

```
In [11]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import mean_squared_error, accuracy_score, classification_report, confusion_matrix

# Create a pipeline
pipe = Pipeline([
    ('preprocessor', OneHotEncoder()),
    ('classifier', LogisticRegression())
])

# Import the Titanic dataset
df = pd.read_csv("C:/Users/91797/Downloads/archive/Titanic-Dataset.csv")

In [12]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   column                non-null count  dtype
--  --
0   Survived              891 non-null  int64
1   Pclass                891 non-null  int64
2   Sex                   891 non-null  object
3   Name                  891 non-null  object
4   Sex                   891 non-null  object
5   Age                   714 non-null  float64
6   SibSp                 891 non-null  int64
7   Parch                891 non-null  int64
8   Ticket                891 non-null  object
9   Fare                  891 non-null  float64
10  Cabin                204 non-null  object
11  Embarked              889 non-null  object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB

In [13]: df.head()
Out[13]: bound method NDFrame.head of      PassengerId  Survived  Pclass \
0             1         0       3
1             2         1       1
2             3         1       3
3             4         1       1
4             5         0       3
..      ..
886            887         0       2
887            888         1       1
888            889         0       3
889            890         1       1
890            891         0       3

      Name                Sex  Age  SibSp  \
0  Braund, Mr. Owen Harris   male  22.0      1
1  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0      1
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)    female  35.0      1
4  Allen, Mr. William Henry   male  35.0      0
..      ..
886  Montvila, Rev. Juozas     male  27.0      0
887  Graham, Miss. Margaret Edith    female  19.0      0
888  Johnston, Miss. Catherine Helen "Carrie"    female  NA      1
889   Behr, Mr. Karl Howell     male  26.0      0
890  Dooley, Mr. Patrick       male  32.0      0

      Parch  Ticket  Fare  Cabin  Embarked
0         0    A/5 21171   7.2500   NaN      S
1         0  C 15999   51.2833   C      C
2         0  STON/O2, 3101282   7.9250   NaN      S
3         0    113803  53.1000  C123   S
4         0   373450   8.0500   NaN      S
..      ..
886         0   213536  13.0000   NaN      S
887         0   112853  30.0000   B42   S
888         2   W./T. 1607  23.0000   NaN      S
889         0   113569  30.0000  C148   C
890         0   378376   7.7500   NaN      Q

[891 rows x 12 columns]

In [14]: duplicate_rows = df[df.duplicated()]
print("Number of duplicate rows: {}".format(duplicate_rows.shape[0]))
print(duplicate_rows)

Number of duplicate rows: 0
Empty DataFrame
Columns: PassengerId, Survived, Pclass, Name, Sex, Age, SibSp, Parch, Ticket, Fare, Cabin, Embarked
Index: []

In [15]: df[df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)]
In [16]: df = df.drop(columns=['PassengerId', 'Name', 'Ticket', 'Fare', 'Cabin'])
print(df.head())

Survived  Pclass    Sex   Age  SibSp  Parch   Fare  Embarked
0         0       3   male  22.0      1      0   7.2500      S
1         1       1   female  38.0      1      0  51.2833      C
2         1       3   female  26.0      0      0   7.9250      S
3         1       1   female  35.0      1      0  53.1000      S
4         0       3   male  35.0      0      0   8.0500      S

In [17]: df.describe()
Out[17]:
   Survived   Pclass    Sex   Age  SibSp  Parch   Fare
count  891.000000  891.000000  714.000000  891.000000  891.000000  891.000000
mean    0.383838  2.308642  29.699118  1.023008  0.381594  32.204208
std    0.486592  0.836071  14.526497  1.102743  0.806057  49.693429
min     0.000000  1.000000  0.420000  0.000000  0.000000  0.000000
25%     0.000000  2.000000  20.125000  0.000000  0.000000  7.910400
50%     0.000000  3.000000  28.000000  0.000000  0.000000  14.454200
75%     1.000000  3.000000  38.000000  1.000000  0.000000  31.000000
max     1.000000  3.000000  80.000000  8.000000  6.000000  512.329200

In [18]: numerical_columns = ['Age', 'SibSp', 'Parch', 'Fare']
plt.figure(figsize=(20, 5))
plt.subplot(1, 4, 1)
for i, column in enumerate(numerical_columns, 1):
    sns.histplot(df[column], kde=True)
    plt.title(f'Histogram of {column}')
    plt.xlabel(column)
    plt.ylabel('Frequency')
plt.tight_layout()
plt.show()

