# SpaceShip Titanic Prediction Report

# M.Junaid 109(4B)

#### **Table of Contents**

Introduction	
Data Loading and Exploration	2
Handling Missing Values	2
Filling Missing Values in Text Columns	2
Filling Missing Values in Number Columns	3
Converting Text to Numbers	3
Preparing Data for Training	3
Training the Model	
Testing Process	4
Loading and Preprocessing Test Data	4
Handling Missing Values in Test Data	4
Filling Missing Values and Encoding	5
Converting Text to Numbers and Scaling	5
Making Predictions on Test Data	6
Model Accuracy	7

#### Introduction

This report explains a machine learning project that predicts which passengers were transported to an alternate dimension during the SpaceShip Titanic collision with a spacetime anomaly. The code was written to preprocess data and create predictions that can be submitted to a Kaggle competition.

#### **Data Loading and Exploration**

First, we load the SpaceShip Titanic dataset and look at its basic information. This helps us understand what data we're working with, including passenger details, cabin information, and other features that might help predict transportation status.



### **Handling Missing Values**

We check which columns have missing values. We need to fill these missing values because machine learning models can't work with incomplete data.

```
null_columns = df.isnull().sum()
   print(f'Total number of columns with null values: {null_columns[null_columns > 0].count()}')
   print(null_columns[null_columns > 0])
Total number of columns with null values: 12
            201
HomePlanet
CryoSleep
               217
Cabin
               199
Destination
              182
              179
Age
VIP
              203
RoomService
            181
FoodCourt
              183
ShoppingMall 208
Spa
               183
VRDeck
               188
               200
Name
```

#### **Filling Missing Values in Text Columns**

For text columns like HomePlanet, CryoSleep, Destination, and VIP status, we fill missing values with the most common value in each column. This is a simple approach that works well for categorical data.

```
# Object columns filled

for column in null_columns[null_columns > 0].index:

    if df[column].dtype == 'object':

        mode_value = df[column].mode()[0]

        df[column] = df[column].fillna(mode_value).astype(str)

Python

College(M. Augusting object) | Complete | 14460| 1202721246 | pur Ext Entroplement | Despecting object dtype appars on filles | ffilled | fills | fills
```

#### **Filling Missing Values in Number Columns**

For number columns like Age, RoomService, FoodCourt, ShoppingMall, Spa, and VRDeck, we use a smarter approach. We first check for outliers (unusual values) using the IQR method. Then we fill missing values with the average of normal values. This helps prevent outliers from affecting our estimates.

```
for column in null_columns[null_columns > 0].index:
   if df[column].dtype in ['int64', 'float64']:
        Q1 = df[column].quantile(0.25)
        Q3 = df[column].quantile(0.75)
        IQR = Q3 - Q1

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

        if df[(df[column] < lower_bound) | (df[column] > upper_bound)].empty:
            mean_value = df[column].mean()
            df[column].fillna(mean_value,inplace=True)
        else:
            values = df[(df[column] >= lower_bound) & (df[column] <= upper_bound)][column]
            df[column].fillna(values.mean(),inplace=True)
            df.loc[(df[column] < lower_bound) | (df[column] > upper_bound), column] = values.mean()
```

#### **Converting Text to Numbers**

Machine learning models need numbers, not text. We convert all text columns to numbers using Label Encoding, which assigns numerical values to categorical features.

```
from sklearn.preprocessing import LabelEncoder

object_columns = df.select_dtypes(include='object').columns.tolist()
label_encoder = LabelEncoder()
for column in object_columns:
    df[column] = label_encoder.fit_transform(df[column].astype(str))
```

#### **Preparing Data for Training**

- Separate the features (X)
- Scale the data so all features have similar ranges, which helps the model learn better

```
from sklearn.preprocessing import StandardScaler
X=df.drop(['PassengerId', 'Transported','CryoSleep','VIP'], axis=1)
scaler = StandardScaler()
df[X.columns]=scaler.fit_transform(df[X.columns])
```

#### **Training the Model**

We use a Random Forest Classifier because it works well for classification problems. The model learns patterns from our training data to predict whether a passenger was transported or not.

