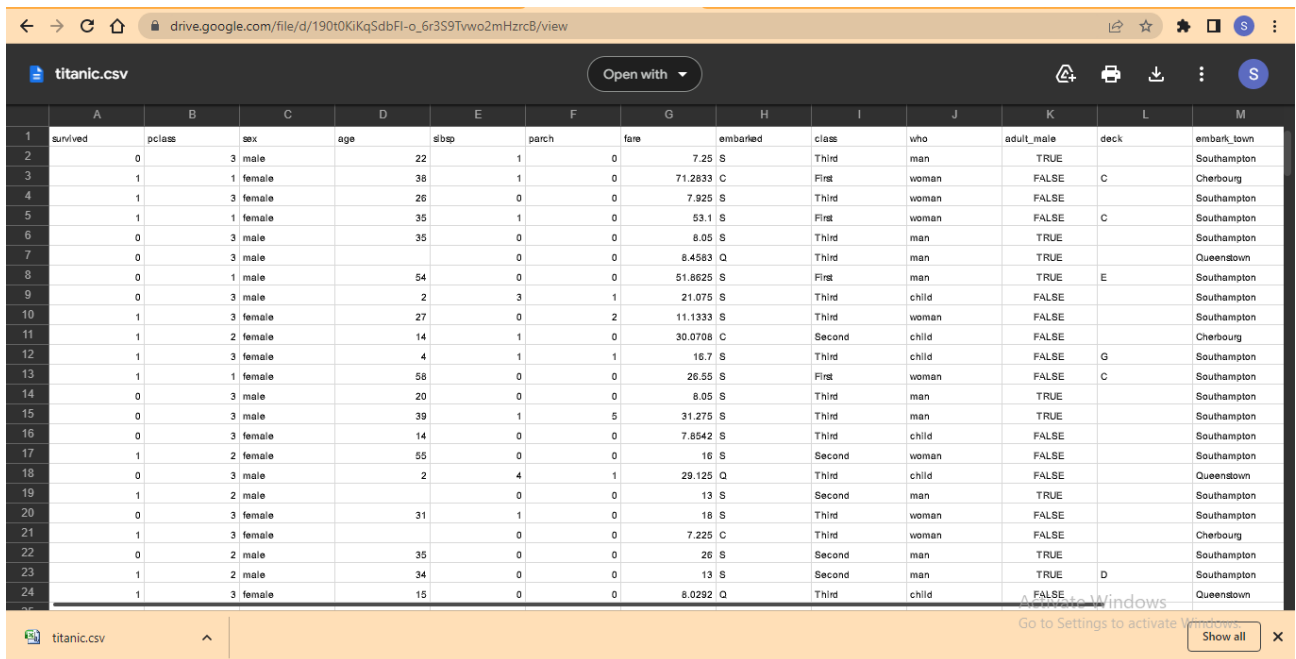
The problem associated with the Titanic dataset is to predict whether a passenger survived the disaster or not. The dataset contains various features such as passenger class, age, gender, cabin, fare, and whether the passenger had any siblings or spouses on board. These features can be used to build a predictive model to determine the likelihood of a passenger surviving the disaster. The dataset offers opportunities for feature engineering, data visualization, and model selection, making it a valuable resource for developing and testing data analysis and machine learning skills.

