

In [50]:

```
1 import numpy as np
2 import math
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 import seaborn as sns
7 import glob
8 import os
9 import io
10 import re
11 import unicodedata
12 from datetime import datetime
13 from itertools import groupby
14 from operator import itemgetter
15 from sklearn.utils.class_weight import compute_class_weight
16 from sklearn.preprocessing import OneHotEncoder, StandardScaler
17 from sklearn.linear_model import LogisticRegression
18 from sklearn.tree import DecisionTreeClassifier
19 from sklearn.ensemble import RandomForestClassifier
20 from sklearn.neighbors import KNeighborsClassifier
21 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
22 from sklearn.metrics import accuracy_score, recall_score, precision_score, f1_score
23 from sklearn.metrics import ConfusionMatrixDisplay
24 from sklearn.metrics import classification_report
25 from sklearn.pipeline import Pipeline
26 from imblearn.pipeline import Pipeline as ImbPipeline
27 from sklearn.decomposition import PCA
28 from imblearn.over_sampling import SMOTE, BorderlineSMOTE
29 from xgboost import XGBClassifier
30 from google.colab import files
31 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.

pivoted\_df, Logistic Regression Model, baseline

In [6]:

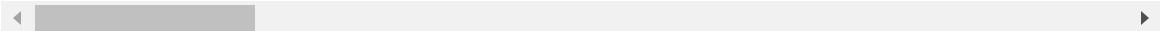
```
1 pivoted_df = pd.read_csv(io.BytesIO(uploaded['~/Documents/Flatiron/Project/flatiron-project-dataset.csv']))
```

In [3]: 1 pivoted\_df.head()

Out[3]:

	season	Age	Throws	Surgery	AB_release_speed_weighted_avg	CH_release_speed_weighted_avg
0	2008	37.0	1	0.0	0.0	8.1
1	2009	38.0	1	0.0	0.0	8.1
2	2010	39.0	1	0.0	0.0	8.1
3	2011	40.0	1	0.0	0.0	8.1
4	2012	41.0	1	0.0	0.0	8.1

5 rows × 130 columns



In [4]: 1 pivoted\_df['Surgery'].value\_counts()

Out[4]:

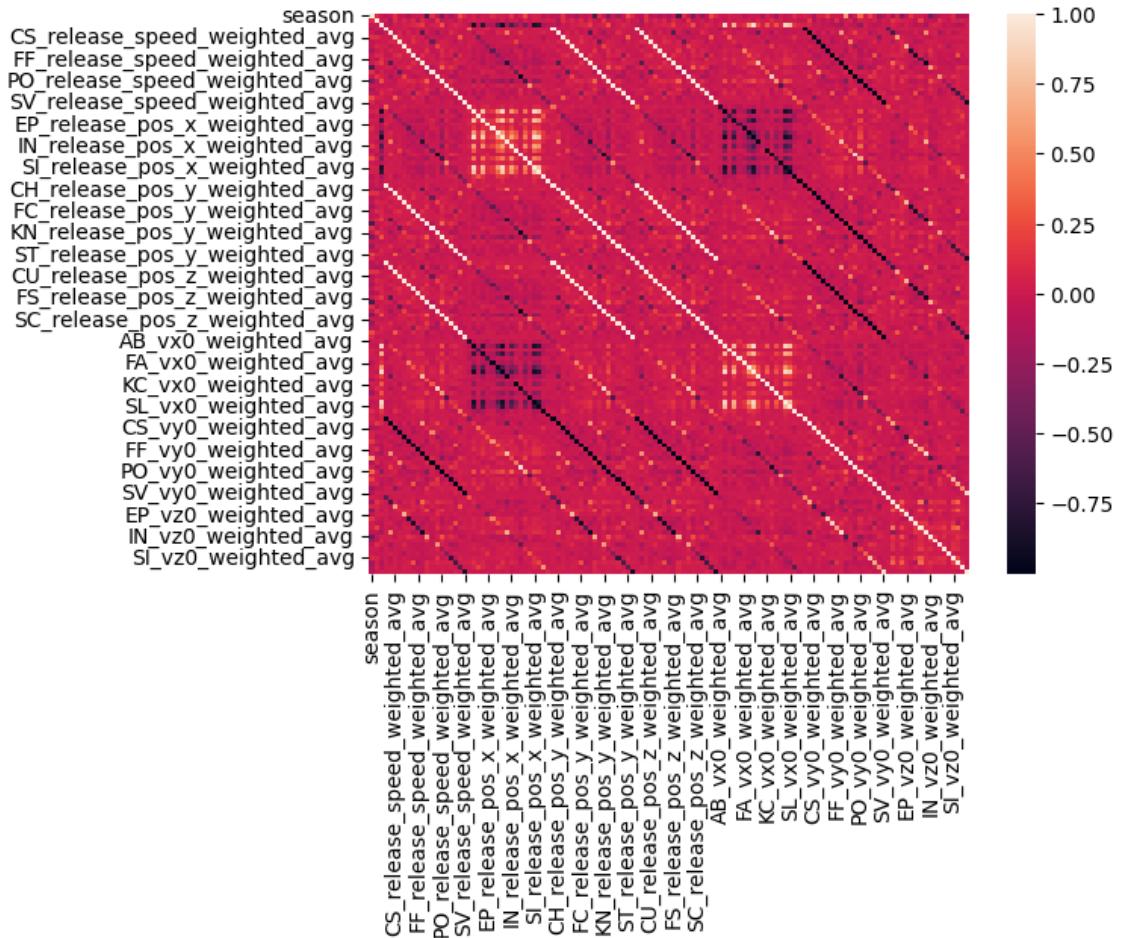
0.0	2772
1.0	916

Name: Surgery, dtype: int64

Can't read this heatmap.

```
In [110]: 1 corr = X_train.corr()
2 sns.heatmap(corr)
```

Out[110]: <Axes: >



```
In [14]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [16]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [17]: ❶ 1 logreg_pipeline = Pipeline([
2     ('scale', StandardScaler()),
3     ('logreg', LogisticRegression(solver='liblinear'))
4 ])
5
6 # Define the parameter grid to search over
7 param_grid = {
8     'logreg_C': [0.001, 0.01, 0.1, 1, 10, 100], # Regularization strength
9     'logreg_penalty': ['l1', 'l2'] # Norm used in the penalization
10 }
11
12 # Initialize GridSearchCV with the pipeline, parameter grid, and desired cross-validation folds
13 grid_search = GridSearchCV(logreg_pipeline, param_grid, cv=5, scoring='accuracy')
14
15 # Assuming X_train and y_train are already defined
16 grid_search.fit(X_train, y_train)
17
18 # Best parameters found
19 print("Best parameters: ", grid_search.best_params_)
20
21 # Best cross-validation score
22 print("Best cross-validation score: {:.2f}".format(grid_search.best_score_))
23
24 # Test set score using the best parameters
25 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: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/svm/_base.py:1244: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
    warnings.warn(
Best parameters: {'logreg_C': 10, 'logreg_penalty': 'l1'}
Best cross-validation score: 0.76
Test set score: 0.77
```

```
In [18]: ❶ 1 logreg_pipeline = Pipeline([
2      ('scale', StandardScaler()),
3      ('logreg', LogisticRegression(penalty='l1', C=10.0, solver='liblinear'))
4 ])
```

```
In [19]: ❶ 1 logreg_pipeline.fit(X_train, y_train)
```

```
Out[19]: Pipeline(steps=[('scale', StandardScaler()),
                           ('logreg',
                            LogisticRegression(C=10.0, penalty='l1', solver='liblinear'))])
```

**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 [20]: ❶ 1 logreg_pipeline.score(X_test, y_test)
```

```
Out[20]: 0.7687443541102078
```

