In [2]: ►

```python
 1  import numpy as np
 2  import pandas as pd
 3  import matplotlib.pyplot as plt
 4  %matplotlib inline
 5  import seaborn as sns
 6  from sklearn.utils.class_weight import compute_class_weight
 7  from sklearn.preprocessing import StandardScaler
 8  from sklearn.linear_model import LogisticRegression
 9  from sklearn.tree import DecisionTreeClassifier
10  from sklearn.ensemble import RandomForestClassifier
11  from sklearn.model_selection import train_test_split, GridSearchCV, cr
12  from sklearn.metrics import accuracy_score, recall_score, precision_sc
13  from sklearn.metrics import ConfusionMatrixDisplay
14  from sklearn.metrics import classification_report
15  from sklearn.pipeline import Pipeline
16  from imblearn.pipeline import Pipeline as ImbPipeline
17  from sklearn.decomposition import PCA
18  from imblearn.over_sampling import SMOTE, BorderlineSMOTE
19  from google.colab import files
20  uploaded = files.upload()
```

Choose Files | No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

```
Saving pivoted_df.csv to pivoted_df.csv
```

Load the dataframe in, inspect the data.

In [5]: ►

```python
 1  pivoted_df = pd.read_csv('~/Documents/Flatiron/Project_5_/data/pivoted
```

In [6]: ►

```python
 1  pivoted_df.head()
```

Out[6]:

| | season | Age | Throws | Surgery | AB_release_speed_weighted_avg | CH_release_speed_weigh |
|---|---|---|---|---|---|---|
| 0 | 2008 | 37.0 | 1 | 0.0 | 0.0 | 8: |
| 1 | 2009 | 38.0 | 1 | 0.0 | 0.0 | 8! |
| 2 | 2010 | 39.0 | 1 | 0.0 | 0.0 | 8 |
| 3 | 2011 | 40.0 | 1 | 0.0 | 0.0 | 8: |
| 4 | 2012 | 41.0 | 1 | 0.0 | 0.0 | 8: |

5 rows × 130 columns

In [7]: ►

```python
 1  pivoted_df.shape
```

Out[7]: (3688, 130)

In [8]:  ▶|  | 1  `pivoted_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3688 entries, 0 to 3687
Columns: 130 entries, season to SV_vz0_weighted_avg
dtypes: float64(128), int64(2)
memory usage: 3.7 MB
```

In [9]:  ▶|  | 1  `pivoted_df['Surgery'].value_counts()`

Out[9]:
```
0.0    2772
1.0     916
Name: Surgery, dtype: int64
```

Time to start modeling! Split target and features and make a baseline model.

In [10]:  ▶|  | 1  `y = pivoted_df['Surgery']`
         | 2  `X = pivoted_df.drop('Surgery', axis=1)`

In [11]:  ▶|  | 1  `X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.`

In [12]:

```python
#Make a pipeline to simplify process
logreg_pipeline = Pipeline([
    ('scale', StandardScaler()),
    ('logreg', LogisticRegression(solver='liblinear'))
])

# Define parameter grid to search
param_grid = {
    'logreg__C': [0.001, 0.01, 0.1, 1, 10, 100],  # Regularization str
    'logreg__penalty': ['l1', 'l2']  # Norm used in the penalization
}

# Initialize GridSearchCV with pipeline, parameter grid, and scoring m
grid_search = GridSearchCV(logreg_pipeline, param_grid, cv=5, scoring=

# Assuming X_train and y_train are already defined
grid_search.fit(X_train, y_train)

# Best parameters found
print("Best parameters: ", grid_search.best_params_)

# Best cross-validation score
print("Best cross-validation score: {:.2f}".format(grid_search.best_sc

# Test set score using the best parameters
print("Test set score: {:.2f}".format(grid_search.score(X_test, y_test
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: Conver
genceWarning: Liblinear failed to converge, increase the number of iterat
ions.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: Conver
genceWarning: Liblinear failed to converge, increase the number of iterat
ions.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: Conver
genceWarning: Liblinear failed to converge, increase the number of iterat
ions.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: Conver
genceWarning: Liblinear failed to converge, increase the number of iterat
ions.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: Conver
genceWarning: Liblinear failed to converge, increase the number of iterat
ions.
  warnings.warn(

