

This code builds a machine learning pipeline using a **Random Forest Classifier** for classifying passengers on the Titanic based on their survival status. Below, I'll break down each part in detail:

## 1. Import Necessary Libraries

```
import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import accuracy_score, classification_report,
confusion_matrix
from sklearn.tree import export_graphviz
```

- **pandas** is used for data manipulation and analysis.
- **RandomForestClassifier** from **sklearn.ensemble** is an ensemble learning method that uses multiple decision trees to make more accurate and robust predictions.
- **LabelEncoder** is used to convert categorical data into numeric format.
- **accuracy\_score**, **classification\_report**, and **confusion\_matrix** are metrics used for model evaluation.
- **export\_graphviz** helps visualize the decision trees within the random forest.

## 2. Load the Training Dataset

```
train_file_path = '../train.csv'
train_data = pd.read_csv(train_file_path)
```

- The **train.csv** file is read into a DataFrame named **train\_data**. This dataset contains features and a target variable (**Survived**) used to train the model.

## 3. Load the Testing Dataset

```
test_file_path = '../test.csv'
test_data = pd.read_csv(test_file_path)
```

- The **test.csv** file is read into a DataFrame named **test\_data**. This dataset includes features but not the **Survived** column, which the model predicts.

## 4. Load the Actual Results

```
gender_submission_file_path = '../gender_submission.csv'
actual_results = pd.read_csv(gender_submission_file_path)
```

- `gender_submission.csv` contains the actual `Survived` values for the test set, used for model evaluation.

## 5. Preprocess the Training Data

```
train_data['Age'].fillna(train_data['Age'].median(), inplace=True)
train_data['Cabin'].fillna('Unknown', inplace=True)
train_data['Embarked'].fillna(train_data['Embarked'].mode()[0], inplace=True)
```

- **Missing value handling:**
  - `Age`: Missing values are replaced with the median age.
  - `Cabin`: Missing cabin values are filled with 'Unknown'.
  - `Embarked`: Missing embarked values are filled with the most common port (mode).

## 6. Label Encoding for Categorical Variables

```
label_encoders = {}
for column in ['gender', 'Cabin', 'Embarked', 'Name', 'Ticket']:
    le = LabelEncoder()
    train_data[column] = le.fit_transform(train_data[column])
    label_encoders[column] = le
```

- Categorical columns (`gender`, `Cabin`, `Embarked`, `Name`, `Ticket`) are converted to numeric using `LabelEncoder`.
- A dictionary (`label_encoders`) stores the encoders for later use on the test data.

## 7. Select Features and Target Variable for Training

```
X_train = train_data[['PassengerId', 'Pclass', 'Name', 'gender', 'Age', 'SibSp',
                      'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked']]
y_train = train_data['Survived']
```

- `X_train`: Features used to train the model.
- `y_train`: The target variable indicating survival (1 for survived, 0 for not).

## 8. Train the Random Forest Classifier

```
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
rf_classifier.fit(X_train, y_train)
```

- A **Random Forest Classifier** with 100 decision trees (`n_estimators=100`) is created and trained using `fit()`.
- `random_state=42` ensures reproducibility by initializing the random number generator.

## 9. Preprocess the Test Data

```
test_data['Age'].fillna(test_data['Age'].median(), inplace=True)
test_data['Cabin'].fillna('Unknown', inplace=True)
test_data['Embarked'].fillna(test_data['Embarked'].mode()[0], inplace=True)
```

- The same preprocessing steps applied to `train_data` are applied to `test_data` to handle missing values.

## 10. Transform Test Data Using LabelEncoders

```
for column in ['gender', 'Cabin', 'Embarked', 'Name', 'Ticket']:
    if column in label_encoders:
        le = label_encoders[column]
        test_data[column] = test_data[column].apply(lambda x: le.transform([x])[0]
        if x in le.classes_ else -1)
```

- The test data columns are transformed using the `LabelEncoders` created earlier. If a value is not seen during training, it is encoded as `-1`.

## 11. Select Features for Test Data

```
X_test = test_data[['PassengerId', 'Pclass', 'Name', 'gender', 'Age', 'SibSp',
'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked']]
```

- `X_test` contains the features to be used for predictions.

## 12. Make Predictions on Test Data

```
y_pred = rf_classifier.predict(X_test)
```

- The trained model makes predictions on `X_test`.

## 13. Evaluate the Model

```
y_test = actual_results['Survived']

print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

- `y_test` contains the actual `Survived` values from `gender_submission.csv`.

- The model's accuracy, precision, recall, F1-score, and confusion matrix are printed.

## 14. Export Decision Trees for Visualization

```
export_graphviz(rf_classifier.estimators_[0],
                out_file="tree.dot",
                feature_names=X_train.columns,
                class_names=['0', '1'],
                filled=True)

export_graphviz(rf_classifier.estimators_[99],
                out_file="tree_99.dot",
                feature_names=X_train.columns,
                class_names=['0', '1'],
                filled=True)
```

- Two trees from the Random Forest (0 and 99) are exported as .dot files for visualization.
- **feature\_names** specify the feature labels for the nodes.
- **filled=True** colors the nodes based on the class they represent.

## Explanation of Random Forest Algorithm

- A **Random Forest** is an ensemble method that builds multiple decision trees using different subsets of the training data and features. The final output is the mode (classification) or average (regression) of all tree predictions.
- The main benefits include **reduced risk of overfitting** and **higher accuracy** due to aggregation.
- Randomness during training (sampling and feature selection) ensures **diverse trees**, which improves generalization.

## Why Random Forest?

- **Resistant to overfitting** compared to individual decision trees.
- **Works well with large datasets** and can handle both numerical and categorical data.
- **Feature importance** is inherently available for model insights.