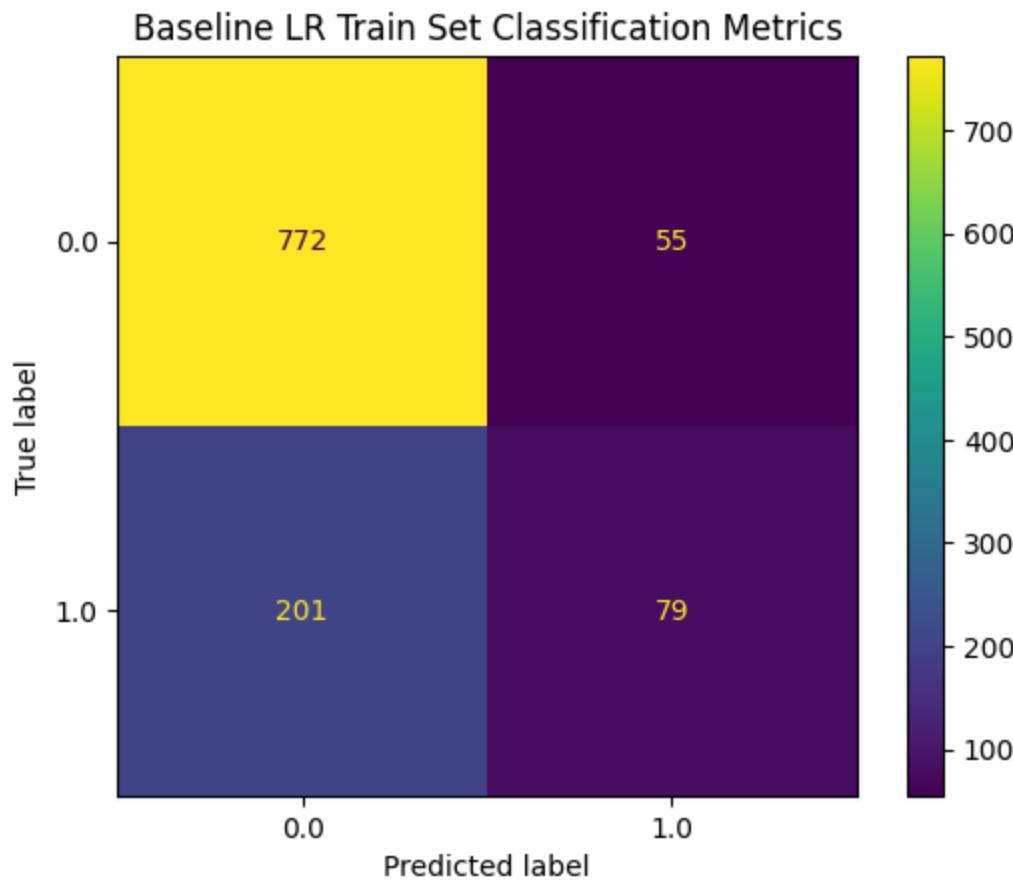
```
In [21]: ❶ 1 y_pred = logreg_pipeline.predict(X_test)
```

In [22]:

```

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))

```



	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

pivoted\_df, Logistic Regression Model, SMOTE

In [23]:

```

1 smote = SMOTE()
2 X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_t
3
4 print(pd.Series(y_train_resampled).value_counts())

```

```

1.0    1945
0.0    1945
Name: Surgery, dtype: int64

```

```
In [24]: 1 logreg_pipeline.fit(X_train_resampled, y_train_resampled)
```

```
Out[24]: Pipeline(steps=[('scale', StandardScaler()), ('logreg', LogisticRegression(C=10.0, penalty='l1', solver='liblinear'))])
```

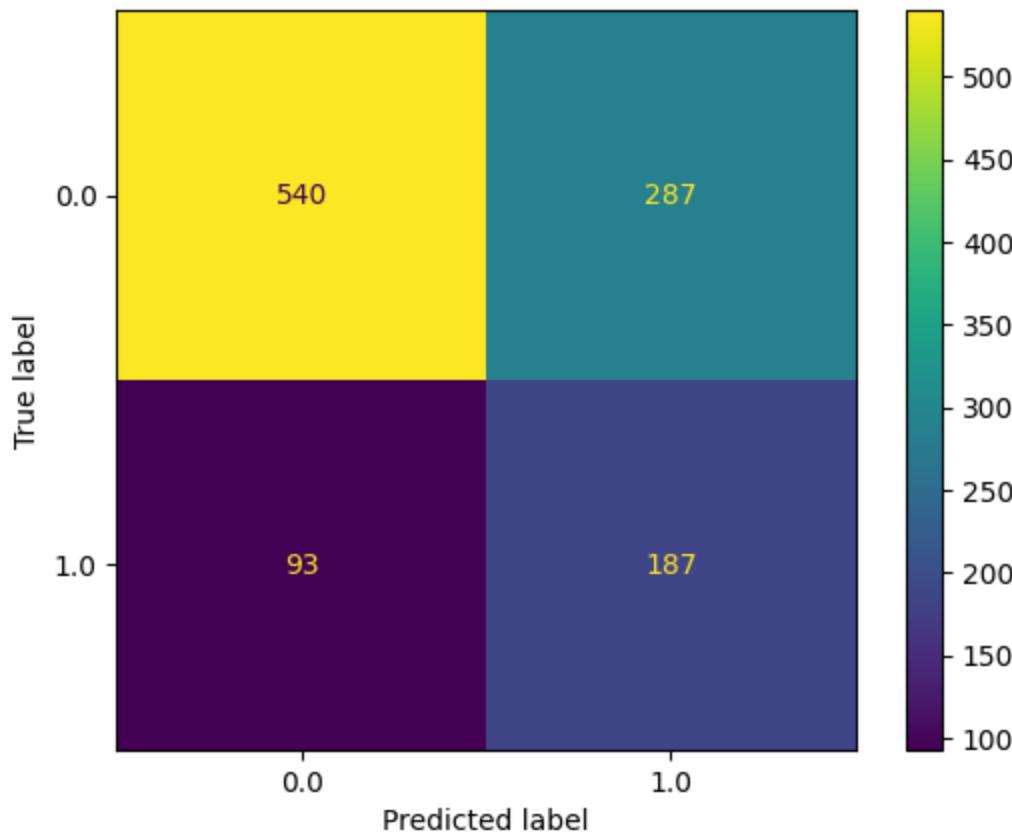
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with [nbviewer.org](#).

```
In [25]: 1 y_pred_resampled = logreg_pipeline.predict(X_test)
```

```
In [26]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred_resampled)
          2 plt.title('Baseline LR Train Set Classification Metrics after SMOTE')
          3 plt.show()
          4 print(classification_report(y_test, y_pred_resampled))
```

Baseline LR Train Set Classification Metrics after SMOTE



	precision	recall	f1-score	support
0.0	0.85	0.65	0.74	827
1.0	0.39	0.67	0.50	280
accuracy			0.66	1107
macro avg	0.62	0.66	0.62	1107
weighted avg	0.74	0.66	0.68	1107

Predicting 1.0 improved, however the rest of the model suffered. Can try condensing the number of features and see how that plays out.

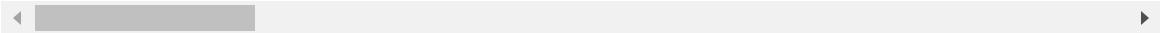
cond\_pivoted\_df, Logistic Regression, Baseline

```
In [23]: 1 cond_pivoted_df = pd.read_csv(io.BytesIO(uploaded['~/Documents/Flatiron...
```

In [84]: 1 cond\_pivoted\_df.head()

	season	Age	Throws	Surgery	BB_release_speed_weighted_avg	FB_release_speed_weighted_avg
0	2008	37.0	1	0.0	75.425843	91.0
1	2009	38.0	1	0.0	78.181818	93.0
2	2010	39.0	1	0.0	74.666667	93.0
3	2011	40.0	1	0.0	76.885714	91.0
4	2012	41.0	1	0.0	76.427273	91.0

5 rows × 39 columns



In [86]: 1 cond\_pivoted\_df.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3688 entries, 0 to 3687
Data columns (total 39 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   season          3688 non-null    int64  
 1   Age              3688 non-null    float64 
 2   Throws           3688 non-null    int64  
 3   Surgery          3688 non-null    float64 
 4   BB_release_speed_weighted_avg  3688 non-null    float64 
 5   FB_release_speed_weighted_avg  3688 non-null    float64 
 6   OS_release_speed_weighted_avg  3688 non-null    float64 
 7   OT_release_speed_weighted_avg  3688 non-null    float64 
 8   SB_release_speed_weighted_avg  3688 non-null    float64 
 9   BB_release_pos_x_weighted_avg 3688 non-null    float64 
 10  FB_release_pos_x_weighted_avg 3688 non-null    float64 
 11  OS_release_pos_x_weighted_avg 3688 non-null    float64 
 12  OT_release_pos_x_weighted_avg 3688 non-null    float64 
 13  SB_release_pos_x_weighted_avg 3688 non-null    float64 
 14  BB_release_pos_y_weighted_avg 3688 non-null    float64 
 15  FB_release_pos_y_weighted_avg 3688 non-null    float64 
 16  OS_release_pos_y_weighted_avg 3688 non-null    float64 
 17  OT_release_pos_y_weighted_avg 3688 non-null    float64 
 18  SB_release_pos_y_weighted_avg 3688 non-null    float64 
 19  BB_release_pos_z_weighted_avg 3688 non-null    float64 
 20  FB_release_pos_z_weighted_avg 3688 non-null    float64 
 21  OS_release_pos_z_weighted_avg 3688 non-null    float64 
 22  OT_release_pos_z_weighted_avg 3688 non-null    float64 
 23  SB_release_pos_z_weighted_avg 3688 non-null    float64 
 24  BB_vx0_weighted_avg         3688 non-null    float64 
 25  FB_vx0_weighted_avg         3688 non-null    float64 
 26  OS_vx0_weighted_avg         3688 non-null    float64 
 27  OT_vx0_weighted_avg         3688 non-null    float64 
 28  SB_vx0_weighted_avg         3688 non-null    float64 
 29  BB_vy0_weighted_avg         3688 non-null    float64 
 30  FB_vy0_weighted_avg         3688 non-null    float64 
 31  OS_vy0_weighted_avg         3688 non-null    float64 
 32  OT_vy0_weighted_avg         3688 non-null    float64 
 33  SB_vy0_weighted_avg         3688 non-null    float64 
 34  BB_vz0_weighted_avg         3688 non-null    float64 
 35  FB_vz0_weighted_avg         3688 non-null    float64 
 36  OS_vz0_weighted_avg         3688 non-null    float64 
 37  OT_vz0_weighted_avg         3688 non-null    float64 
 38  SB_vz0_weighted_avg         3688 non-null    float64 
dtypes: float64(37), int64(2)
memory usage: 1.1 MB
```

In [31]: 1 cond\_pivoted\_df['Surgery'].value\_counts()

