

1. Import Necessary Libraries

In [26]:

```
# Import necessary libraries
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.preprocessing import LabelEncoder
# Import necessary libraries for evaluation
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.tree import export_graphviz
import pandas as pd
```

- pandas (pd) : Used for data manipulation and analysis, particularly for loading and working with data in DataFrame format.
- DecisionTreeClassifier : A machine learning algorithm from sklearn used to classify data by learning simple decision rules inferred from the features.
- LabelEncoder : A utility from sklearn to convert categorical string labels into numerical values.
- accuracy_score, classification_report, confusion_matrix : These are metrics to evaluate the performance of the classifier.
- export_graphviz : A function to export the decision tree structure for visualization.

2. Load the Training Dataset

In [27]:

```
# Load the training dataset
train_file_path = '../train.csv'
train_df = pd.read_csv(train_file_path)
```

- The training data is loaded from a CSV file using pandas.read_csv . The data is expected to be a CSV file, and the DataFrame train_df stores it.

3. Load the Testing Dataset

In [28]:

```
# Load the testing dataset
test_file_path = '../test.csv' # Update this path with your actual test file path if needed
test_df = pd.read_csv(test_file_path)
```

- Similarly, the testing dataset is loaded from a CSV file into test_df . This dataset will be used to test the trained model.

4. Data Preprocessing for Training Set

In [29]:

```
# Data preprocessing for training set
# Fill missing age values with the median
train_df['Age'].fillna(train_df['Age'].median(), inplace=True)
```

In [30]:

```
# Fill missing embarked values with the most common port
train_df['Embarked'].fillna(train_df['Embarked'].mode()[0], inplace=True)
```

In [31]:

```
# Drop 'Cabin' due to too many missing values
train_df.drop('Cabin', axis=1, inplace=True)
```

- train_df['Age'].fillna(train_df['Age'].median(), inplace=True) : Missing values in the 'Age' column are filled with the median age from the training data. This is done to handle missing values.
- train_df['Embarked'].fillna(train_df['Embarked'].mode()[0], inplace=True) : Missing values in the 'Embarked' column are replaced with the most common port (mode).
- train_df.drop('Cabin', axis=1, inplace=True) : The 'Cabin' column is dropped because it contains too many missing values, making it unsuitable for analysis.

5. Encode Categorical Variables

In [32]:

```
# Encode categorical variables ('Sex' and 'Embarked')
label_encoder = LabelEncoder()
train_df['Sex'] = label_encoder.fit_transform(train_df['Sex'])
train_df['Embarked'] = label_encoder.fit_transform(train_df['Embarked'])
```

- `LabelEncoder` : Converts categorical string values into numerical values. This is necessary because machine learning algorithms typically work with numerical data.
- `train_df['Sex']` and `train_df['Embarked']` are encoded to numeric values (0 or 1).

6. Define Features and Target Variable for Training

In [33]:

```
# Features and target for training
features = ['PassengerId', 'Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked']
X_train = train_df[features]
y_train = train_df['Survived']
```

- `features` : A list of the columns to be used as features (inputs) for the model. These columns represent characteristics of passengers that could influence their survival.
- `X_train` : Contains the feature values from the training data.
- `y_train` : The target variable, 'Survived', indicates whether the passenger survived (1) or not (0).

7. Initialize and Train the Decision Tree Classifier

In [34]:

```
# Initialize and train the Decision Tree Classifier
clf = DecisionTreeClassifier(random_state=42)
clf.fit(X_train, y_train)
```

Out[34]:

```
DecisionTreeClassifier
DecisionTreeClassifier(random_state=42)
(https://scikit-learn.org/1.4/modules/generated/sklearn.tree.DecisionTreeClassifier.html)
```

- `DecisionTreeClassifier` : The decision tree model is initialized. `random_state=42` ensures reproducibility of the model's results.
- `clf.fit(X_train, y_train)` : The classifier is trained on the feature set `X_train` and the target variable `y_train`.

8. Data Preprocessing for the Test Set

In [35]:

```
# Data preprocessing for the test set
# Fill missing age values with the median from the training set
test_df['Age'].fillna(train_df['Age'].median(), inplace=True)
```

In [36]:

```
# Fill missing embarked values with the most common port from the training set
test_df['Embarked'].fillna(train_df['Embarked'].mode()[0], inplace=True)
```

In [37]:

```
# Drop 'Cabin' as done with the training set
test_df.drop('Cabin', axis=1, inplace=True)
```

In [38]:

```
# Encode categorical variables in the test set
test_df['Sex'] = label_encoder.fit_transform(test_df['Sex'])
test_df['Embarked'] = label_encoder.fit_transform(test_df['Embarked'])
```

- The test set is preprocessed in a similar manner to the training set:
 - Missing values in 'Age' are filled with the median from the training set.
 - Missing values in 'Embarked' are filled with the mode from the training set.
 - The 'Cabin' column is dropped.
 - Categorical columns ('Sex' and 'Embarked') are encoded using the same `LabelEncoder` used for the training set to maintain consistency.