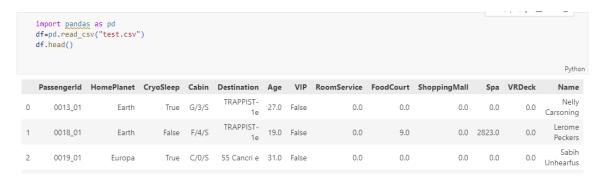
```
from sklearn.linear model import LogisticRegression
import joblib
train_data=df.drop(['PassengerId', 'Transported'], axis=1)
Y=df['Transported']

model = LogisticRegression()
model.fit(train_data, Y)
joblib.dump(model, 'Mymodel.pkl')
```

# **Testing Process**

#### **Loading and Preprocessing Test Data**

We load the test data and apply the same preprocessing steps as we did with the training data.



#### **Handling Missing Values in Test Data**

The test data also has missing values that need to be filled.

```
null_columns = df.isnull().sum()
  print(f'Total number of columns with null values: {null_columns[null_columns > 0].count()}')
  print(null_columns[null_columns > 0])
Total number of columns with null values: 12
HomePlanet 87
CryoSleep
              93
            100
Cabin
Destination 92
              91
Age
VIP
RoomService
              82
ShoppingMall 98
          101
Spa
VRDeck
```

#### Filling Missing Values and Encoding

We apply the same methods to fill missing values and encode categorical variables as we did with the training data.

```
for column in null_columns[null_columns > 0].index:
    if df[column].dtype == 'object':
        mode_value = df[column].mode()[0]
        df[column] = df[column].fillna(mode_value)
        df[column] = df[column].astype(str)
```

```
for column in null_columns[null_columns > 0].index:
    if df[column].dtype in ['int64', 'float64']:
        Q1 = df[column].quantile(0.25)
        Q3 = df[column].quantile(0.75)
        IQR = Q3 - Q1

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

        if df[(df[column] < lower_bound) | (df[column] > upper_bound)].empty:

            mean_value = df[column].mean()
            df[column].fillna(mean_value,inplace=True)
        else:

        values = df[(df[column] >= lower_bound) & (df[column] <= upper_bound)][column]
        df[column].fillna(values.mean(),inplace=True)
        df.loc[(df[column] < lower_bound) | (df[column] > upper_bound), column] = values.mean()
```

#### **Converting Text to Numbers and Scaling**

Machine learning models need numbers, not text. We convert all text columns to numbers using Label Encoding, which assigns numerical values to categorical features and scale the data.

```
from sklearn.preprocessing import LabelEncoder

object_columns = df.select_dtypes(include='object').columns.tolist()
label_encoder = LabelEncoder()
for column in object_columns:
    df[column] = label_encoder.fit_transform(df[column].astype(str))
```

```
from sklearn.preprocessing import StandardScaler
X=df.drop(['PassengerId','CryoSleep','VIP'], axis=1)
scaler = StandardScaler()
df[X.columns]=scaler.fit_transform(df[X.columns])
test_data=df.drop('PassengerId', axis=1)
```

#### **Making Predictions on Test Data**

- Loads our saved model
- Uses it to predict transportation status for the test data
- Creates a submission file with passenger IDs and predicted status
- Saves this file for uploading to Kaggle

```
import joblib
model = joblib.load('Mymodel.pkl')
```

```
predictions = model.predict(test_data)
submission_df = pd.DataFrame({'PassengerId': df['PassengerId'], 'Transported': predictions})
submission_df.to_csv('submission_file.csv', index=False)
```

# **Model Accuracy**

The model's performance was evaluated using metrics like Accuracy, Precision, Recall, and F1-Score. These metrics help us understand how well our model is predicting passenger transportation status.