```
Out[31]: 0.0      2772
        1.0      916
Name: Surgery, dtype: int64
```

```
In [83]: 1 cond_pivoted_df.drop(columns=['pitcher'], inplace=True)
```

```
In [18]: 1 y = cond_pivoted_df['Surgery']
2 X = cond_pivoted_df.drop('Surgery', axis=1)
```

```
In [19]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [23]: 1 logreg_pipeline = Pipeline([
2     ('scale', StandardScaler()),
3     ('logreg', LogisticRegression(solver='liblinear'))
4 ])
5
6 # Define the parameter grid to search over
7 param_grid = {
8     'logreg__C': [0.001, 0.01, 0.1, 1, 10, 100], # Regularization strength
9     'logreg__penalty': ['l1', 'l2'] # Norm used in the penalization
10 }
11
12 # Initialize GridSearchCV with the pipeline, parameter grid, and desired cross-validation splits
13 grid_search = GridSearchCV(logreg_pipeline, param_grid, cv=5, scoring='accuracy')
14
15 # Assuming X_train and y_train are already defined
16 grid_search.fit(X_train_resampled, y_train_resampled)
17
18 # Best parameters found
19 print("Best parameters: ", grid_search.best_params_)
20
21 # Best cross-validation score
22 print("Best cross-validation score: {:.2f}".format(grid_search.best_score_))
23
24 # Test set score using the best parameters
25 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: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.

```
warnings.warn(
```

```
Best parameters: {'logreg__C': 10, 'logreg__penalty': 'l1'}
Best cross-validation score: 0.67
Test set score: 0.64
```

```
In [26]: 1 logreg_pipeline = Pipeline([
2     ('scale', StandardScaler()),
3     ('logreg', LogisticRegression(penalty='l1', C=10.0, solver='liblinear'))
4 ])
```

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

```
Out[14]: Pipeline(steps=[('scale', StandardScaler()), ('logreg', LogisticRegression(penalty='l1', solver='liblinear'))])
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**

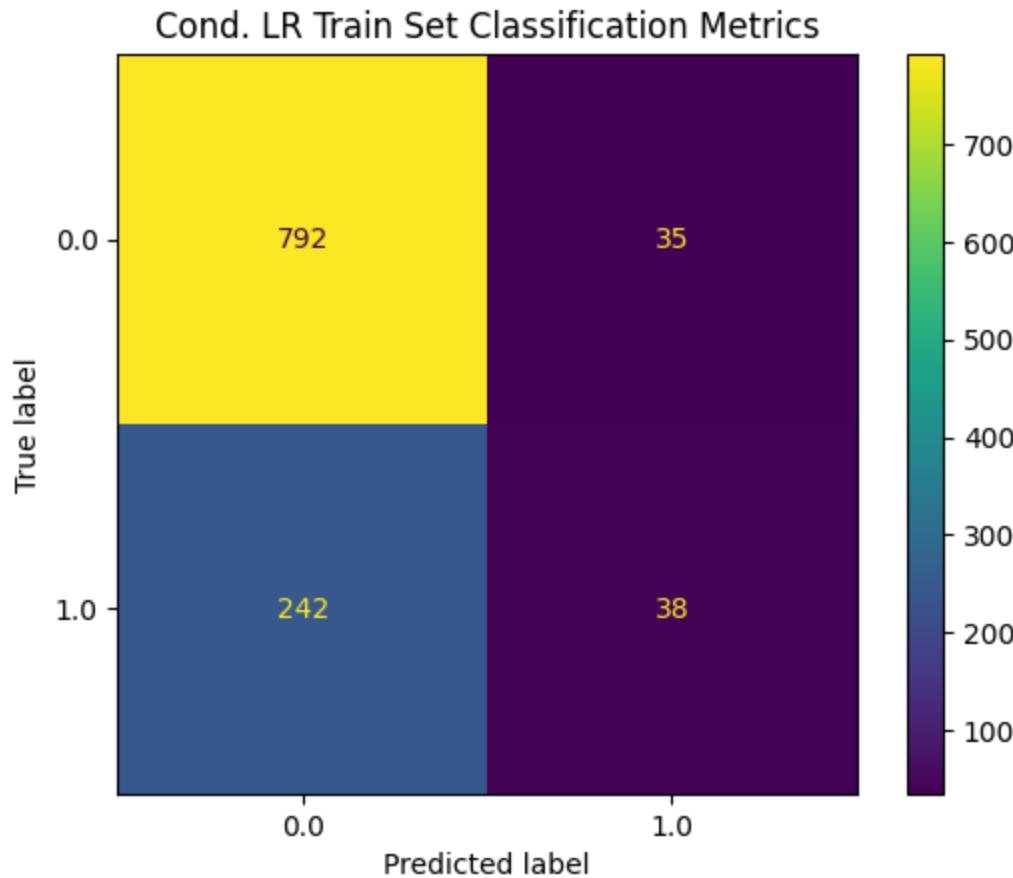
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

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

```
Out[15]: 0.7497741644083108
```

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

```
In [17]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('Cond. LR Train Set Classification Metrics')
3 plt.show()
4 print(classification_report(y_test, y_pred))
```



	precision	recall	f1-score	support
0.0	0.77	0.96	0.85	827
1.0	0.52	0.14	0.22	280
accuracy			0.75	1107
macro avg	0.64	0.55	0.53	1107
weighted avg	0.70	0.75	0.69	1107

Score for 1.0 predicting Surgery too low. Now try with SMOTE.

cond\_pivoted\_df, Logistic Regression Model, SMOTE

```
In [25]: 1 smote = SMOTE()
2 X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_t
3
4 print(pd.Series(y_train_resampled).value_counts())
```

```
1.0    1945
0.0    1945
Name: Surgery, dtype: int64
```

```
In [27]: ⏎ 1 logreg_pipeline.fit(X_train_resampled, y_train_resampled)
```

```
Out[27]: Pipeline(steps=[('scale', StandardScaler()), ('logreg', LogisticRegression(C=10.0, penalty='l1', solver='liblinear'))])
```

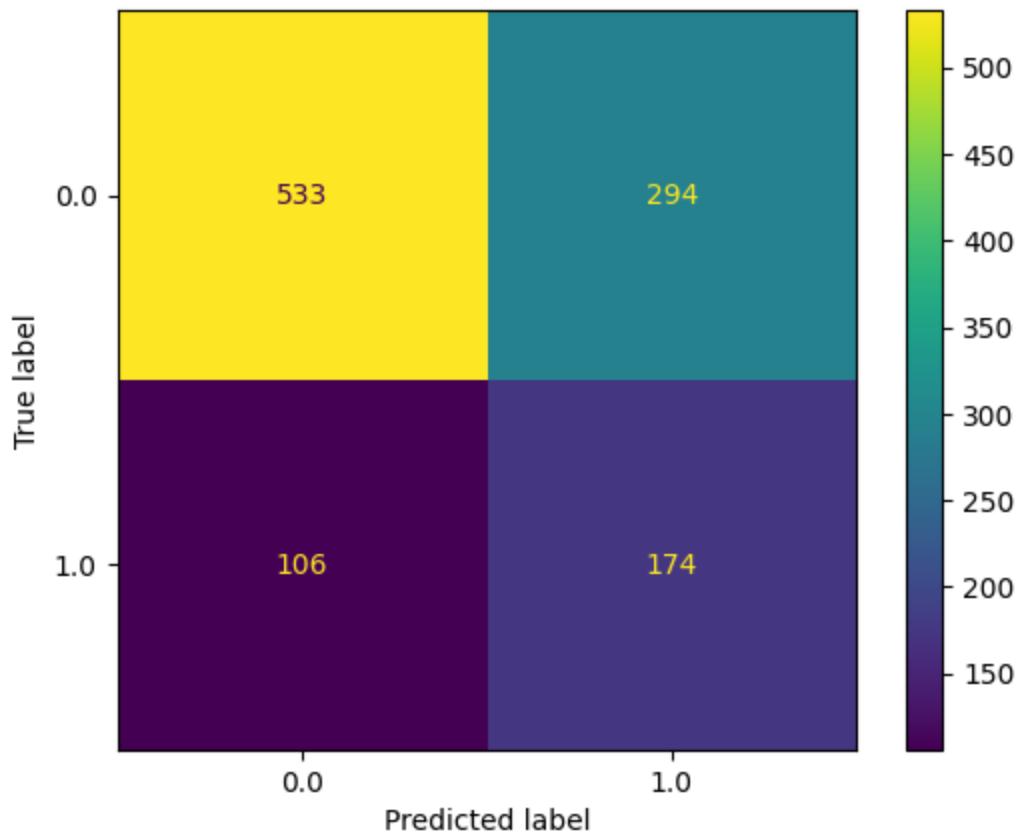
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with [nbviewer.org](#).