Best parameters:  {'logreg__C': 10, 'logreg__penalty': 'l1'}
Best cross-validation score: 0.76
Test set score: 0.77
```

In [13]: ▶|
```
1  logreg_pipeline = Pipeline([
2      ('scale', StandardScaler()),
3      ('logreg', LogisticRegression(penalty='l1', C=10.0, solver='liblin
4  ])
```

In [14]: ▶|
```
1  logreg_pipeline.fit(X_train, y_train)
```

Out[14]: Pipeline(steps=[('scale', StandardScaler()),
                         ('logreg',
                          LogisticRegression(C=10.0, penalty='l1', solver='libline
ar'))])

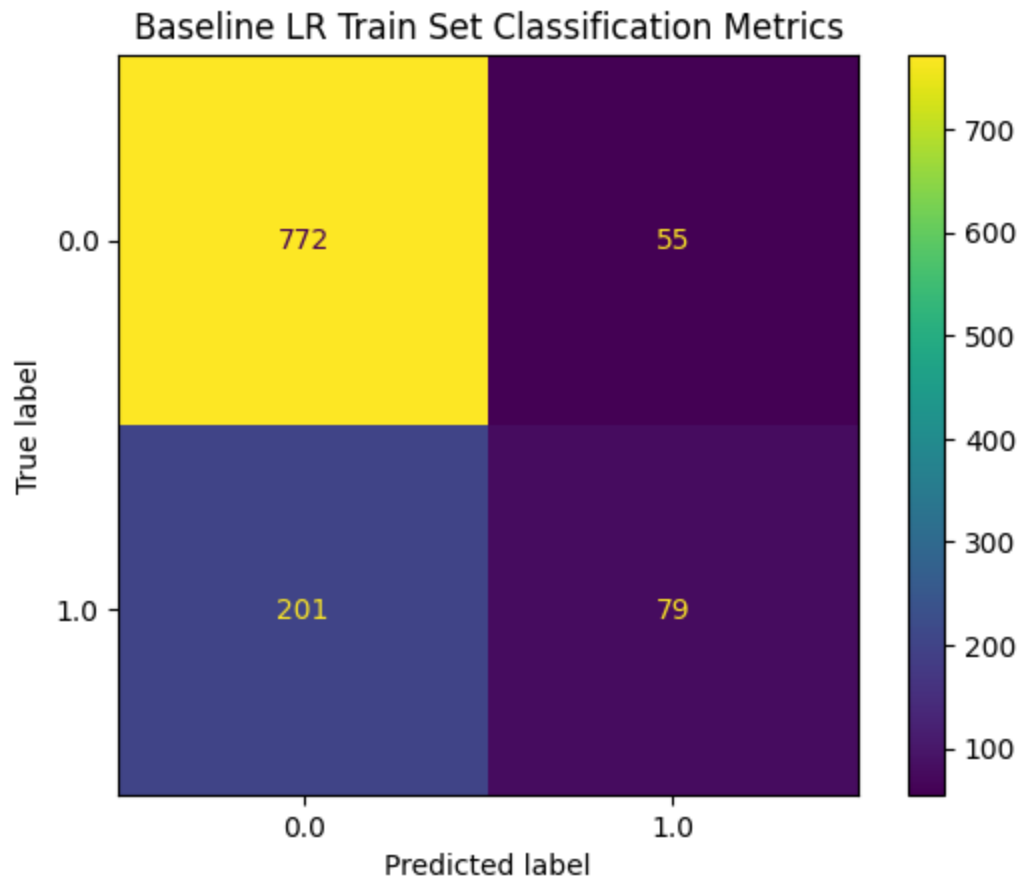**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

In [15]: ▶|
```
1  logreg_pipeline.score(X_test, y_test)
```

Out[15]: 0.7687443541102078

In [16]: ▶|
```
1  y_pred = logreg_pipeline.predict(X_test)
```

In [17]:

```python
1  ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2  plt.title('Baseline LR Train Set Classification Metrics')
3  plt.show()
4  print(classification_report(y_test, y_pred))
```

## Baseline LR Train Set Classification Metrics



```
                precision    recall  f1-score   support

         0.0        0.79      0.93      0.86       827
         1.0        0.59      0.28      0.38       280

    accuracy                            0.77      1107
   macro avg        0.69      0.61      0.62      1107
weighted avg        0.74      0.77      0.74      1107
```

Dataset is imbalanced, need to adjust. Should also focus on Recall score since this is a medical issue (better to have False Positive than True Negative!)

In [18]:

```python
1  y = pivoted_df['Surgery']
2  X = pivoted_df.drop('Surgery', axis=1)
```

In [19]:

```python
1  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
```

In [24]:
```python
1  # Set up pipeline
2  weight_logreg_pipeline = Pipeline([
3      ('scale', StandardScaler()),
4      ('logreg', LogisticRegression(solver='liblinear'))
5  ])
6
7  # Define the parameter grid to search over, including class weights
8  param_grid = {
9      'logreg__C': [0.01, 0.1, 1, 10],
10     'logreg__penalty': ['l1', 'l2'],
11     'logreg__class_weight': [None, 'balanced', {0: 1, 1: 2}, {0: 1, 1:
12     'logreg__max_iter': [5000],
13     'logreg__tol': [0.01]
14 }
15
16 # Create a scoring function that focuses on recall for the positive cl
17 recall_scorer = make_scorer(recall_score, pos_label=1)
18
19 # Initialize GridSearch with pipeline, param grid, and recall
20 grid_search = GridSearchCV(weight_logreg_pipeline, param_grid, cv=5, s
21
22 # Fit the grid search to the data
23 grid_search.fit(X_train, y_train)
24
25 # Print the best parameters found and the best recall score
26 print("Best parameters: ", grid_search.best_params_)
27 print("Best cross-validation recall score: {:.2f}".format(grid_search.
28
29 # Evaluate the best model on the test set
30 best_model = grid_search.best_estimator_
31 y_pred = best_model.predict(X_test)
32 print("Test set recall score: {:.2f}".format(recall_score(y_test, y_pr
```

