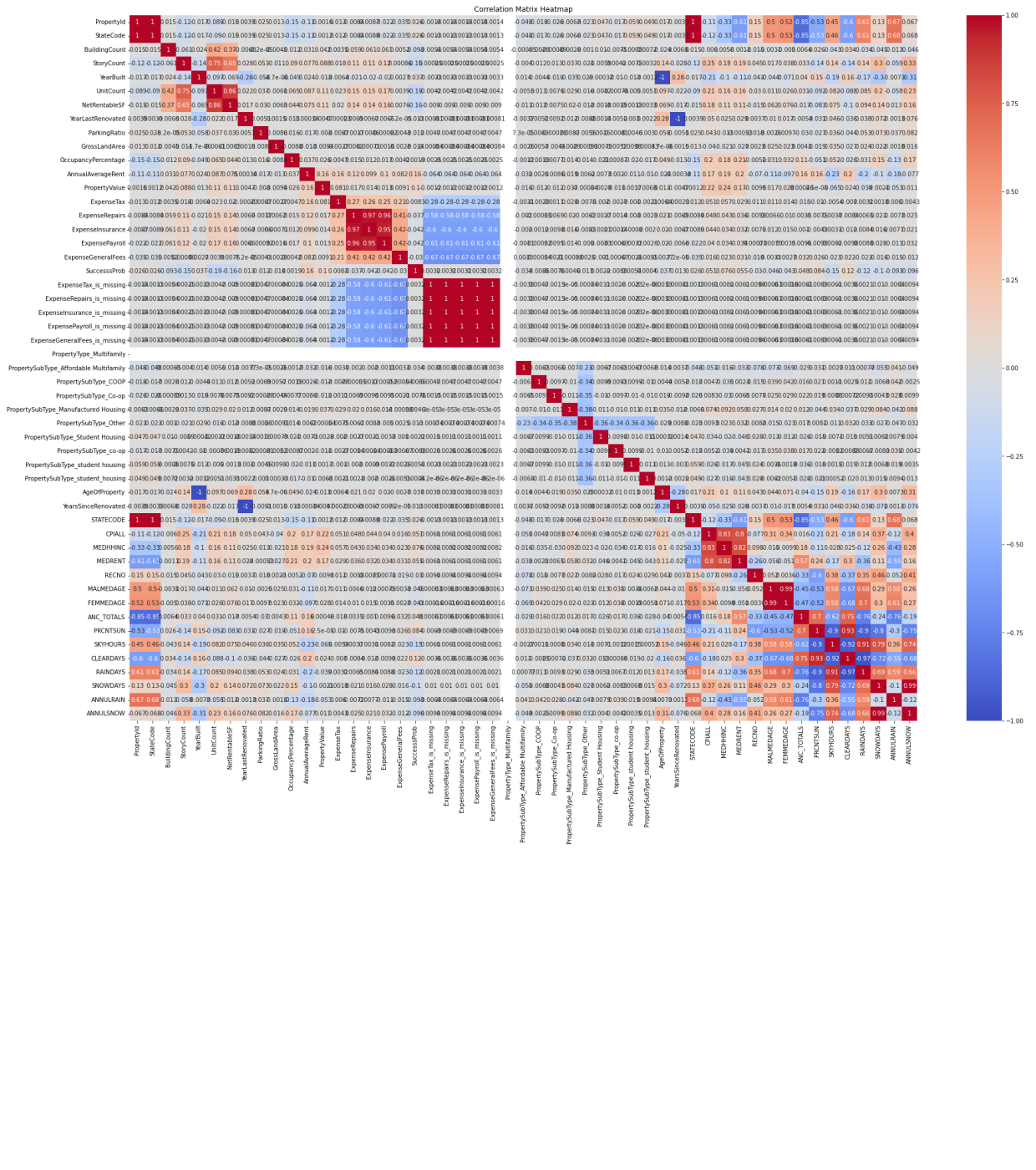
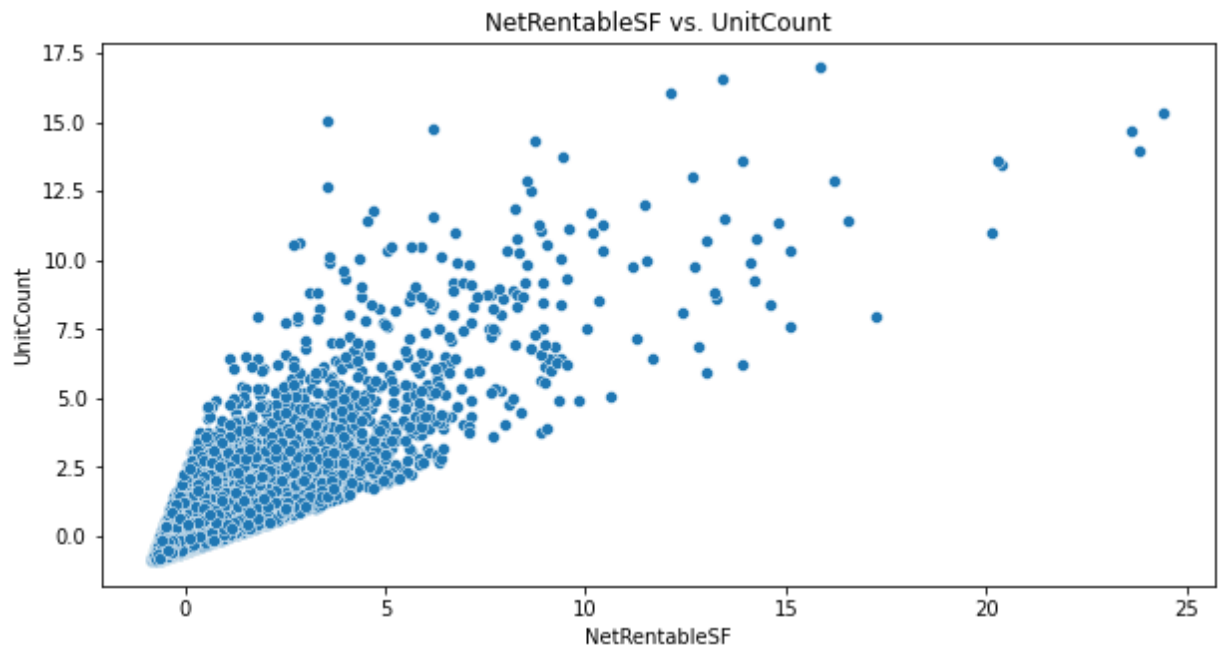


```
3 plt.title('Correlation Matrix Heatmap')
4 plt.show()
5
```



```
In [238]: 1 plt.figure(figsize=(10, 5))
2 sns.scatterplot(x='NetRentableSF', y='UnitCount', data=df)
3 plt.title('NetRentableSF vs. UnitCount')
4 plt.show()
5
```



```
In [239]: 1 #Model building
2 from sklearn.model_selection import train_test_split
3
4 # Define your features and target variable
5 X = df.drop('SuccesssProb', axis=1) # assuming 'SuccesssProb' is your targe
6 y = df['SuccesssProb']
7
8 # Convert 'SuccesssProb' into a binary variable
9 y = (df['SuccesssProb'] > 0.5).astype(int)
10
11 # Split the data into a training set and a test set
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, ran
13
14 # Create a Logistic regression model
15 model = LogisticRegression()
16
17 # Train the model on the training data
18 model.fit(X_train, y_train)
19
20
```