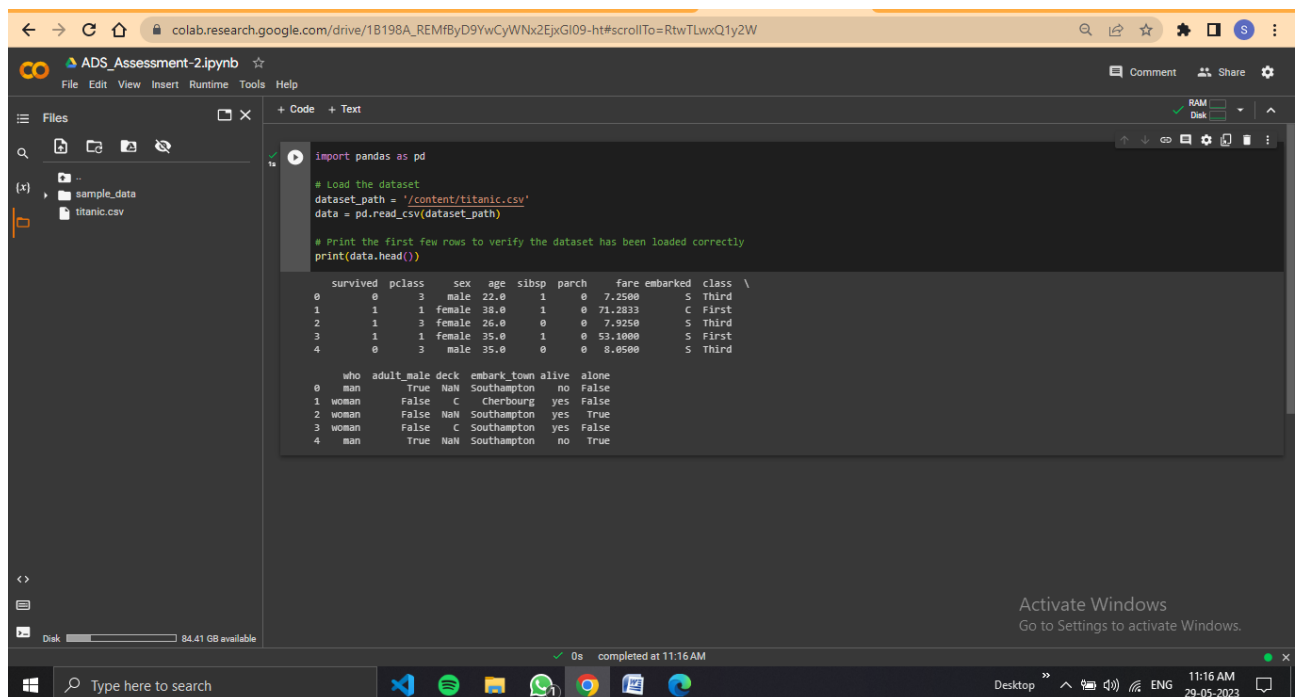## **Drive Link to Colab File**:

Link

## Tasks:

## 1.Downloading the dataset:



## 2.Load the dataset.

# 3.Perform Below Visualizations.



## • Univariate Analysis

● **Bi - Variate Analysis**

- **Multi - Variate Analysis**



## 4.Perform descriptive statistics on the dataset.



## 5.Handle the Missing values.

# 6.Find the outliers and replace the outliers

7.

```python
# Replace outliers in Fare column with median
Q1 = data['fare'].quantile(0.25)
Q3 = data['fare'].quantile(0.75)
IQR = Q3 - Q1

lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

data['fare'] = data['fare'].apply(lambda x: data['fare'].median() if x < lower_bound or x > upper_bound else x)
```

```python
import pandas as pd

# Load the dataset
dataset_path = '/content/titanic.csv'  # Update with the actual path to the dataset
data = pd.read_csv(dataset_path)

# Find outliers using IQR method for numeric columns
numeric_columns = ['age', 'fare']  # Update with the actual numeric columns in your dataset

for column in numeric_columns:
    Q1 = data[column].quantile(0.25)
    Q3 = data[column].quantile(0.75)
    IQR = Q3 - Q1

    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR

    outliers = data[(data[column] < lower_bound) | (data[column] > upper_bound)]

    print(f"Outliers in {column}:")
    print(outliers)
```

```
Outliers in age:
     survived  pclass   sex   age  sibsp  parch     fare embarked  class  \
33          0       2  male  66.0      0      0  10.5000        S  Second
54          0       1  male  65.0      0      1  61.9792        C   First
96          0       1  male  71.0      0      0  34.6542        C   First
116         0       3  male  70.5      0      0   7.7500        Q   Third
280         0       3  male  65.0      0      0   7.7500        Q   Third
456         0       1  male  65.0      0      0  26.5500        S   First
493         0       1  male  71.0      0      0  49.5042        C   First
630         1       1  male  80.0      0      0  30.0000        S   First
672         0       2  male  70.0      0      0  10.5000        S  Second
745         0       1  male  70.0      1      1  71.0000        S   First
851         0       3  male  74.0      0      0   7.7750        S   Third

     who  adult_male deck  embark_town alive  alone
```

```python
# Replace outliers in age column with median
Q1 = data['age'].quantile(0.25)
Q3 = data['age'].quantile(0.75)
IQR = Q3 - Q1

lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

data['age'] = data['age'].apply(lambda x: data['age'].median() if x < lower_bound or x > upper_bound else x)
```
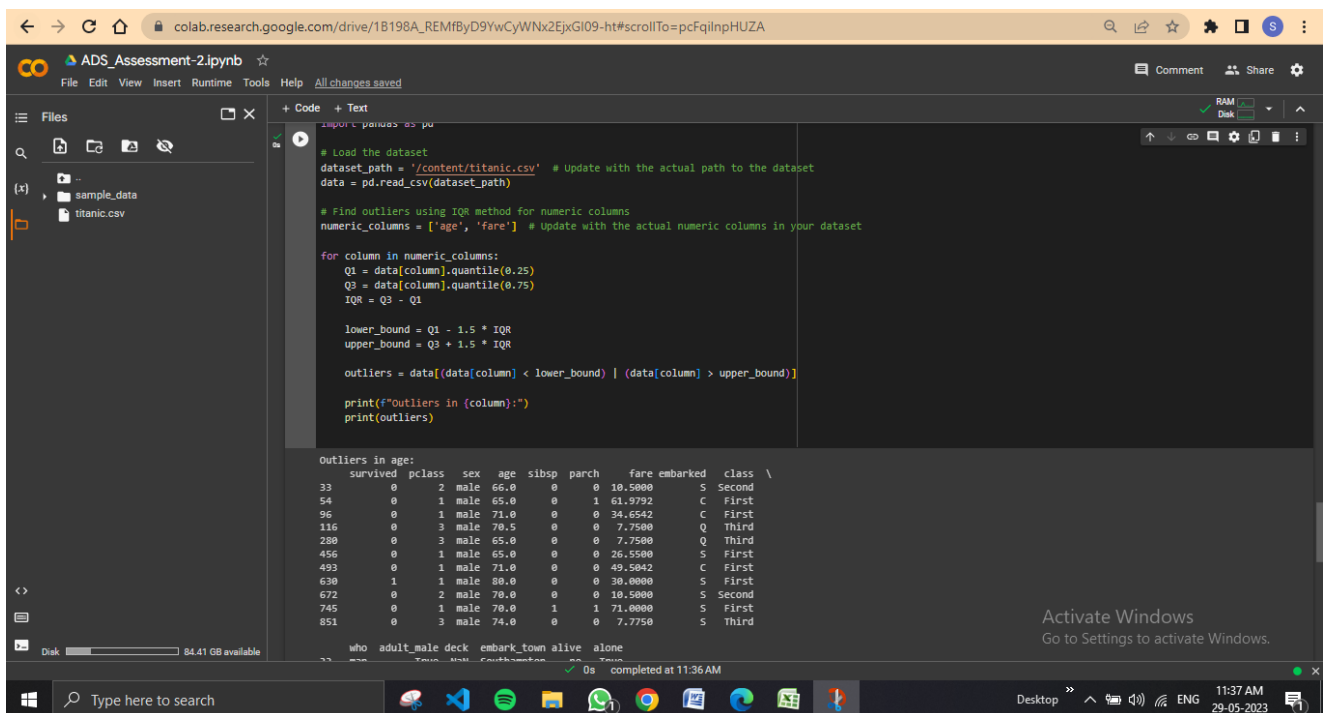
# 7.Check for Categorical columns and perform encoding.