```
Best parameters:  {'logreg__C': 0.01, 'logreg__class_weight': {0: 1, 1:
5}, 'logreg__max_iter': 5000, 'logreg__penalty': 'l1', 'logreg__tol': 0.0
1}
Best cross-validation recall score: 0.87
Test set recall score: 0.86
```

In [25]:
```python
1  best_model.fit(X_train, y_train)
```

Out[25]:
```
Pipeline(steps=[('scale', StandardScaler()),
                ('logreg',
                 LogisticRegression(C=0.01, class_weight={0: 1, 1: 5},
                                    max_iter=5000, penalty='l1',
                                    solver='liblinear', tol=0.01))])
```

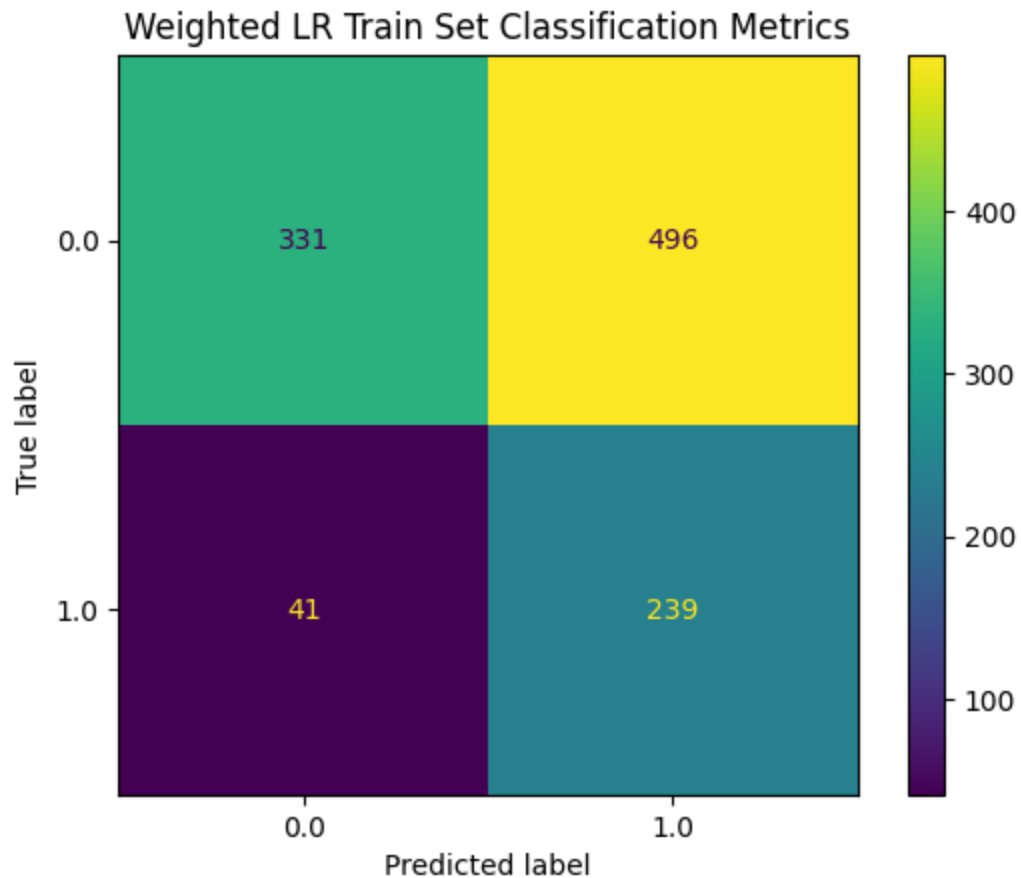**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [26]:    ▶    1  best_model.score(X_test, y_test)
```

Out[26]: 0.5149051490514905

```
In [28]:    ▶    1  y_pred = best_model.predict(X_test)
```

```
In [29]:    ▶    1  ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
                  2  plt.title('Weighted LR Train Set Classification Metrics')
                  3  plt.show()
                  4  print(classification_report(y_test, y_pred))
```

### Weighted LR Train Set Classification Metrics



```
              precision    recall  f1-score   support

         0.0       0.89      0.40      0.55       827
         1.0       0.33      0.85      0.47       280

    accuracy                           0.51      1107
   macro avg       0.61      0.63      0.51      1107
weighted avg       0.75      0.51      0.53      1107
```

Much better model. False Negatives is low, other classes much higher.

Decision Tree Classifier, baseline model.