Out[239]: LogisticRegression()

```
In [240]: 1 from sklearn.linear_model import LogisticRegression
2
3 # Create a logistic regression model
4 model_lr = LogisticRegression()
5
6 # Train the model on the training data
7 model_lr.fit(X_train, y_train)
8
```

Out[240]: LogisticRegression()

```
In [241]: 1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f
2
3 # Make predictions on the test set
4 y_pred = model_lr.predict(X_test)
5
6 # Calculate metrics
7 accuracy = accuracy_score(y_test, y_pred)
8 precision = precision_score(y_test, y_pred)
9 recall = recall_score(y_test, y_pred)
10 f1 = f1_score(y_test, y_pred)
11 roc_auc = roc_auc_score(y_test, y_pred)
12
13 # Print metrics
14 print(f'Accuracy: {accuracy}')
15 print(f'Precision: {precision}')
16 print(f'Recall: {recall}')
17 print(f'F1 Score: {f1}')
18 print(f'ROC AUC Score: {roc_auc}')
19
```

Accuracy: 0.7747813411078717

Precision: 0.7751460767946577

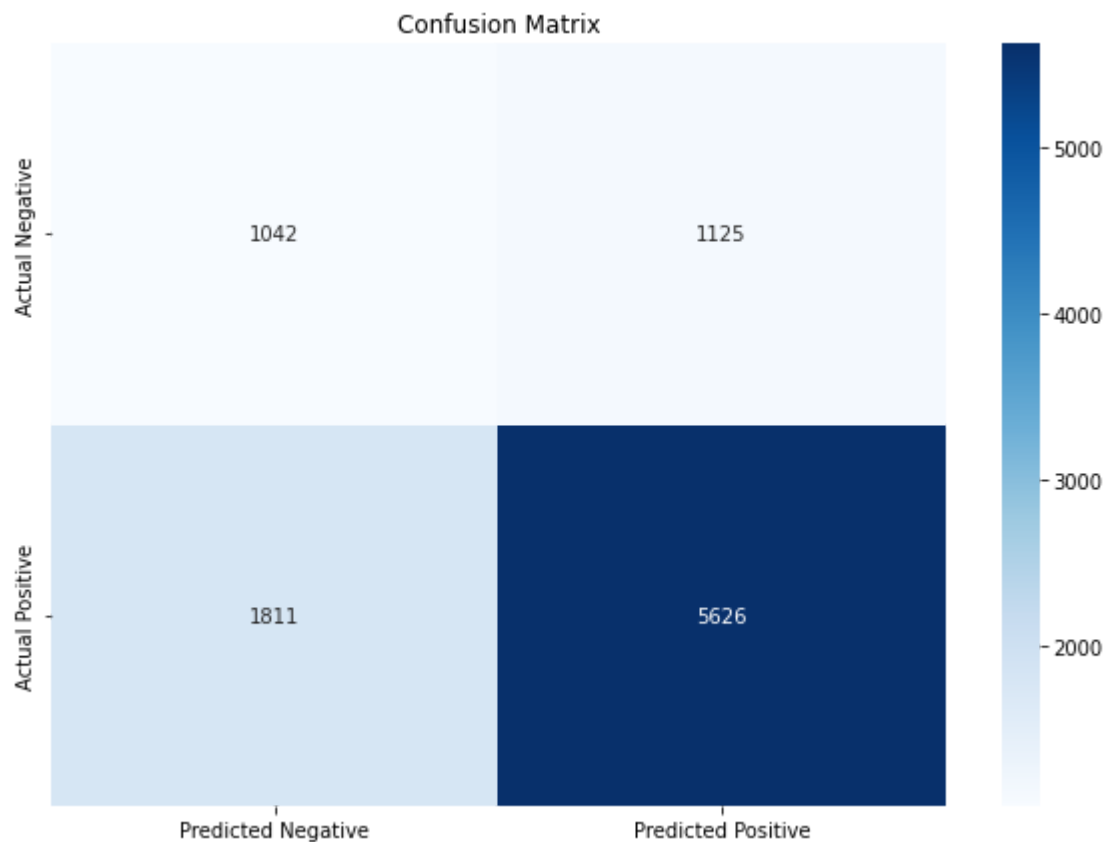
Recall: 0.9989242974317601

F1 Score: 0.8729216849773809

ROC AUC Score: 0.5022309535151417

In [242]:

```
1 from sklearn.metrics import confusion_matrix
2
3 # Make probability predictions on the test set
4 y_prob = model.predict_proba(X_test)[: , 1]
5
6 # Classify properties as successful or unsuccessful based on the optimal dec
7 y_pred = (y_prob > min_threshold).astype(int)
8
9 # Create a confusion matrix
10 cm = confusion_matrix(y_test, y_pred)
11
12 # Create a DataFrame from the confusion matrix for better visualization
13 cm_df = pd.DataFrame(cm, columns=['Predicted Negative', 'Predicted Positive']
14
15 plt.figure(figsize=(10, 7))
16
17 sns.heatmap(cm_df, annot=True, fmt='d', cmap='Blues')
18 plt.title('Confusion Matrix')
19 plt.show()
20
```