```
In [28]: ⏎ 1 y_pred_resampled = logreg_pipeline.predict(X_test)
```

```
In [29]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred_resampled)
2 plt.title('Cond. Baseline LR Train Set Classification Metrics after SMOTE')
3 plt.show()
4 print(classification_report(y_test, y_pred_resampled))
```

Cond. Baseline LR Train Set Classification Metrics after SMOTE



	precision	recall	f1-score	support
0.0	0.83	0.64	0.73	827
1.0	0.37	0.62	0.47	280
accuracy			0.64	1107
macro avg	0.60	0.63	0.60	1107
weighted avg	0.72	0.64	0.66	1107

pivot\_df is better at predicting TP & TN. cond\_pivot\_df is better at FP & FN. Trade-off between two models is slim. SMOTE helped, can I address class imbalance further, or will that lead to issues? Still need to see improvement. Need to try PCA, Random Forest, XG Boost

Try to see how GridSearch works when I add in class balance. Will try this out for both pivot\_df and cond\_pivot\_df

pivot\_df, Logistic Regression Model, adjust class weights

```
In [14]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [15]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [13]: 1 # Define the parameter grid to search over, including class weights
2 class_weights = [None, 'balanced', {0: 1, 1: 2}, {0: 1, 1: 3}, {0: 1,
3 param_grid = {
4     'logreg_C': [0.01, 0.1, 1, 10],
5     'logreg_penalty': ['l1', 'l2'],
6     'logreg_class_weight': class_weights,
7     'logreg_max_iter': [5000],
8     'logreg_tol': [0.01]
9 }
10
11 # Create a scoring function that focuses on recall for the positive class
12 recall_scorer = make_scorer(recall_score, pos_label=1)
13
14 # Initialize GridSearchCV with the pipeline, parameter grid, and recall
15 grid_search = GridSearchCV(weight_logreg_pipeline, param_grid, cv=5, s
16
17 # Fit the grid search to the data
18 grid_search.fit(X_train, y_train)
19
20 # Print the best parameters found and the best recall score
21 print("Best parameters: ", grid_search.best_params_)
22 print("Best cross-validation recall score: {:.2f}".format(grid_search.
23
24 # Evaluate the best model on the test set
25 best_model = grid_search.best_estimator_
26 y_pred = best_model.predict(X_test)
27 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': 0.01}
Best cross-validation recall score: 0.87
Test set recall score: 0.85
```

```
In [16]: 1 weight_logreg_pipeline = Pipeline([
2     ('scale', StandardScaler()),
3     ('logreg', LogisticRegression(penalty='l1', C=0.01, class_weight={0:
4 })]
```

In [17]: 1 weight\_logreg\_pipeline.fit(X\_train, y\_train)

Out[17]: Pipeline(steps=[('scale', StandardScaler()), ('logreg', LogisticRegression(C=0.01, class\_weight={0: 1, 1: 5}, penalty='l1', solver='liblinear'))])

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

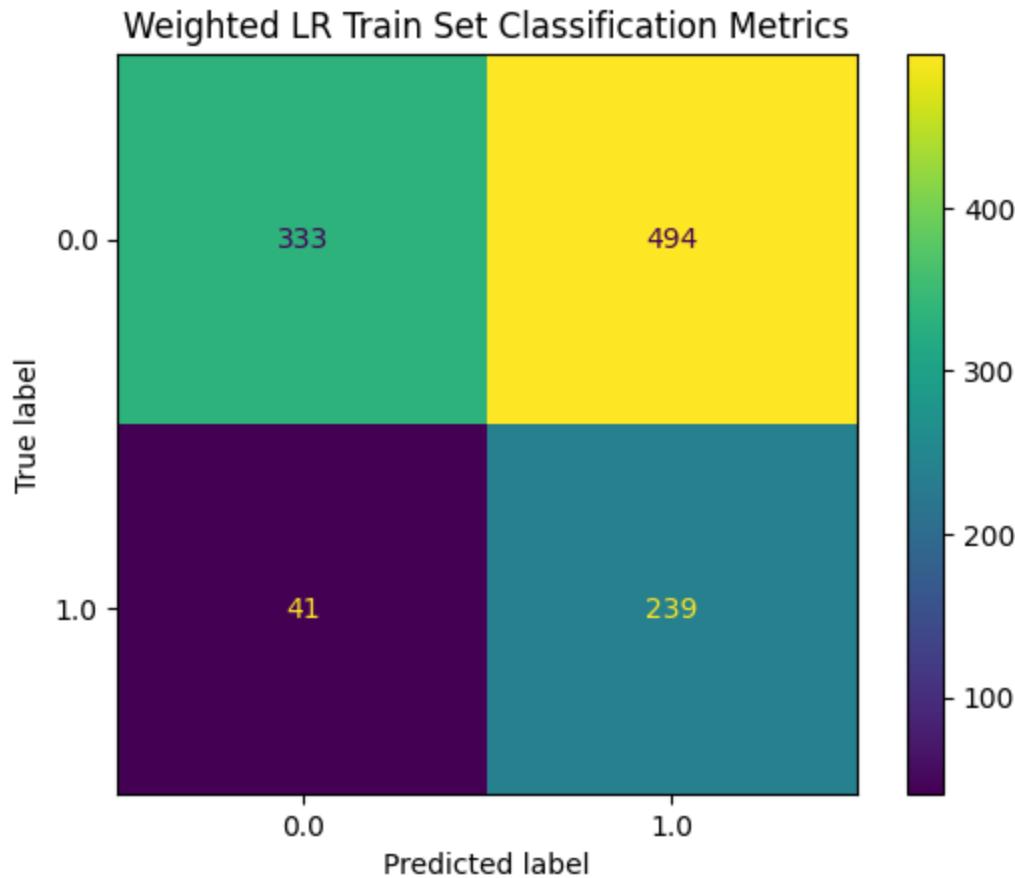
On GitHub, the HTML representation is unable to render, please try loading this page with [nbviewer.org](#).

In [18]: 1 weight\_logreg\_pipeline.score(X\_test, y\_test)

Out[18]: 0.5167118337850045

In [20]: 1 y\_pred = weight\_logreg\_pipeline.predict(X\_test)

```
In [21]: 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))
```



	precision	recall	f1-score	support
0.0	0.89	0.40	0.55	827
1.0	0.33	0.85	0.47	280
accuracy			0.52	1107
macro avg	0.61	0.63	0.51	1107
weighted avg	0.75	0.52	0.53	1107

This looks better. The number of False negatives has decreased dramatically. Would like to see further improvement to be more accurate all around. Now to test cond\_pivoted\_df

```
In [74]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [75]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [10]: Pipeline([
    ('scale', StandardScaler()),
    ('pca', PCA()),
    ('logreg', LogisticRegression(solver='liblinear'))
])
class_weights = [None, 'balanced', {0: 1, 1: 2}, {0: 1, 1: 3}, {0: 1,
param_grid = {
    'pca_n_components': [5, 10, 20, 30],
    'logreg_C': [0.01, 0.1, 1, 10],
    'logreg_penalty': ['l1', 'l2'],
    'logreg_class_weight': class_weights,
    'logreg_max_iter': [5000],
    'logreg_tol': [0.01]
}
recall_scorer = make_scorer(recall_score, pos_label=1)
grid_search = GridSearchCV(logreg_pipeline, param_grid, cv=5, scoring=
grid_search.fit(X_train, y_train)
print("Best parameters: ", grid_search.best_params_)
print("Best cross-validation recall score: {:.2f}".format(grid_search.
best_model = grid_search.best_estimator_
y_pred = best_model.predict(X_test)
print("Test set recall score: {:.2f}".format(recall_score(y_test, y_pr
```

```
Best parameters: {'logreg_C': 0.1, 'logreg_class_weight': {0: 1, 1: 5}, 'logreg_max_iter': 5000, 'logreg_penalty': 'l1', 'logreg_tol': 0.01, 'pca_n_components': 5}
Best cross-validation recall score: 0.92
Test set recall score: 0.90
```

```
In [76]: Pipeline([
    ('scale', StandardScaler()),
    ('pca', PCA(n_components=5)),
    ('logreg', LogisticRegression(penalty='l1', C=0.1, class_weight={0
])
```

```
In [77]: logreg_pipeline.fit(X_train, y_train)
```

```
Out[77]: Pipeline(steps=[('scale', StandardScaler()), ('pca', PCA(n_components=
5)),
    ('logreg',
        LogisticRegression(C=0.1, 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.**