In [30]: ▶|  
```
1  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
```

In [31]: ▶|  
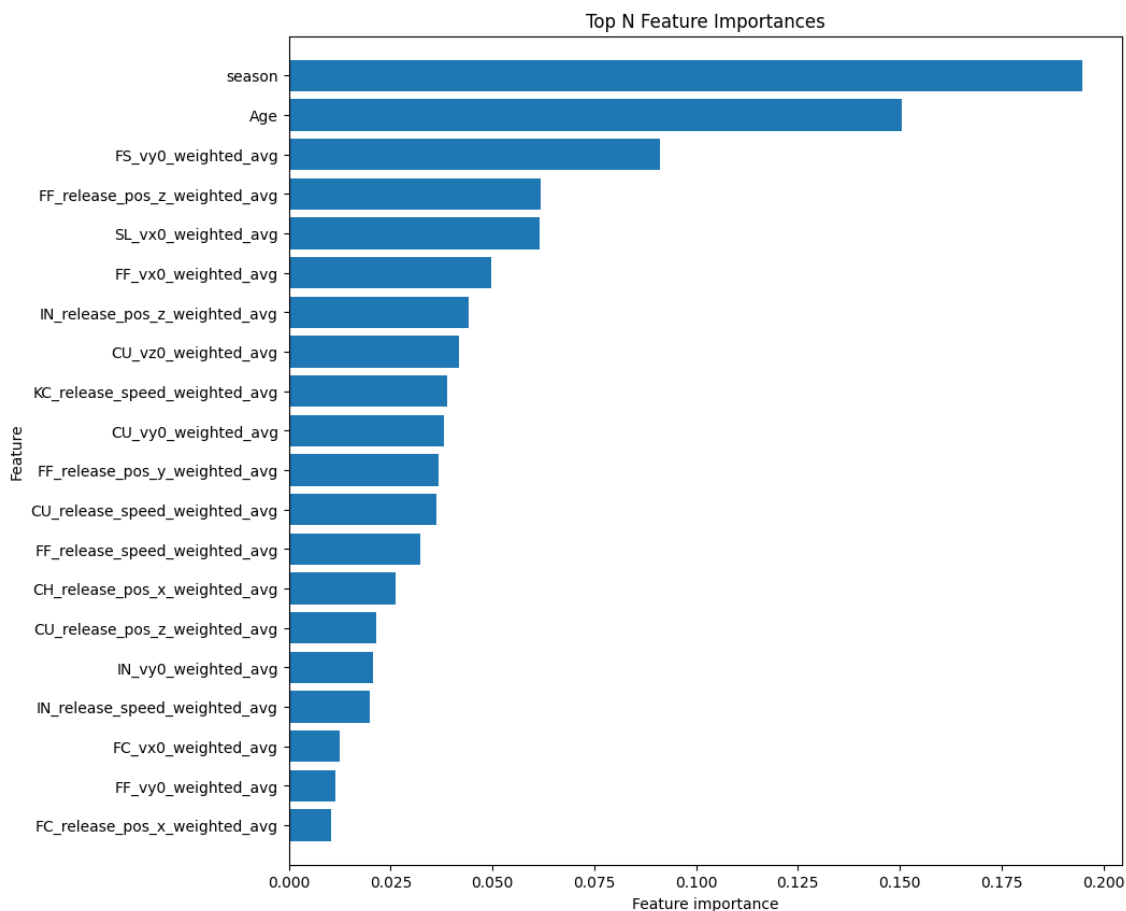```
1  tree_clf = DecisionTreeClassifier(criterion='gini', max_depth=5)
2  tree_clf.fit(X_train, y_train)
```

Out[31]: DecisionTreeClassifier(max_depth=5)

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**
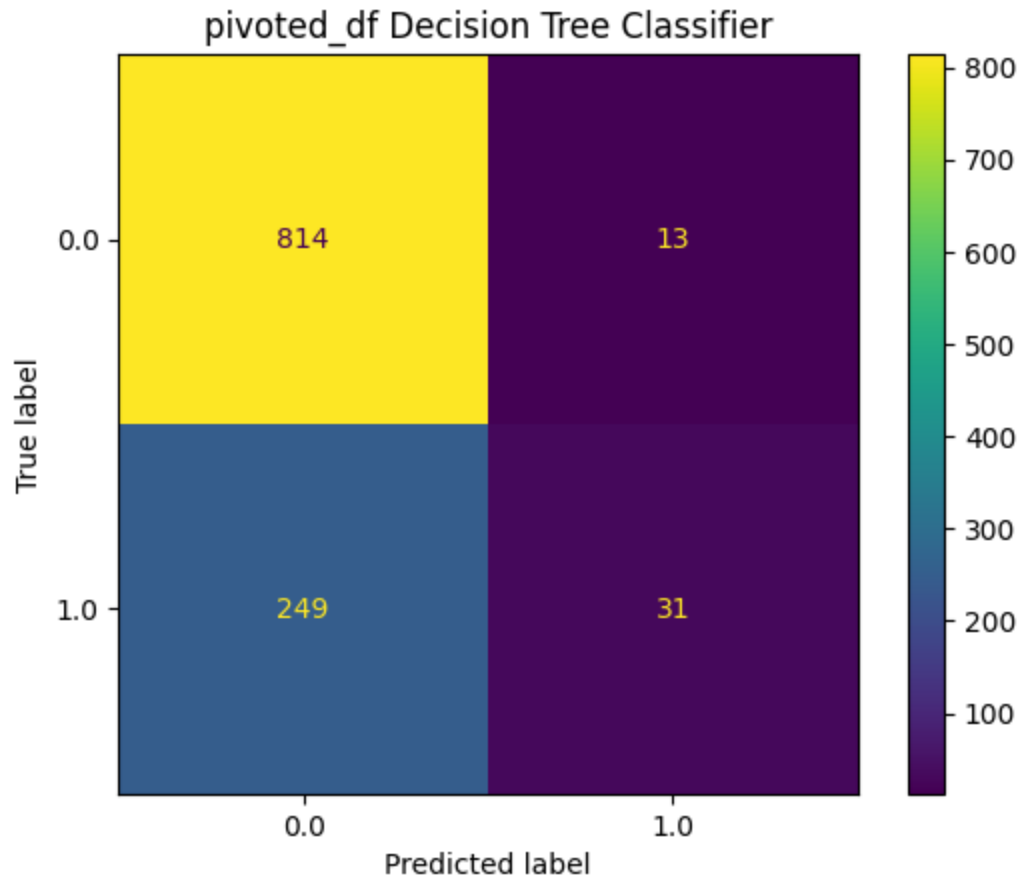
In [32]: ▶|  
```
1  def plot_feature_importances(model, n_top_features=20):
2      importances = model.feature_importances_
3      indices = np.argsort(importances)[-n_top_features:]
4      plt.figure(figsize=(10,10))
5      plt.title('Top N Feature Importances')
6      plt.barh(range(n_top_features), importances[indices], align='cente
7      plt.yticks(range(n_top_features), [X_train.columns[i] for i in ind
8      plt.xlabel('Feature importance')
9      plt.ylabel('Feature')
10     plt.ylim(-1, n_top_features)
11
12 plot_feature_importances(tree_clf, n_top_features=20)
13 plt.show()
```

In [33]: ▶|
```python
1  pred = tree_clf.predict(X_test)
```

In [34]: ▶|
```python
1  pred = tree_clf.predict(X_test)
2
3  ConfusionMatrixDisplay.from_predictions(y_test, pred)
4  plt.title('pivoted_df Decision Tree Classifier')
5  plt.show()
6  print(classification_report(y_test, pred))
```