In [243]:

```
1  #cost benifit analysis
2  import numpy as np
3
4  # Define the costs and benefits
5  benefit = 1000000
6  cost = 3000000
7
8  # Calculate the probabilities of the positive class
9  y_prob = model.predict_proba(X_test)[:, 1]
10
11 # Calculate the expected cost for different decision thresholds
12 thresholds = np.linspace(0, 1, 100)
13 expected_costs = []
14 for threshold in thresholds:
15     y_pred = (y_prob > threshold).astype(int)
16     fp = np.sum((y_test == 0) & (y_pred == 1))
17     tp = np.sum((y_test == 1) & (y_pred == 1))
18     expected_cost = fp * cost - tp * benefit
19     expected_costs.append(expected_cost)
20
21 # Find the threshold with the minimum expected cost
22 min_cost = min(expected_costs)
23 min_threshold = thresholds[expected_costs.index(min_cost)]
24
25 print(f'Minimum Expected Cost: {min_cost}')
26 print(f'Optimal Decision Threshold: {min_threshold}')
27
```

Minimum Expected Cost: -2251000000

Optimal Decision Threshold: 0.7575757575757577

In [244]:

```
1  #A negative value indicates that the benefits outweigh the costs.
```

The threshold value of approximately 0.76. is the optimal decision threshold for classifying a property as successful or unsuccessful, based on minimizing the expected cost.

This means that if the model predicts a success probability of greater than 0.76 for a property, we should classify it as successful; otherwise, you should classify it as unsuccessful.

Generally, the accuracy is 50%, but based on our conditions of 1:3 ratio of profit to loss, we ended up at 0.76 as threshold

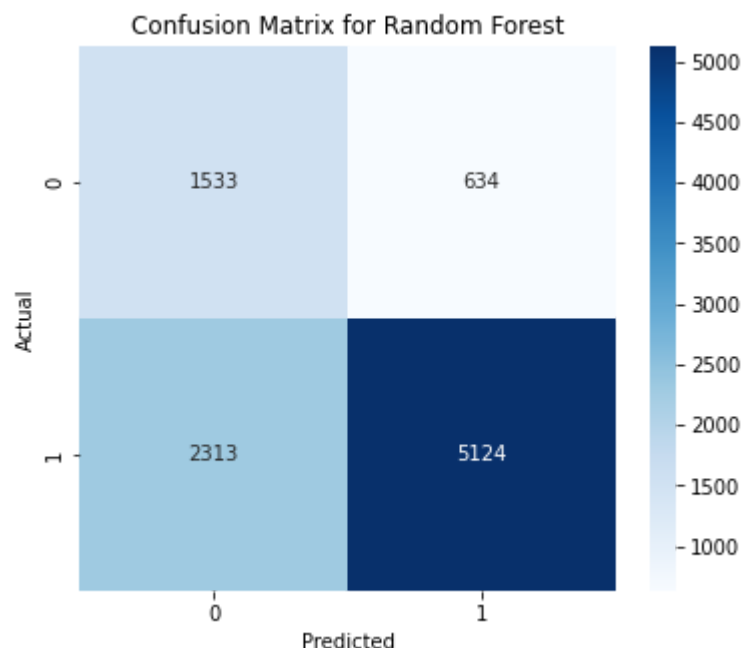
In [245]:

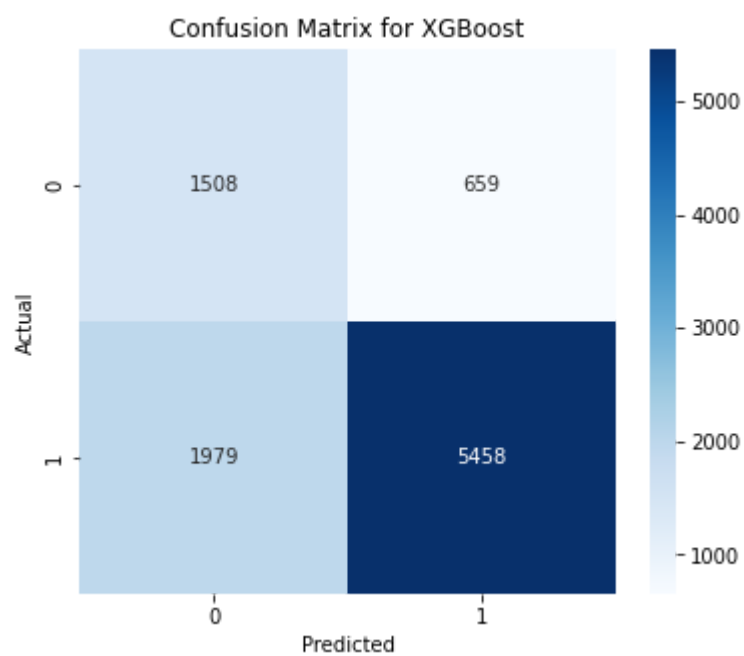
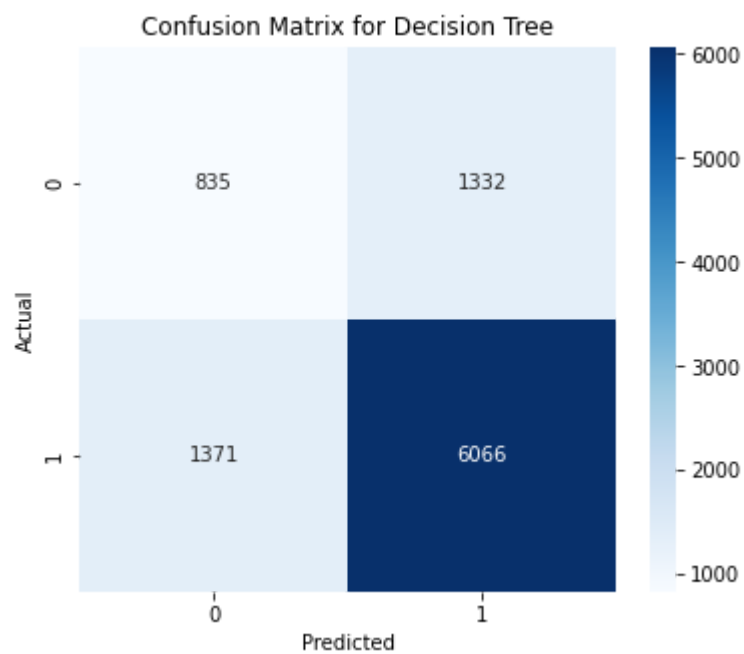
```
1  from sklearn.ensemble import RandomForestClassifier
2
3  # Initialize the model
4  rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
5
6  # Train the model
7  rf_model.fit(X_train, y_train)
8
9  # Make predictions
10 y_pred_rf = rf_model.predict(X_test)
11
12
13 from sklearn.tree import DecisionTreeClassifier
14
15 # Initialize the model
16 dt_model = DecisionTreeClassifier(random_state=42)
17
18 # Train the model
19 dt_model.fit(X_train, y_train)
20
21 # Make predictions
22 y_pred_dt = dt_model.predict(X_test)
23
24
25 import xgboost as xgb
26
27 # Initialize the model
28 xgb_model = xgb.XGBClassifier(random_state=42)
29
30 # Train the model
31 xgb_model.fit(X_train, y_train)
32
33 # Make predictions
34 y_pred_xgb = xgb_model.predict(X_test)
35
```

[14:49:57] WARNING: ../src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

In [246]:

```
1 from sklearn.metrics import confusion_matrix
2
3 # Make probability predictions on the test set
4 y_prob_rf = rf_model.predict_proba(X_test)[: , 1]
5 y_prob_dt = dt_model.predict_proba(X_test)[: , 1]
6 y_prob_xgb = xgb_model.predict_proba(X_test)[: , 1]
7
8 # Classify properties as successful or unsuccessful based on the optimal dec
9 y_pred_rf = (y_prob_rf > min_threshold).astype(int)
10 y_pred_dt = (y_prob_dt > min_threshold).astype(int)
11 y_pred_xgb = (y_prob_xgb > min_threshold).astype(int)
12
13 # Create confusion matrices
14 cm_rf = confusion_matrix(y_test, y_pred_rf)
15 cm_dt = confusion_matrix(y_test, y_pred_dt)
16 cm_xgb = confusion_matrix(y_test, y_pred_xgb)
17
18 import seaborn as sns
19 import matplotlib.pyplot as plt
20
21 # Create a list of models and their confusion matrices
22 models = ['Random Forest', 'Decision Tree', 'XGBoost']
23 cms = [cm_rf, cm_dt, cm_xgb]
24
25 # Plot the confusion matrix for each model
26 for model, cm in zip(models, cms):
27     plt.figure(figsize=(6, 5))
28     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
29     plt.title(f'Confusion Matrix for {model}')
30     plt.xlabel('Predicted')
31     plt.ylabel('Actual')
32     plt.show()
33
```