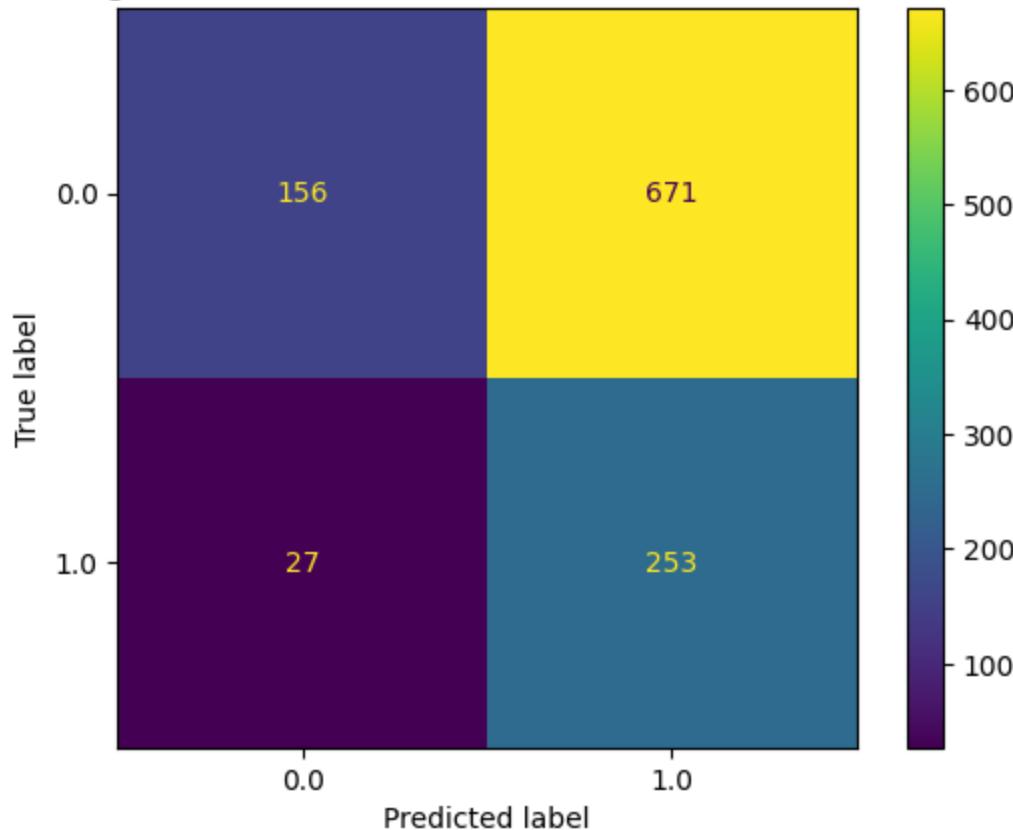
```
In [78]: 1 logreg_pipeline.score(X_test, y_test)
```

```
Out[78]: 0.36946702800361336
```

```
In [79]: 1 y_pred = logreg_pipeline.predict(X_test)
```

```
In [80]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('Weighted LR with PCA Train Set Classification Metrics')
3 plt.show()
4 print(classification_report(y_test, y_pred))
```

Weighted LR with PCA Train Set Classification Metrics



	precision	recall	f1-score	support
0.0	0.85	0.19	0.31	827
1.0	0.27	0.90	0.42	280
accuracy			0.37	1107
macro avg	0.56	0.55	0.36	1107
weighted avg	0.71	0.37	0.34	1107

Too strong on recall info. Trues weakened significantly compared to the previous model. Will try a weak version of SMOTE.

```
In [123]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [124]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [131]: 1 logreg_pipeline = ImbPipeline([
2     ('scale', StandardScaler()),
3     ('smote', SMOTE()),
4     ('pca', PCA()),
5     ('logreg', LogisticRegression(solver='liblinear'))
6 ])
7
8 param_grid = {
9     'pca_n_components': [5, 10, 15, 20],
10    'smote_sampling_strategy': [0.55, 0.6, 0.65],
11    'logreg_C': [0.01, 0.05, 0.1, 1, 10],
12    'logreg_penalty': ['l1', 'l2'],
13    'logreg_class_weight': [None, 'balanced', {0: 1, 1: 2}],
14    'logreg_max_iter': [5000],
15    'logreg_tol': [0.01]
16 }
17
18 recall_scorer = make_scorer(recall_score, pos_label=1)
19 grid_search = GridSearchCV(logreg_pipeline, param_grid, cv=5, scoring=
20 grid_search.fit(X_train, y_train)
21
22 print("Best parameters: ", grid_search.best_params_)
23 print("Best cross-validation recall score: {:.2f}".format(grid_search.
24
25 best_logreg = grid_search.best_estimator_
26 y_pred = best_logreg.predict(X_test)
27 print("Test set recall score: {:.2f}".format(recall_score(y_test, y_pr
```

```
Best parameters: {'logreg_C': 0.1, 'logreg_class_weight': {0: 1, 1: 2}, 'logreg_max_iter': 5000, 'logreg_penalty': 'l2', 'logreg_tol': 0.01, 'pca_n_components': 5, 'smote_sampling_strategy': 0.65}
Best cross-validation recall score: 0.82
Test set recall score: 0.78
```

```
In [132]: 1 best_logreg.fit(X_train, y_train)
```

```
Out[132]: Pipeline(steps=[('scale', StandardScaler()),
                           ('smote', SMOTE(sampling_strategy=0.65)),
                           ('pca', PCA(n_components=5)),
                           ('logreg',
                            LogisticRegression(C=0.1, class_weight={0: 1, 1: 2},
                                               max_iter=5000, 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.**

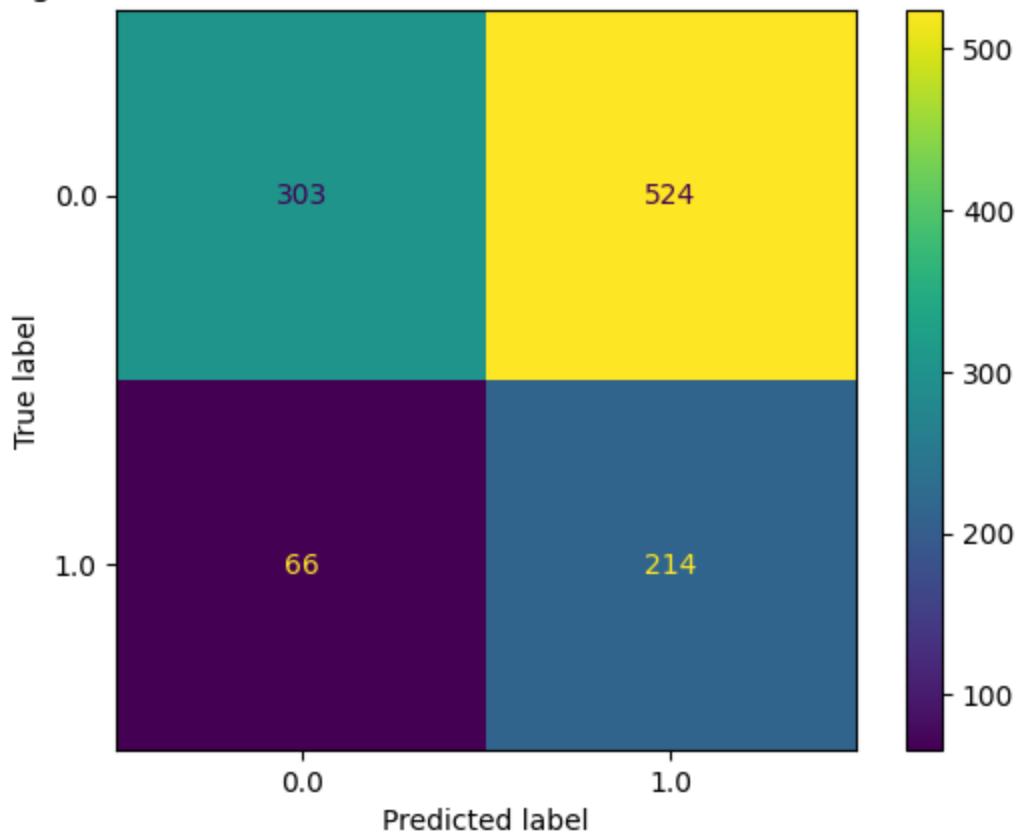
```
In [133]: 1 best_logreg.score(X_test, y_test)
```

```
Out[133]: 0.4670280036133695
```

```
In [134]: 1 y_pred = best_logreg.predict(X_test)
```

```
In [135]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('Weighted LR PCA Weak SMOTE Train Set Classification Metrics')
3 plt.show()
4 print(classification_report(y_test, y_pred))
```

Weighted LR PCA Weak SMOTE Train Set Classification Metrics



	precision	recall	f1-score	support
0.0	0.82	0.37	0.51	827
1.0	0.29	0.76	0.42	280
accuracy			0.47	1107
macro avg	0.56	0.57	0.46	1107
weighted avg	0.69	0.47	0.48	1107

Can try BorderlineSMOTE