C:\Users\91797\anaconda3\lib\site-packages\seaborn\oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
with pd.option_context('mode.use_inf_as_na', True):
C:\Users\91797\anaconda3\lib\site-packages\seaborn\oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
with pd.option_context('mode.use_inf_as_na', True):
C:\Users\91797\anaconda3\lib\site-packages\seaborn\oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
with pd.option_context('mode.use_inf_as_na', True):
C:\Users\91797\anaconda3\lib\site-packages\seaborn\oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
with pd.option_context('mode.use_inf_as_na', True):

Histogram of Age
Histogram of SibSp
Histogram of Parch
Histogram of Fare

In [19]: numerical_columns = ['Age', 'SibSp', 'Parch', 'Fare']
plt.figure(figsize=(20, 5))
for i, column in enumerate(numerical_columns, 1):
    plt.subplot(1, 4, i)
    sns.boxplot(df[column])
    plt.title(f'Boxplot of {column}')
plt.tight_layout()
plt.show()

Boxplot of Age
Boxplot of SibSp
Boxplot of Parch
Boxplot of Fare

In [20]: categorical_columns = ['Sex', 'Embarked', 'Survived', 'Pclass']
plt.figure(figsize=(20, 5)) # Adjust figure size as needed
for i, column in enumerate(categorical_columns, 1):
    plt.subplot(1, 4, i) # 1 row, 4 columns, i-th plot
    sns.countplot(df[column])
    ax = sns.countplot(df[column], ax=plt.gca()) # Use countplot
    ax = sns.countplot(df[column], ax=plt.gca()) # Store the axes object

    # Add labels to the bars
    for p in ax.patches:
        ax.annotate((p.get_height() * 100, p.get_x() + p.get_width() / 2., p.get_height()),
                    ha='center', va='center', fontsize=10, color='black', xytext=(0, 5),
                    textcoords='offset points')

    plt.title(f'Countplot of {column}')
    plt.xlabel(column)
    plt.ylabel('Frequency')
plt.tight_layout()
plt.show()

Countplot of Sex
Countplot of Embarked
Countplot of Survived
Countplot of Pclass

In [21]: bins_age = [0, 15, 60, np.inf]
labels = ['Child', 'Adult', 'Elderly']

# Create a new column 'AgeGroup' with binned age values
df['AgeGroup'] = pd.cut(df['Age'], bins=bins_age, labels=labels)
print(df[['Age', 'AgeGroup']].head())

Age  AgeGroup
0  22.0  Adult
1  38.0  Adult
2  26.0  Adult
3  35.0  Adult
4  35.0  Adult

In [22]: df['SibSp'] = df['SibSp'].clip(upper=4)
# Keep Parch at 4
df['Parch'] = df['Parch'].clip(upper=4)
# Keep Fare at 200
df['Fare'] = df['Fare'].clip(upper=200)
print(df[['SibSp', 'Parch', 'Fare']].head())

SibSp  Parch  Fare
0         1      0   7.2500
1         1      0  51.2833
2         0      0   7.9250
3         1      0  53.1000
4         0      0   8.0500

In [23]: df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
# Create 'FamilySize' column by binning the 'Fare' column into 4 equal-sized bins
df['FareGroup'] = pd.cut(df['Fare'], q=4, labels=['Low', 'Mid', 'High', 'Very High'])
print(df[['SibSp', 'Parch', 'FamilySize', 'FareGroup']].head())

SibSp  Parch  FamilySize  FareGroup
0         1      0         2      Low
1         1      0         2  Very High
2         0      0         1      Mid
3         1      0         2  Very High
4         0      0         1      Mid

In [24]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 11 columns):
#   column                non-null count  dtype
--  --
0   Survived              891 non-null  int64
1   Pclass                891 non-null  int64
2   Sex                   891 non-null  object
3   Age                   714 non-null  float64
4   SibSp                 891 non-null  int64
5   Parch                891 non-null  int64
6   Fare                  891 non-null  float64
7   Embarked              891 non-null  object
8   AgeGroup              891 non-null  category
9   FamilySize            891 non-null  int64
10  FareGroup              891 non-null  category
dtypes: category(2), float64(2), int64(5), object(2)
memory usage: 64.8+ KB

In [25]: df.describe()
Out[25]:
   Survived   Pclass    Sex   Age  SibSp  Parch   Fare  FamilySize
count  891.000000  891.000000  714.000000  891.000000  891.000000  891.000000
mean    0.383838  2.308642  29.699118  1.023008  0.381594  32.204208
std    0.486592  0.836071  14.526497  1.102743  0.806057  49.693429
min     0.000000  1.000000  0.420000  0.000000  0.000000  0.000000  1.000000
25%     0.000000  2.000000  20.125000  0.000000  0.000000  7.910400  1.000000
50%     0.000000  3.000000  28.000000  0.000000  0.000000  14.454200  1.000000
75%     1.000000  3.000000  38.000000  1.000000  0.000000  31.000000  2.000000