Check for Categorical columns and perform encoding:

```python
[15]  # One-hot encoding for Survived column
      data = pd.get_dummies(data, columns=['survived'], drop_first=True)
```

```python
import pandas as pd

# Define the mapping dictionary
sex_mapping = {'male': 0, 'female': 1}

# Apply max-min encoding to the 'Sex' column
data['Sex_encoded'] = data['sex'].map(sex_mapping)

# Print the updated 'Sex_encoded' column
print(data['Sex_encoded'])
```
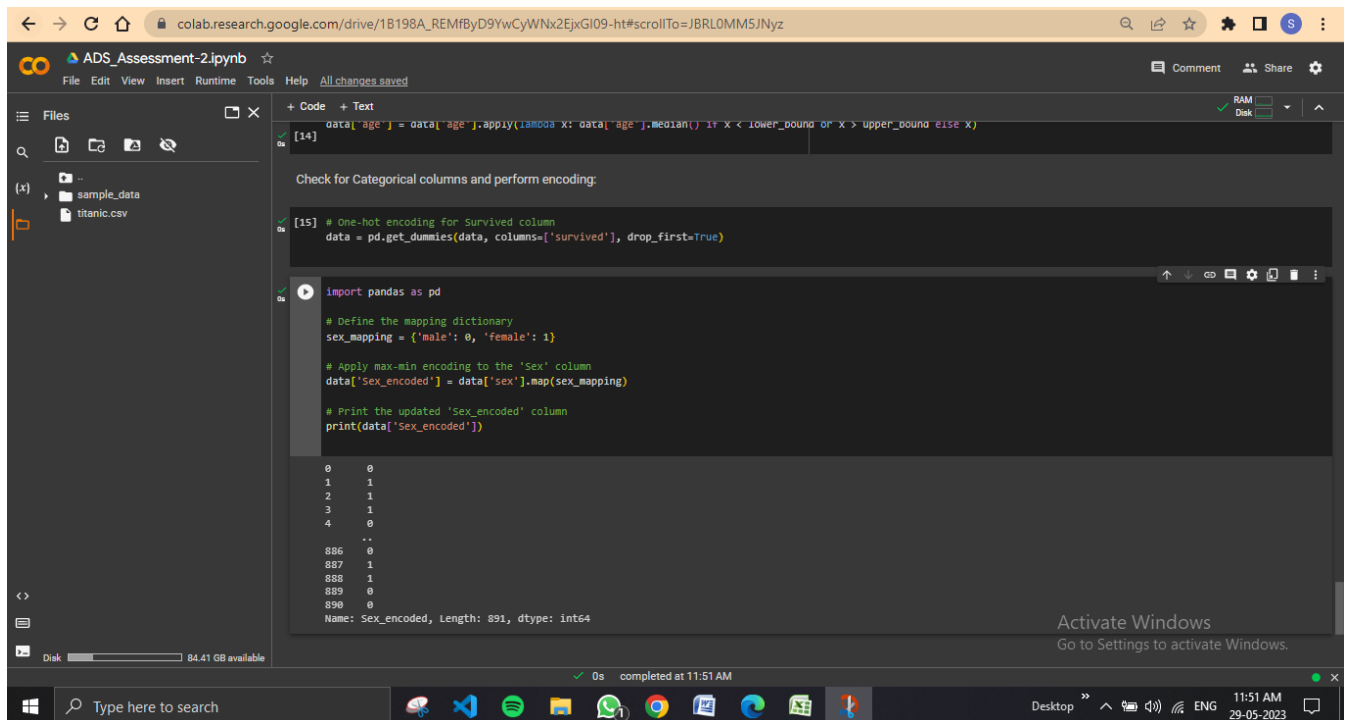
```
0       0
1       1
2       1
3       1
4       0
       ..
886     0
887     1
888     1
889     0
890     0
Name: Sex_encoded, Length: 891, dtype: int64
```



```python
from sklearn.preprocessing import LabelEncoder

# Create an instance of LabelEncoder
label_encoder = LabelEncoder()

# Perform label encoding on the 'alone' column
data['alone_encoded'] = label_encoder.fit_transform(data['alone'])
print(data['alone_encoded'])
```

```
0       0
1       0
2       1
3       0
4       1
       ..
886     1
887     1
888     0
889     1
890     1
Name: alone_encoded, Length: 891, dtype: int64
```