In [247]:

```
1  from sklearn.metrics import accuracy_score, precision_score, recall_score
2
3
4
5  # Calculate metrics for Random Forest
6  tn_rf, fp_rf, fn_rf, tp_rf = cm_rf.ravel()
7  specificity_rf = tn_rf / (tn_rf + fp_rf)
8
9  # Calculate metrics for Decision Tree
10 tn_dt, fp_dt, fn_dt, tp_dt = cm_dt.ravel()
11 specificity_dt = tn_dt / (tn_dt + fp_dt)
12
13 # Calculate metrics for XGBoost
14 tn_xgb, fp_xgb, fn_xgb, tp_xgb = cm_xgb.ravel()
15 specificity_xgb = tn_xgb / (tn_xgb + fp_xgb)
16
17
18
19 # Calculate metrics for Random Forest
20 accuracy_rf = accuracy_score(y_test, y_pred_rf)
21 precision_rf = precision_score(y_test, y_pred_rf)
22 recall_rf = recall_score(y_test, y_pred_rf)
23
24 # Calculate metrics for Decision Tree
25 accuracy_dt = accuracy_score(y_test, y_pred_dt)
26 precision_dt = precision_score(y_test, y_pred_dt)
27 recall_dt = recall_score(y_test, y_pred_dt)
28
29 # Calculate metrics for XGBoost
30 accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
31 precision_xgb = precision_score(y_test, y_pred_xgb)
32 recall_xgb = recall_score(y_test, y_pred_xgb)
33
34
35
36 # Print the metrics for each model
37 print('\nMetrics for Random Forest:')
38 print(f'Accuracy: {accuracy_rf:.4f}')
39 print(f'Precision: {precision_rf:.4f}')
40 print(f'Recall (Sensitivity): {recall_rf:.4f}')
41 print(f'Specificity: {specificity_rf:.4f}')
42
43 print('\nMetrics for Decision Tree:')
44 print(f'Accuracy: {accuracy_dt:.4f}')
45 print(f'Precision: {precision_dt:.4f}')
46 print(f'Recall (Sensitivity): {recall_dt:.4f}')
47 print(f'Specificity: {specificity_dt:.4f}')
48
49 print('\nMetrics for XGBoost:')
50 print(f'Accuracy: {accuracy_xgb:.4f}')
51 print(f'Precision: {precision_xgb:.4f}')
52 print(f'Recall (Sensitivity): {recall_xgb:.4f}')
53 print(f'Specificity: {specificity_xgb:.4f}')
```

Metrics for Random Forest:

Accuracy: 0.6931

Precision: 0.8899

Recall (Sensitivity): 0.6890

Specificity: 0.7074

Metrics for Decision Tree:

Accuracy: 0.7186

Precision: 0.8200

Recall (Sensitivity): 0.8157

Specificity: 0.3853

Metrics for XGBoost:

Accuracy: 0.7253

Precision: 0.8923

Recall (Sensitivity): 0.7339

Specificity: 0.6959

In [248]:

```
1 model_rf = rf_model
2 model_dt = dt_model
3 model_xgb = xgb_model
```