max     1.000000  3.000000  80.000000  8.000000  6.000000  512.329200  7.000000

In [26]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 11 columns):
#   column                non-null count  dtype
--  --
0   Survived              891 non-null  int64
1   Pclass                891 non-null  int64
2   Sex                   891 non-null  object
3   Age                   714 non-null  float64
4   SibSp                 891 non-null  int64
5   Parch                891 non-null  int64
6   Fare                  891 non-null  float64
7   Embarked              891 non-null  object
8   AgeGroup              714 non-null  category
9   FamilySize            891 non-null  int64
10  FareGroup              891 non-null  category
dtypes: category(2), float64(2), int64(5), object(2)
memory usage: 64.8+ KB

In [27]: numerical_features = ['Fare', 'FamilySize']
categorical_features = ['Pclass', 'Sex', 'Embarked', 'AgeGroup', 'FareGroup']
print('Numerical Features:', numerical_features)
print('Categorical Features:', categorical_features)

Numerical features: ['Fare', 'FamilySize']
Categorical features: ['Pclass', 'Sex', 'Embarked', 'AgeGroup', 'FareGroup']

In [28]: numerical_features = ['Fare', 'FamilySize']
categorical_features = ['Pclass', 'Sex', 'Embarked', 'AgeGroup', 'FareGroup']
target = 'Survived'

# Split the data into training and testing sets
X = df[numerical_features + categorical_features]
y = df[target]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Preprocess the data
preprocessor = ColumnTransformer(
    transformers=[
        ('num', Pipeline([('scaler', StandardScaler())]), ('passthrough', 'passthrough'))], numerical_features),
        ('cat', OneHotEncoder(handle_unknown='ignore')), categorical_features)
])

# Pipeline with SMOTE and Random Forest ---
pipeline_rf_smote = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('smote', SMOTE(random_state=42)), # Apply SMOTE for oversampling
    ('classifier', RandomForestClassifier(random_state=42)) # Random Forest Classifier
])

# Train the model ---
pipeline_rf_smote.fit(X_train, y_train)

# Make Predictions ---
y_pred = pipeline_rf_smote.predict(X_test)

# Evaluate the Model ---
print("Random Forest Classifier:")
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy}')
print(classification_report(y_test, y_pred))

cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", cm)

Logistic Regression:
Accuracy: 0.78709602726739
precision    recall  f1-score   support

   0    0.87   0.75   0.81    165
   1    0.78   0.84   0.77    174

accuracy    0.79   0.80   0.79    179
macro avg   0.79   0.80   0.79    179
weighted avg 0.80   0.79   0.79    179

Confusion Matrix:
[[79 26]
 [24 42]]

In [29]: X = df[numerical_features + categorical_features]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Preprocess the data
preprocessor = ColumnTransformer(
    transformers=[
        ('num', Pipeline([('scaler', StandardScaler())]), ('passthrough', 'passthrough'))], numerical_features),
        ('cat', OneHotEncoder(handle_unknown='ignore')), categorical_features)
])

# Pipeline with SMOTE and Random Forest ---
pipeline_rf_smote = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('smote', SMOTE(random_state=42)), # Apply SMOTE for oversampling
    ('classifier', RandomForestClassifier(random_state=42)) # Random Forest Classifier
])

# Train the model ---
pipeline_rf_smote.fit(X_train, y_train)

# Make Predictions ---
y_pred = pipeline_rf_smote.predict(X_test)

# Evaluate the Model ---
print("Random Forest Classifier:")
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy}')
print(classification_report(y_test, y_pred))

cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", cm)

Random Forest Classifier:
Accuracy: 0.82129602726739
precision    recall  f1-score   support

   0    0.84   0.86   0.85    165
   1    0.79   0.77   0.78    174

accuracy    0.82   0.81   0.82    179
macro avg   0.82   0.82   0.82    179
weighted avg 0.82   0.82   0.82    179

Confusion Matrix:
[[80 15]
 [14 42]]

In [30]: pipeline_gbm_smote = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('smote', SMOTE(random_state=42)), # Apply SMOTE
    ('classifier', XGBClassifier(random_state=42)) # XGB classifier
])

# Train the model ---
pipeline_gbm_smote.fit(X_train, y_train)

# Make Predictions ---
y_pred = pipeline_gbm_smote.predict(X_test)

# Evaluate the Model ---
print("XGBM:")
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy}')
print(classification_report(y_test, y_pred))

cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", cm)