## 8.Split the data into dependent and independent variables.
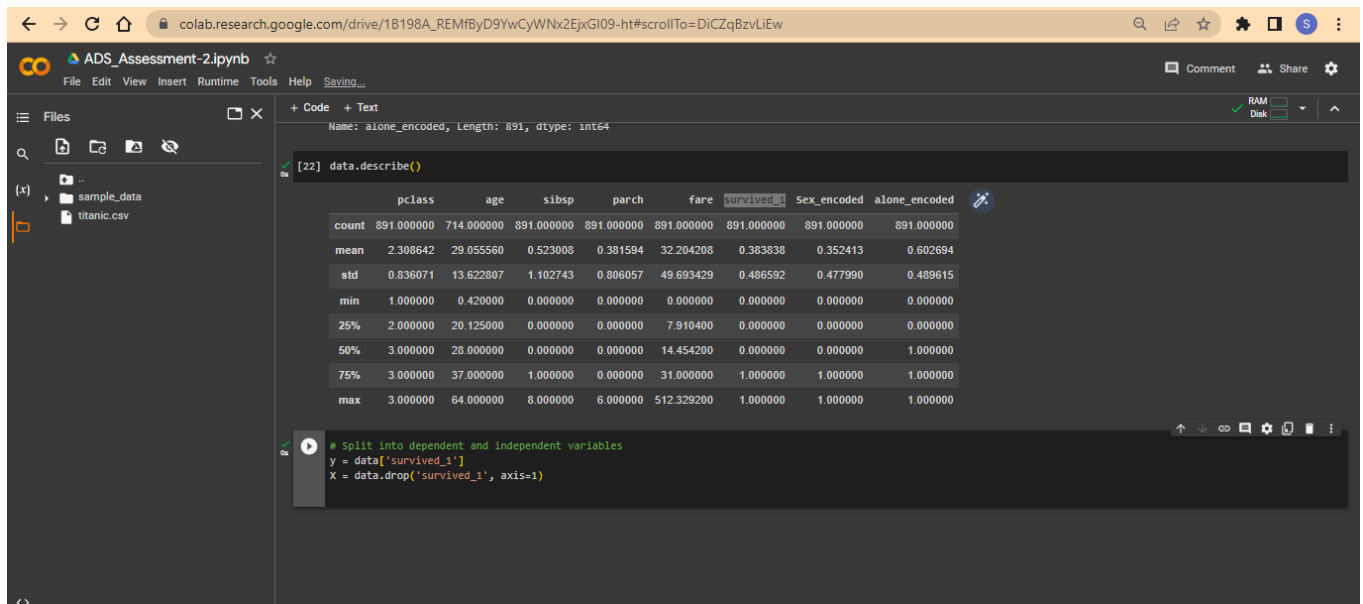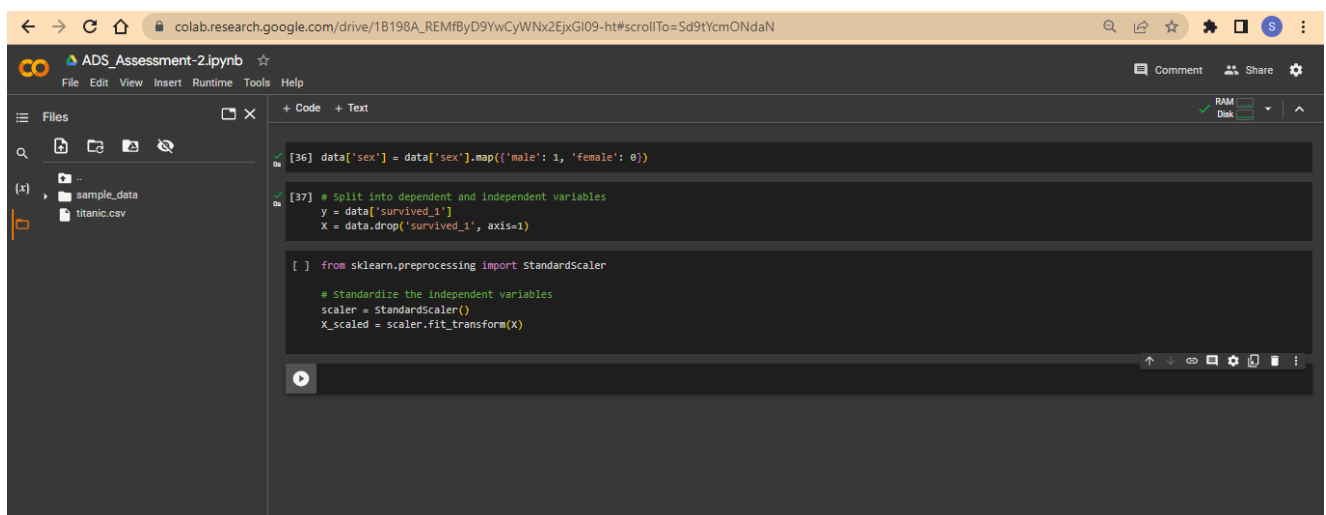


```
[22] data.describe()
```