```
In [136]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [137]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [138]: 1 logreg_pipeline = ImbPipeline([
2     ('scale', StandardScaler()),
3     ('blsmote', BorderlineSMOTE()),
4     ('pca', PCA()),
5     ('logreg', LogisticRegression(solver='liblinear'))
6 ])
7
8 param_grid = {
9     'pca_n_components': [5, 10, 15, 20],
10    'blsmote_sampling_strategy': [0.55, 0.6, 0.65],
11    'logreg_C': [0.01, 0.05, 0.1, 1, 10],
12    'logreg_penalty': ['l1', 'l2'],
13    'logreg_class_weight': [None, 'balanced', {0: 1, 1: 2}],
14    'logreg_max_iter': [5000],
15    'logreg_tol': [0.01]
16 }
17
18 recall_scorer = make_scorer(recall_score, pos_label=1)
19 grid_search = GridSearchCV(logreg_pipeline, param_grid, cv=5, scoring=
20 grid_search.fit(X_train, y_train)
21
22 print("Best parameters: ", grid_search.best_params_)
23 print("Best cross-validation recall score: {:.2f}".format(grid_search.
24
25 best_logreg = grid_search.best_estimator_
26 y_pred = best_logreg.predict(X_test)
27 print("Test set recall score: {:.2f}".format(recall_score(y_test, y_pr
```

```
Best parameters: {'blsmote_sampling_strategy': 0.65, 'logreg_C': 1, 'logreg_class_weight': {0: 1, 1: 2}, 'logreg_max_iter': 5000, 'logreg_penalty': 'l2', 'logreg_tol': 0.01, 'pca_n_components': 5}
Best cross-validation recall score: 0.82
Test set recall score: 0.77
```

```
In [139]: 1 best_logreg.fit(X_train, y_train)
```

```
Out[139]: Pipeline(steps=[('scale', StandardScaler()),
                           ('blsmote', BorderlineSMOTE(sampling_strategy=0.65)),
                           ('pca', PCA(n_components=5)),
                           ('logreg',
                            LogisticRegression(C=1, class_weight={0: 1, 1: 2},
                                               max_iter=5000, solver='liblinear',
                                               tol=0.01))])
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**

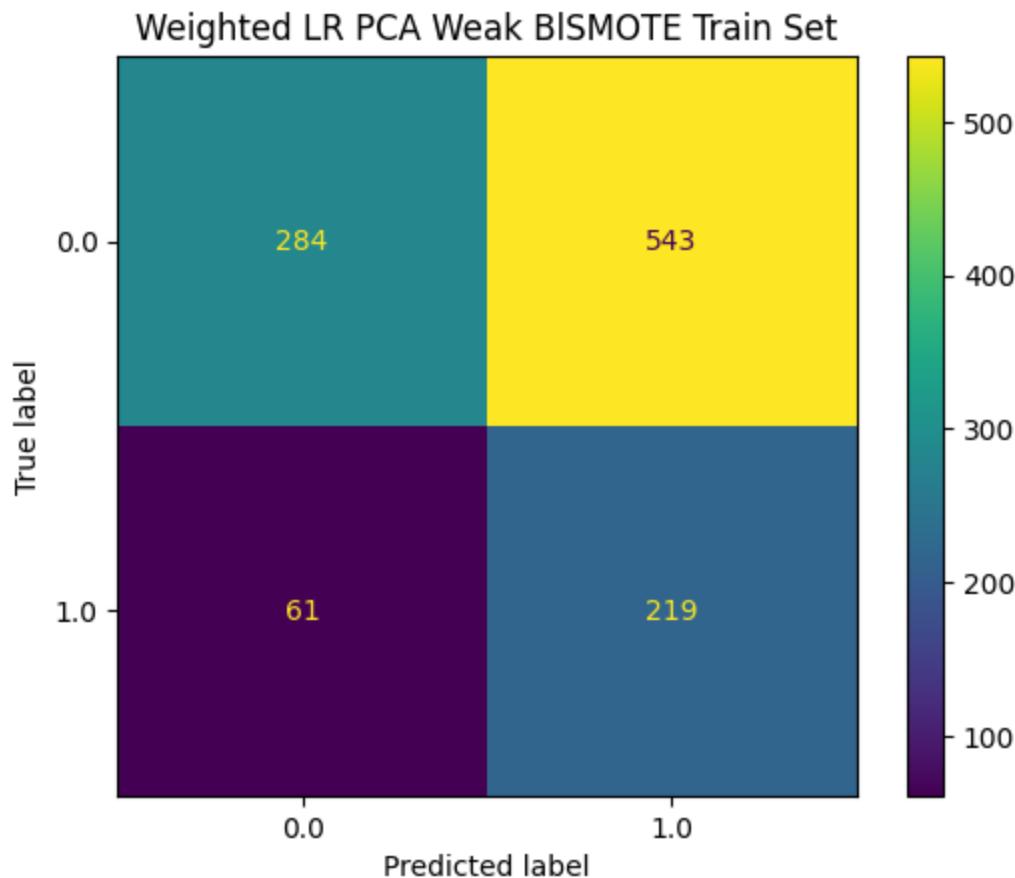
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [140]: 1 best_logreg.score(X_test, y_test)
```

```
Out[140]: 0.45438121047877145
```

```
In [141]: 1 y_pred = best_logreg.predict(X_test)
```

```
In [142]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('Weighted LR PCA Weak BISMOTE Train Set')
3 plt.show()
4 print(classification_report(y_test, y_pred))
```



	precision	recall	f1-score	support
0.0	0.82	0.34	0.48	827
1.0	0.29	0.78	0.42	280
accuracy			0.45	1107
macro avg	0.56	0.56	0.45	1107
weighted avg	0.69	0.45	0.47	1107

cond\_pivoted\_df, Logistic Regression Model, adjust class weights.

```
In [41]: 1 y = cond_pivoted_df['Surgery']
2 X = cond_pivoted_df.drop('Surgery', axis=1)
```

```
In [42]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [32]: 1 class_weights = [None, 'balanced', {0: 1, 1: 2}, {0: 1, 1: 3}, {0: 1,
2 param_grid = {
3     'logreg_C': [0.01, 0.1, 1, 10, 100],
4     'logreg_penalty': ['l1', 'l2'],
5     'logreg_class_weight': class_weights,
6     'logreg_max_iter': [5000],
7     'logreg_tol': [0.01]
8 }
9
10 recall_scorer = make_scorer(recall_score, pos_label=1)
11 grid_search = GridSearchCV(weight_logreg_pipeline, param_grid, cv=5, s
12 grid_search.fit(X_train_resampled, y_train_resampled)
13
14 print("Best parameters: ", grid_search.best_params_)
15 print("Best cross-validation recall score: {:.2f}".format(grid_search.
16
17 best_model = grid_search.best_estimator_
18 y_pred = best_model.predict(X_test)
19 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: 2}, 'logreg\_max\_iter': 5000, 'logreg\_penalty': 'l1', 'logreg\_tol': 0.01}  
 Best cross-validation recall score: 1.00  
 Test set recall score: 0.99

```
In [43]: 1 smote = SMOTE()
2 X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_t
3
4 print(pd.Series(y_train_resampled).value_counts())
```

1.0 1945  
 0.0 1945  
 Name: Surgery, dtype: int64

```
In [44]: 1 weight_logreg_pipeline = Pipeline([
2     ('scale', StandardScaler()),
3     ('logreg', LogisticRegression(penalty='l1', C=0.01, class_weight={
4 ])
```

```
In [45]: 1 weight_logreg_pipeline.fit(X_train_resampled, y_train_resampled)
```

```
Out[45]: Pipeline(steps=[('scale', StandardScaler()), ('logreg', LogisticRegression(C=0.01, class_weight={0: 1, 1: 3}, penalty='l1', solver='liblinear'))])
```

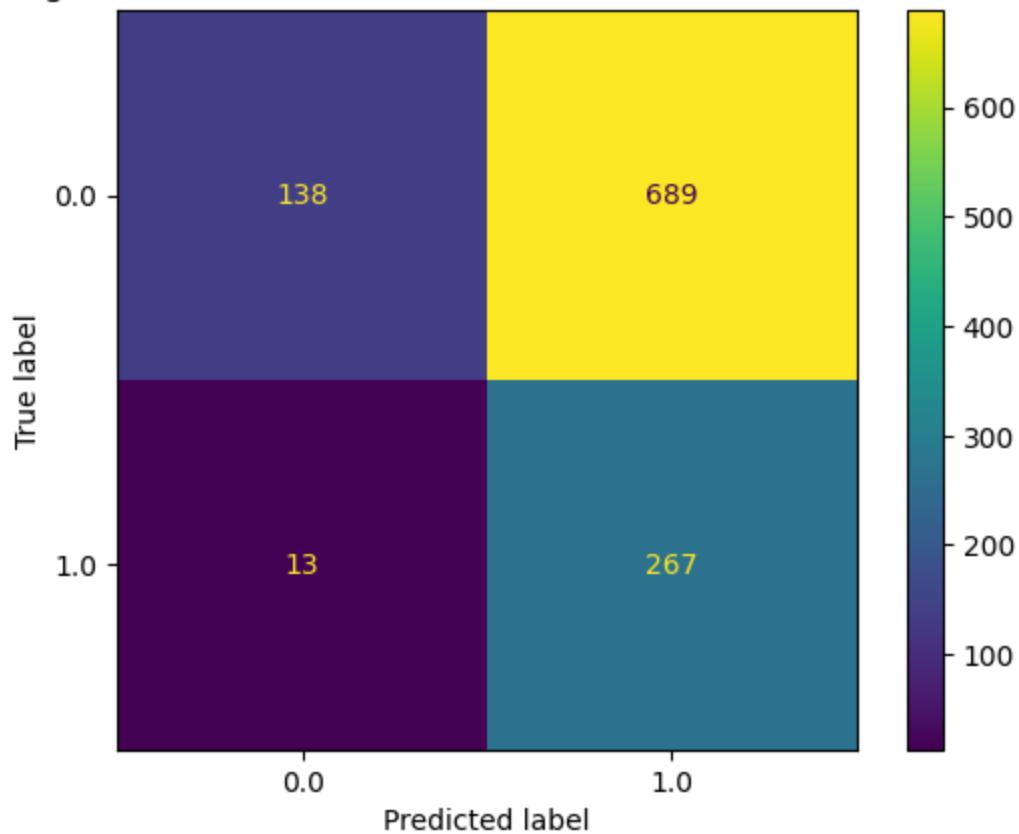
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with [nbviewer.org](http://nbviewer.org).

```
In [46]: 1 y_pred_resampled = weight_logreg_pipeline.predict(X_test)
```

```
In [47]: ⏷ 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred_resampled)
2 plt.title('Weight Cond. LR Train Set Classification Metrics after SMOTE')
3 plt.show()
4 print(classification_report(y_test, y_pred_resampled))
```