```
              precision    recall  f1-score   support

         0.0       0.77      0.98      0.86       827
         1.0       0.70      0.11      0.19       280

    accuracy                           0.76      1107
   macro avg       0.74      0.55      0.53      1107
weighted avg       0.75      0.76      0.69      1107
```

Terrible for TP and FP. Need to adjust. Features are interesting. Mostly fastball, curveball, some slider and split-finger.

In [35]: ▶|
```python
1  y = pivoted_df['Surgery']
2  X = pivoted_df.drop('Surgery', axis=1)
```

In [36]: ▶|   1  `X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.`

In [37]: ▶|
```python
 1  param_grid = {
 2      'criterion': ['gini', 'entropy'],
 3      'max_depth': [5, 10, 15, 20],
 4      'min_samples_split': [2, 5, 10],
 5      'min_samples_leaf': [1, 2, 4],
 6      'class_weight': ['balanced', {0:1, 1:2}, {0:1, 1:3}]
 7  }
 8
 9  tree_clf = DecisionTreeClassifier()
10  scorer = make_scorer(recall_score)
11  grid_search = GridSearchCV(estimator=tree_clf, param_grid=param_grid,
12  grid_search.fit(X_train, y_train)
13
14  print("Best parameters:", grid_search.best_params_)
15  print("Best score:", grid_search.best_score_)
16
17  best_tree = grid_search.best_estimator_
18  y_pred = best_tree.predict(X_test)
19  print("Test recall score:", recall_score(y_test, y_pred))
```