|       | pclass     | age        | sibsp      | parch      | fare       | survived_1 | Sex_encoded | alone_encoded |
|-------|------------|------------|------------|------------|------------|------------|-------------|---------------|
| count | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000000 | 891.000000 | 891.000000  | 891.000000    |
| mean  | 2.308642   | 29.055560  | 0.523008   | 0.381594   | 32.204208  | 0.383838   | 0.352413    | 0.602694      |
| std   | 0.836071   | 13.622807  | 1.102743   | 0.806057   | 49.693429  | 0.486592   | 0.477990    | 0.489615      |
| min   | 1.000000   | 0.420000   | 0.000000   | 0.000000   | 0.000000   | 0.000000   | 0.000000    | 0.000000      |
| 25%   | 2.000000   | 20.125000  | 0.000000   | 0.000000   | 7.910400   | 0.000000   | 0.000000    | 0.000000      |
| 50%   | 3.000000   | 28.000000  | 0.000000   | 0.000000   | 14.454200  | 0.000000   | 0.000000    | 1.000000      |
| 75%   | 3.000000   | 37.000000  | 1.000000   | 0.000000   | 31.000000  | 1.000000   | 1.000000    | 1.000000      |
| max   | 3.000000   | 64.000000  | 8.000000   | 6.000000   | 512.329200 | 1.000000   | 1.000000    | 1.000000      |

```python
# Split into dependent and independent variables
y = data['survived_1']
X = data.drop('survived_1', axis=1)
```

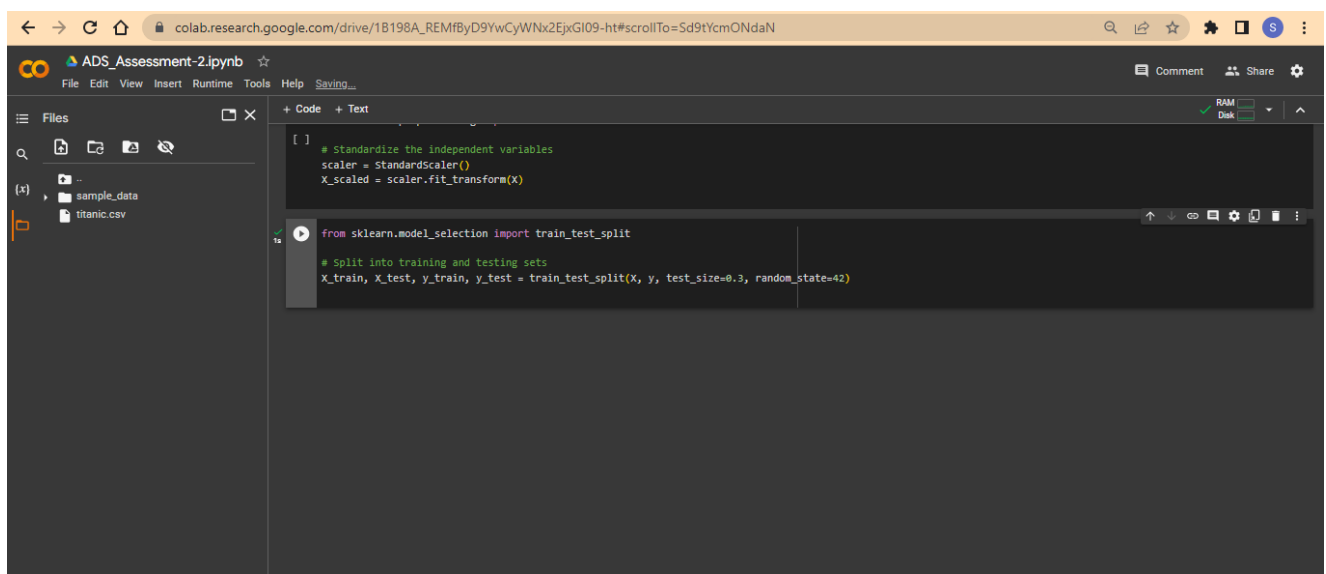## 9.Scale the independent variables.



```python
[36] data['sex'] = data['sex'].map({'male': 1, 'female': 0})
```

```python
[37] # Split into dependent and independent variables
y = data['survived_1']
X = data.drop('survived_1', axis=1)
```

```python
[ ] from sklearn.preprocessing import StandardScaler

# Standardize the independent variables
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

## 10.    Split the data into training and testing



```python
[ ]    # Standardize the independent variables
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

```python
from sklearn.model_selection import train_test_split

# Split into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```