In [249]:

```
1  #Probability thresholds for various models:
2
3  import numpy as np
4
5  # Define the costs and benefits
6  benefit = 1000000
7  cost = 3000000
8
9  # Calculate the probabilities of the positive class for each model
10
11 y_prob_rf = model_rf.predict_proba(X_test)[: , 1]
12 y_prob_dt = model_dt.predict_proba(X_test)[: , 1]
13 y_prob_xgb = model_xgb.predict_proba(X_test)[: , 1]
14
15 # Calculate the expected cost for different decision thresholds for each model
16 thresholds = np.linspace(0, 1, 100)
17
18
19
20 # Random Forest
21 expected_costs_rf = []
22 for threshold in thresholds:
23     y_pred_rf = (y_prob_rf > threshold).astype(int)
24     fp_rf = np.sum((y_test == 0) & (y_pred_rf == 1))
25     tp_rf = np.sum((y_test == 1) & (y_pred_rf == 1))
26     expected_cost_rf = fp_rf * cost - tp_rf * benefit
27     expected_costs_rf.append(expected_cost_rf)
28
29 # Decision Tree
30 expected_costs_dt = []
31 for threshold in thresholds:
32     y_pred_dt = (y_prob_dt > threshold).astype(int)
33     fp_dt = np.sum((y_test == 0) & (y_pred_dt == 1))
34     tp_dt = np.sum((y_test == 1) & (y_pred_dt == 1))
35     expected_cost_dt = fp_dt * cost - tp_dt * benefit
36     expected_costs_dt.append(expected_cost_dt)
37
38 # XGBoost
39 expected_costs_xgb = []
40 for threshold in thresholds:
41     y_pred_xgb = (y_prob_xgb > threshold).astype(int)
42     fp_xgb = np.sum((y_test == 0) & (y_pred_xgb == 1))
43     tp_xgb = np.sum((y_test == 1) & (y_pred_xgb == 1))
44     expected_cost_xgb = fp_xgb * cost - tp_xgb * benefit
45     expected_costs_xgb.append(expected_cost_xgb)
46
47 # Find the threshold with the minimum expected cost for each model
48
49
50 min_cost_rf = min(expected_costs_rf)
51 min_threshold_rf = thresholds[expected_costs_rf.index(min_cost_rf)]
52
53 min_cost_dt = min(expected_costs_dt)
54 min_threshold_dt = thresholds[expected_costs_dt.index(min_cost_dt)]
55
56 min_cost_xgb = min(expected_costs_xgb)
57 min_threshold_xgb = thresholds[expected_costs_xgb.index(min_cost_xgb)]
```

```

58
59 # Print the minimum expected costs and optimal decision thresholds for each
60
61 print('Random Forest:')
62 print(f'Minimum Expected Cost: {min_cost_rf}')
63 print(f'Optimal Decision Threshold: {min_threshold_rf}\n')
64
65 print('Decision Tree:')
66 print(f'Minimum Expected Cost: {min_cost_dt}')
67 print(f'Optimal Decision Threshold: {min_threshold_dt}\n')
68
69 print('XGBoost:')
70 print(f'Minimum Expected Cost: {min_cost_xgb}')
71 print(f'Optimal Decision Threshold: {min_threshold_xgb}\n')
72
73

```

Random Forest:  
Minimum Expected Cost: -3229000000  
Optimal Decision Threshold: 0.7474747474747475

Decision Tree:  
Minimum Expected Cost: -2070000000  
Optimal Decision Threshold: 0.0

XGBoost:  
Minimum Expected Cost: -3509000000  
Optimal Decision Threshold: 0.7777777777777778



## Conclusion

In [253]:

```

1 import pandas as pd
2
3 # Create a dictionary with the model metrics
4 model_metrics = {
5     'Model': ['Logistic Regression', 'Random Forest', 'Decision Tree', 'XGBo
6     'Accuracy': [accuracy, accuracy_rf, accuracy_dt, accuracy_xgb],
7     'Optimal Threshold': [min_threshold, min_threshold_rf, min_threshold_dt,
8 }
9
10 # Create a DataFrame from the dictionary
11 metrics_table = pd.DataFrame(model_metrics)
12
13 # Print the table
14 print(metrics_table)
15

```

	Model	Accuracy	Optimal Threshold
0	Logistic Regression	0.774781	0.757576
1	Random Forest	0.693149	0.747475
2	Decision Tree	0.718555	0.000000
3	XGBoost	0.725323	0.777778



Based on the EDA and the modeling contingent on the cost benigit contriants provided in the Question

Here are the best models:

Logistic Regression: 77% accuracy, threshold is 75%

XGboost: 72% accuracy , threshold is 77%

In [ ]:

1

In [ ]:

1