```
Best parameters: {'class_weight': 'balanced', 'criterion': 'gini', 'max_d
epth': 5, 'min_samples_leaf': 2, 'min_samples_split': 10}
Best score: 0.6524852362204725
Test recall score: 0.7607142857142857
```

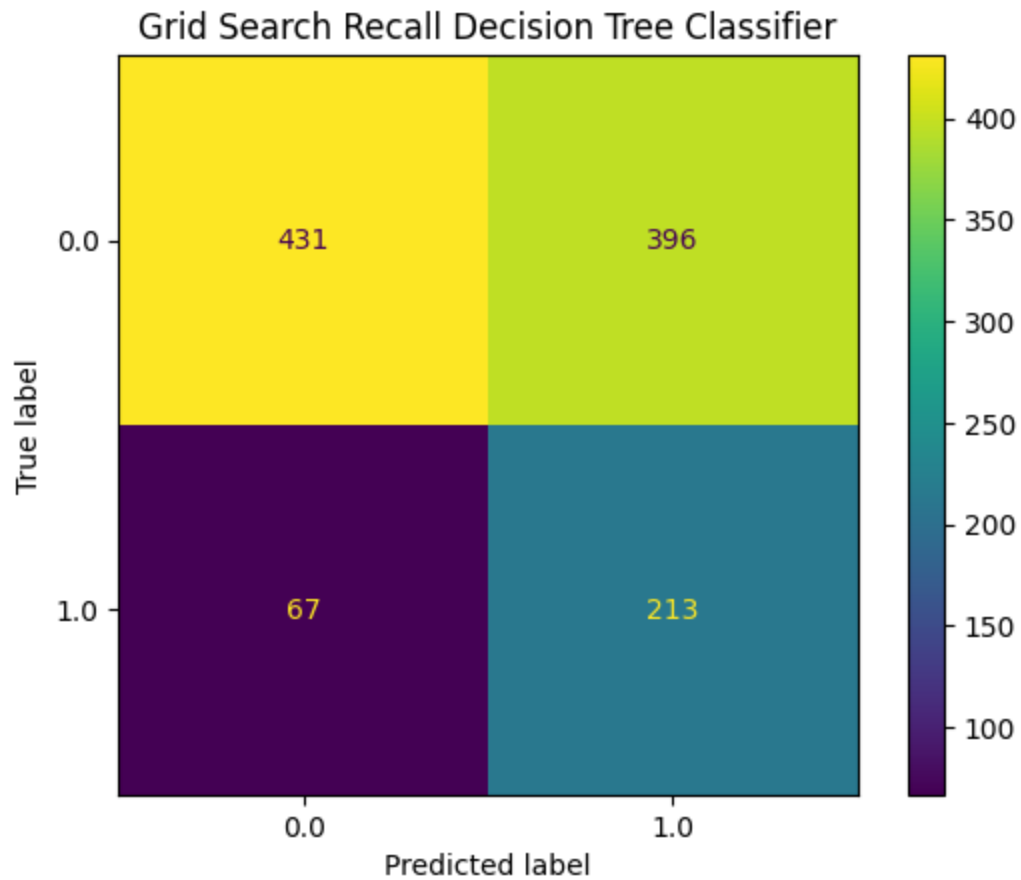In [38]: ▶|   1  `best_tree.fit(X_train, y_train)`

Out[38]: DecisionTreeClassifier(class_weight='balanced', max_depth=5, min_samples_
leaf=2,
                              min_samples_split=10)

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

In [39]: ▶|   1  `pred = best_tree.predict(X_test)`

In [40]: ►|
```python
1  ConfusionMatrixDisplay.from_predictions(y_test, pred)
2  plt.title('Grid Search Recall Decision Tree Classifier')
3  plt.show()
4  print(classification_report(y_test, pred))
```



Grid Search Recall Decision Tree Classifier

```
                precision    recall  f1-score   support

         0.0       0.87      0.52      0.65       827
         1.0       0.35      0.76      0.48       280

    accuracy                           0.58      1107
   macro avg       0.61      0.64      0.56      1107
weighted avg       0.74      0.58      0.61      1107
```

The Logistic Regression model with adjusted class weights performed the best.

In [ ]: ►|
```
1
```