

Lab 2

Task: Kaggle Competition: Spaceship Titanic

Here's a step-by-step explanation of the code:

- **Import necessary libraries:**
 - Libraries such as `pandas`, `numpy`, `matplotlib`, `seaborn`, and various `scikit-learn` components are imported to handle data manipulation, visualization, model building, and evaluation.
- **Load the datasets:**
 - `train_df` and `test_df` are loaded from CSV files, representing the training and test datasets, respectively.
- **Data Exploration:**
 - The `head()`, `info()`, and `describe()` functions are used to get a quick look at the data, check for any missing data, and understand basic statistics like mean, median, etc.
- **Check for missing values:**
 - The `isnull().sum()` method checks if there are any missing (null) values in the columns of the training dataset.
- **Drop unnecessary columns:**
 - Unnecessary columns such as `PassengerId`, `Name`, and `Ticket` are dropped from both the training and test datasets, as they won't contribute to model prediction.
- **Handle missing values:**
 - Missing `Age` values are filled with the median of the column using `fillna()`.
 - Missing values for `Embarked` are filled with the most frequent value (mode) using `fillna()` as well.
- **One-Hot Encoding for categorical variables:**
 - The categorical variables `Pclass`, `Sex`, and `Embarked` are transformed into binary columns using one-hot encoding (via `pd.get_dummies()`), with `drop_first=True` to avoid multicollinearity.
- **Split the data into features and target:**
 - Features (X) are all the columns except `Transported`, while the target variable (y) is the `Transported` column.
- **Split the dataset into training and validation sets:**
 - The dataset is split into training (80%) and validation (20%) sets using `train_test_split()`.
- **Feature Scaling (Optional):**
 - Although `RandomForest` doesn't require feature scaling, `StandardScaler` is used to scale the features to a similar range (for consistency and future use with other models).
- **Initialize and train the model:**
 - A `RandomForestClassifier` model is created with 100 trees and trained on the training set (`x_train` and `y_train`).
- **Make predictions on the validation set:**

- ## OUTPUT

```
new_axis = axis.drop(labels, errors='errors')
^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^

File "C:\Users\Usman Ghani\AppData\Roaming\Python\Python312\site-packages\pandas\core\indexes\base.py", line 7070, in drop
    raise KeyError(f"{labels[mask].tolist()} not found in axis")
KeyError: "['PassengerId', 'Name', 'Ticket'] not found in axis"

PS C:\Users\Usman Ghani\Desktop\myworld> PassengerId  Pclass  Sex  Age  SibSp  Parch  Fare  Embarked  Transported
>> 0          1      3  male  22.0      1      0  7.25      S          0
>> 1          2      1 female  38.0      1      0 71.2833  C          1
>> 2          3      3 female  26.0      0      0  7.925   S          0
>> 3          4      1 female  35.0      1      0  53.1    S          0
>> 4          5      3  male  35.0      0      0  8.05    S          0
>>
```