9. Select Features for the Test Set

In [39]:

```
# Select features for the test set
X_test = test_df[features]
```

- The features for the test set (`X_test`) are selected in the same way as for the training set.

10. Make Predictions on the Test Set

In [40]:

```
# Predict on the test set
y_pred = clf.predict(X_test)
```

- `y_pred` : The model makes predictions on the test data (`X_test`). These predictions are the model's estimations of whether each passenger survived or not.

11. Create a DataFrame for Output

In [41]:

```
# Create a DataFrame for the output
output = pd.DataFrame({
    'PassengerId': test_df['PassengerId'],
    'Survived': y_pred
})
```

- An output DataFrame is created, combining the `PassengerId` from the test set and the predicted survival labels (`y_pred`).

12. Display and Save the Prediction Results

In [42]:

```
# Display the prediction results
print(output)
```

	PassengerId	Survived
0	892	0
1	893	0
2	894	0
3	895	0
4	896	0
..
413	1305	0
414	1306	1
415	1307	0
416	1308	0
417	1309	1

[418 rows x 2 columns]

In [43]:

```
# Save the output to a CSV file
output.to_csv('titanic_predictions.csv', index=False)
```

- The `output` DataFrame is printed to the console and saved as a CSV file (`titanic_predictions.csv`) for later use.

13. Load the Actual Results for Evaluation

In [44]:

```
# Load the actual results from 'gender_submission.csv'
actual_results_file_path = '../gender_submission.csv' # Update this path if necessary
actual_results_df = pd.read_csv(actual_results_file_path)
```

- The actual results (gender_submission.csv) are loaded to compare against the model's predictions.

14. Align the Actual Results with the Predictions

In [45]:

```
# Ensure that the 'PassengerId' is used to align predictions and actual results
y_test = actual_results_df.set_index('PassengerId')['Survived'].reindex(X_test['PassengerId']).values
```

- y_test : The actual survival outcomes are retrieved using the PassengerId as the index and reindexed to align with the test set's predictions.

In [46]:

```
# Evaluate the model
print("Accuracy Score:", accuracy_score(y_test, y_pred))
```

Accuracy Score: 0.7990430622009569

In [47]:

```
# Generate the classification report
print("\nClassification Report:\n", classification_report(y_test, y_pred))
```

Classification Report:					
	precision	recall	f1-score	support	
0	0.82	0.88	0.85	266	
1	0.76	0.65	0.70	152	
accuracy			0.80	418	
macro avg	0.79	0.77	0.78	418	
weighted avg	0.80	0.80	0.80	418	

- accuracy_score : Computes the accuracy, i.e., the proportion of correct predictions out of all predictions.
- classification_report : Provides a detailed report with precision, recall, and F1-score for each class (0 or 1), along with the overall accuracy.

16. Generate the Decision Tree Visualization

In [48]:

```
# Export the decision tree to a .dot file for visualization
export_graphviz(clf,
                out_file="tree.dot",
                feature_names=features,
                class_names=['0', '1'],
                filled=True)

# Note: You can use tools like Graphviz or online viewers to visualize the 'tree.dot' file.
```

- export_graphviz : Exports the decision tree as a .dot file, which can be used with Graphviz to visualize the tree. The filled=True option colors the nodes based on the predicted class.

17. Print the Confusion Matrix

In [49]:

```
# Print the confusion matrix
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

Confusion Matrix:

```
[[235  31]
 [ 53  99]]
```

In [50]:

```
# Export the decision tree to a .dot file for visualization
export_graphviz(clf,
                out_file="tree.dot",
                feature_names=features,
                class_names=['0', '1'],
                filled=True)

# Note: You can use tools like Graphviz or online viewers to visualize the 'tree.dot' file.
```

- `confusion_matrix` : Computes a confusion matrix, which compares the actual vs predicted values. It shows how many true positives, true negatives, false positives, and false negatives the model produced.

Conclusion

This code demonstrates the process of training a decision tree model on a Titanic dataset to predict survival outcomes based on various passenger features. After training, the model's performance is evaluated using accuracy, classification reports, and confusion matrices, and the decision tree is exported for visualization.