GBM:
Accuracy: 0.8156420180537
precision    recall  f1-score   support

   0    0.85   0.84   0.84    165
   1    0.77   0.78   0.78    174

accuracy    0.82   0.81   0.82    179
macro avg   0.81   0.81   0.81    179
weighted avg 0.82   0.82   0.82    179

Confusion Matrix:
[[88 17]
 [16 42]]

In [31]: pipeline_gb_smote = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('smote', SMOTE(random_state=42)), # Use the same preprocessor
    ('classifier', GradientBoostingClassifier(random_state=42))
])

pipeline_gb_smote.fit(X_train, y_train)
accuracy_gb = accuracy_score(y_test, y_pred_gb)
print("Gradient Boosting:")
print(f'Accuracy: {accuracy_gb}')
print(classification_report(y_test, y_pred_gb))
cm_gb = confusion_matrix(y_test, y_pred_gb)
print("Confusion Matrix:\n", cm_gb)

Gradient Boosting:
Accuracy: 0.80162811713844
precision    recall  f1-score   support

   0    0.88   0.86   0.87    165
   1    0.81   0.84   0.82    174

accuracy    0.84   0.85   0.85    179
macro avg   0.84   0.85   0.85    179
weighted avg 0.85   0.85   0.85    179

Confusion Matrix:
[[88 17]
 [16 42]]

In [32]: pipeline_gb_smote = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('smote', SMOTE(random_state=42)), # Use the same preprocessor
    ('classifier', XGBClassifier(random_state=42), max_iter=500, early_stopping=True)) # Increased max_iter, added early_stopping
])

pipeline_gb_smote.fit(X_train, y_train)
y_pred_gb = pipeline_gb_smote.predict(X_test)

accuracy_gb = accuracy_score(y_test, y_pred_gb)
print("XGBM:")
print(f'Accuracy: {accuracy_gb}')
print(classification_report(y_test, y_pred_gb))

cm_gb = confusion_matrix(y_test, y_pred_gb)
print("Confusion Matrix:\n", cm_gb)

XGBM:
Accuracy: 0.79765163128491
precision    recall  f1-score   support

   0    0.85   0.71   0.78    165
   1    0.67   0.82   0.74    174

accuracy    0.76   0.77   0.76    179
macro avg   0.76   0.77   0.76    179
weighted avg 0.78   0.76   0.76    179

Confusion Matrix:
[[79 38]
 [13 41]]

In [33]: param_grid = {
    'classifier__n_estimators': [50, 100], # Number of trees
    'classifier__max_depth': [None, 10], # Maximum depth of trees
    'classifier__min_samples_split': [2, 5], # Minimum samples to split a node
    'classifier__min_samples_leaf': [1, 3], # Minimum samples in a leaf node
    'classifier__criterion': ['gini', 'entropy'] # Splitting criteria
}

grid_search = GridSearchCV(pipeline_rf_smote, param_grid, cv=5, scoring='accuracy', n_jobs=-1, verbose=1) # 5-fold CV, use all cores

grid_search.fit(X_train, y_train)

# Best Model and Best Hyperparameters ---
best_model = grid_search.best_estimator_
best_params = grid_search.best_params_
print("Best Hyperparameters: ", best_params)

# Evaluate the Best Model on the Test Set ---
y_pred = best_model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy (Best Model): {accuracy}')
print(classification_report(y_test, y_pred))

cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", cm)

Cross-Validation Score (Optional) ---
cv_scores = cross_val_score(best_model, X_train, y_train, cv=5, scoring='accuracy', n_jobs=-1)
print(f'Cross-Validation Score (Best Model): {cv_scores.mean():.3f}')
print(f'Mean CV Accuracy (Best Model): {cv_scores.mean():.3f}')

Fitting 5 folds for each of 32 candidates, totalling 160 fits
Best Hyperparameters: ('classifier__criterion', 'gini', 'classifier__max_depth', 10, 'classifier__min_samples_split', 5, 'classifier__n_estimators', 50)
Accuracy (Best Model): 0.8132276174526955
precision    recall  f1-score   support

   0    0.83   0.89   0.86    165
   1    0.82   0.74   0.78    174

accuracy    0.83   0.83   0.83    179
macro avg   0.83   0.81   0.82    179
weighted avg 0.83   0.83   0.83    179

Confusion Matrix:
[[93 12]
 [10 41]]

Cross-Validation Score (Best Model): [0.8841958 0.8841958 0.8169141 0.7887329 0.85211268]
Mean CV Accuracy (Best Model): 0.8132276174526955

In [34]: df.head()
Out[34]:
   Survived  Pclass    Sex   Age  SibSp  Parch   Fare  Embarked  AgeGroup  FamilySize  FareGroup
0         0       3   male  22.0      1      0   7.2500      S      Adult         2      Low
1         1       1   female  38.0      1      0  51.2833      C      Adult         2  Very High
2         1       3   female  26.0      0      0   7.9250      S      Adult         1      Mid
3         1       1   female  35.0      1      0  53.1000      S      Adult         2  Very High
4         0       3   male  35.0      0      0   8.0500      S      Adult         1      Mid
```