Weight Cond. LR Train Set Classification Metrics after SMOTE



	precision	recall	f1-score	support
0.0	0.91	0.17	0.28	827
1.0	0.28	0.95	0.43	280
accuracy			0.37	1107
macro avg	0.60	0.56	0.36	1107
weighted avg	0.75	0.37	0.32	1107

Can revisit the Logistic Regression model later. Time to try Random Forest, XG Boost, maybe CatBoost (not sure if it applies..)

pivoted\_df, Decision Tree (Baseline Model)

```
In [50]: ⏷ 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [51]: 1 | X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.
```

```
In [52]: 1 tree_clf = DecisionTreeClassifier(criterion='gini', max_depth=5)
          2 tree_clf.fit(X_train, y_train)
```

Out[52]: DecisionTreeClassifier(max\_depth=5)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

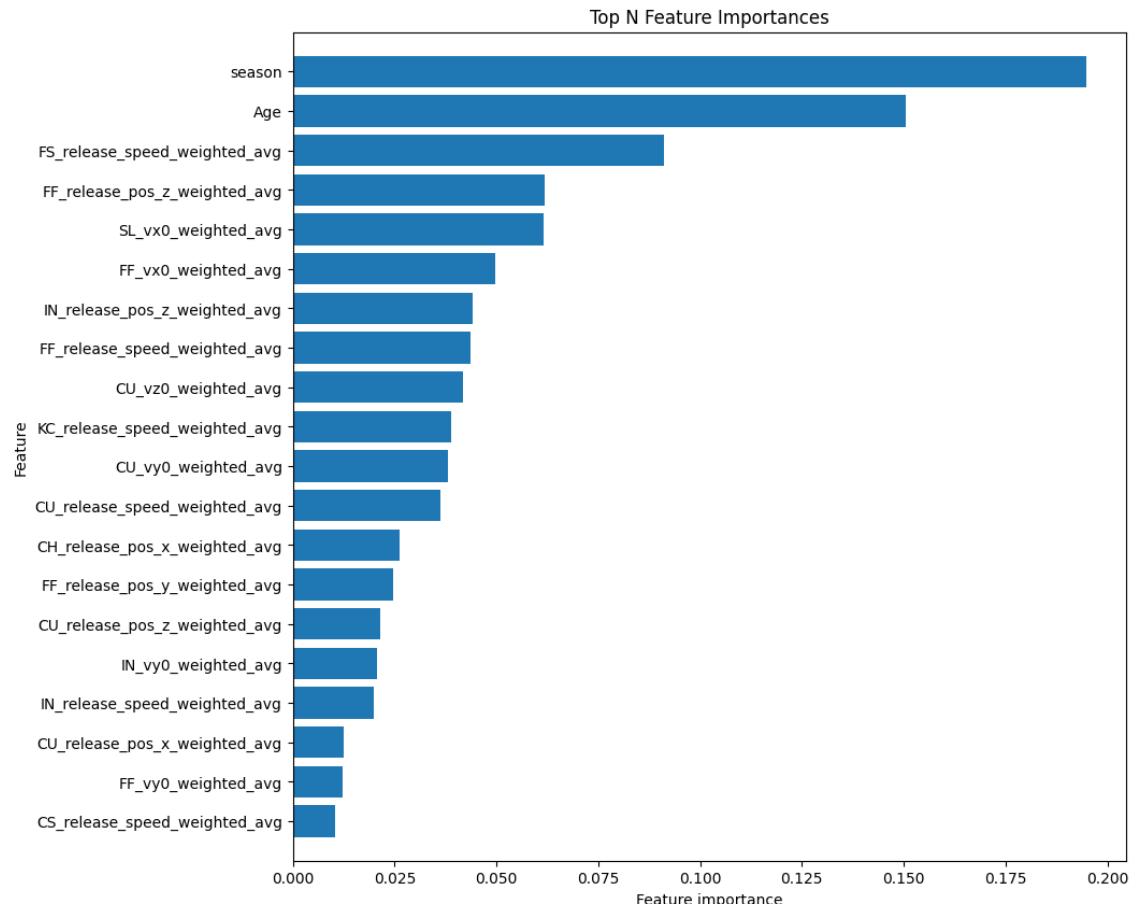
On GitHub, the HTML representation is unable to render, please try loading this page with [nbviewer.org](http://nbviewer.org).

```
In [53]: 1 tree_clf.feature_importances_
```

```
Out[53]: array([0.19477055, 0.15049678, 0. , 0. , 0. ,
   0.01033349, 0.03611204, 0. , 0. , 0. ,
   0.04363543, 0.09102768, 0.01991462, 0.03879671, 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0.02612076, 0. , 0.01246673,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0.02465828, 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0.02142519, 0. , 0. , 0. , 0.06182032,
   0. , 0.04426144, 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0.04978206, 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0.06168103, 0. , 0. , 0. , 0. ,
   0. , 0.0381441 , 0. , 0. , 0. ,
   0.0121275 , 0. , 0.02063147, 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0.04179382,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 0. ,
   0. , 0. , 0. , 0. , 1])
```

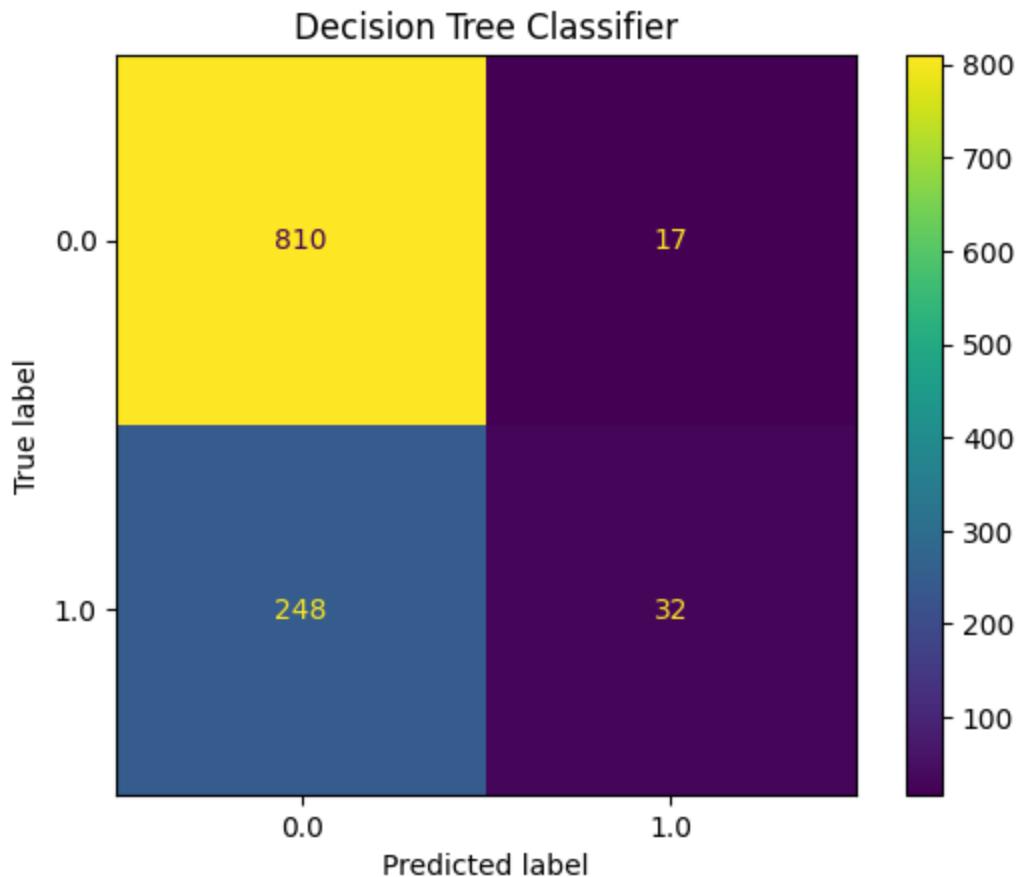
In [57]:

```
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='center')
7     plt.yticks(range(n_top_features), [X_train.columns[i] for i in indices])
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 [59]:

```
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	0.65	0.11	0.19	280
accuracy			0.76	1107
macro avg	0.71	0.55	0.53	1107
weighted avg	0.74	0.76	0.69	1107

```
In [66]: 1 param_grid = {
2     'criterion': ['gini', 'entropy'],
3     'max_depth': [None, 5, 10, 15, 20],
4     'min_samples_split': [2, 5, 10],
5     'min_samples_leaf': [1, 2, 4],
6     'class_weight': [None, '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
13 grid_search.fit(X_train, y_train)
14
15 print("Best parameters:", grid_search.best_params_)
16 print("Best score:", grid_search.best_score_)
17
18 best_tree = grid_search.best_estimator_
19 y_pred = best_tree.predict(X_test)
20 print("Test set accuracy:", accuracy_score(y_test, y_pred))
```

Best parameters: {'class\_weight': 'balanced', 'criterion': 'gini', 'max\_depth': 5, 'min\_samples\_leaf': 4, 'min\_samples\_split': 10}  
 Best score: 0.6540600393700788  
 Test set accuracy: 0.5817524841915086

```
In [61]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [62]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

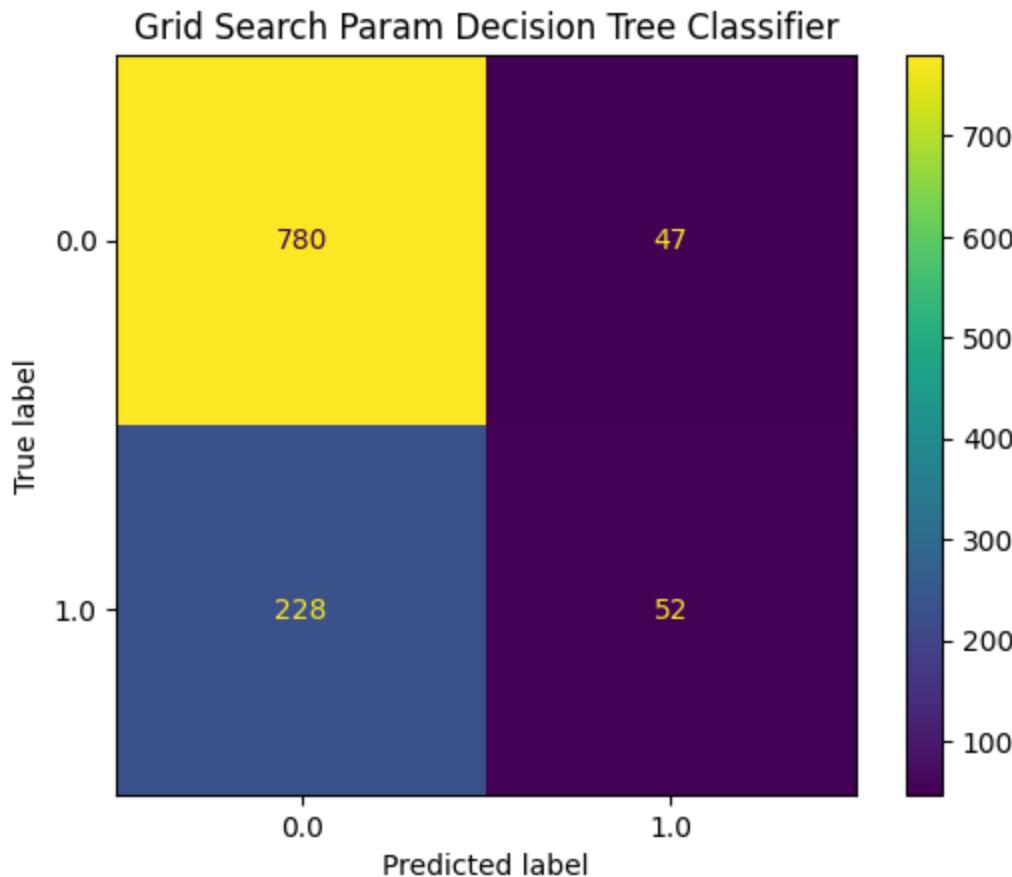
```
In [63]: 1 tree_clf = DecisionTreeClassifier(criterion='entropy', max_depth=5, mi
2 tree_clf.fit(X_train, y_train)
```

**Out[63]:** DecisionTreeClassifier(criterion='entropy', max\_depth=5, 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 [64]: pred = tree_clf.predict(X_test)
2
3 ConfusionMatrixDisplay.from_predictions(y_test, pred)
4 plt.title('Grid Search Param Decision Tree Classifier')
5 plt.show()
6 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.77	0.94	0.85	827
1.0	0.53	0.19	0.27	280
accuracy			0.75	1107
macro avg	0.65	0.56	0.56	1107
weighted avg	0.71	0.75	0.70	1107

Grid Search improved on TN & FN but worse on TP & TN. Would rather improve on TN than FN. Can try rerunning.

```
In [176]: y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [177]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [182]: ► 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_depth': 5, 'min_samples_leaf': 4, 'min_samples_split': 2}  
Best score: 0.6540600393700788  
Test recall score: 0.7607142857142857
```

```
In [183]: ► 1 best_tree.fit(X_train, y_train)
```

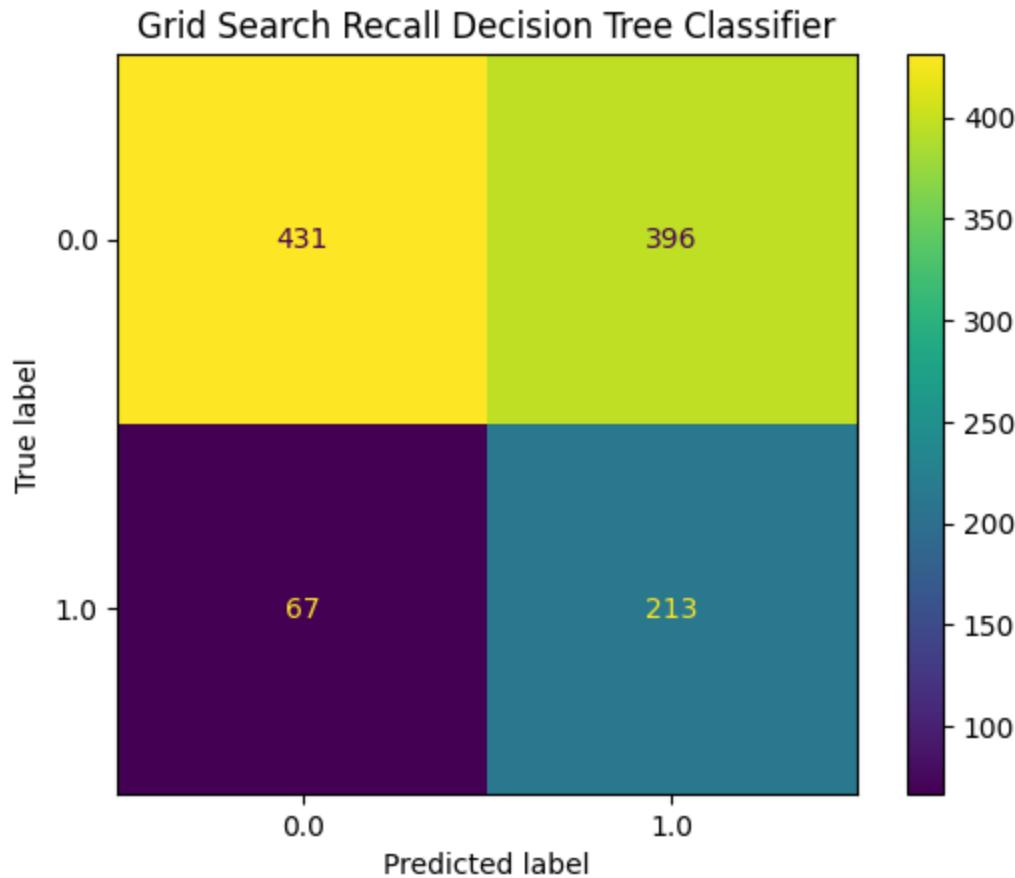
Out[183]: DecisionTreeClassifier(class\_weight='balanced', max\_depth=5, min\_samples\_leaf=4)

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 [184]: ► 1 pred = best_tree.predict(X_test)
```

```
In [186]: 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))
```



	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

```
In [ ]: ❶ 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))
```

```
In [187]: ❶ tree_clf = DecisionTreeClassifier(criterion='gini', class_weight='balan  
2 tree_clf.fit(X_train, y_train)
```

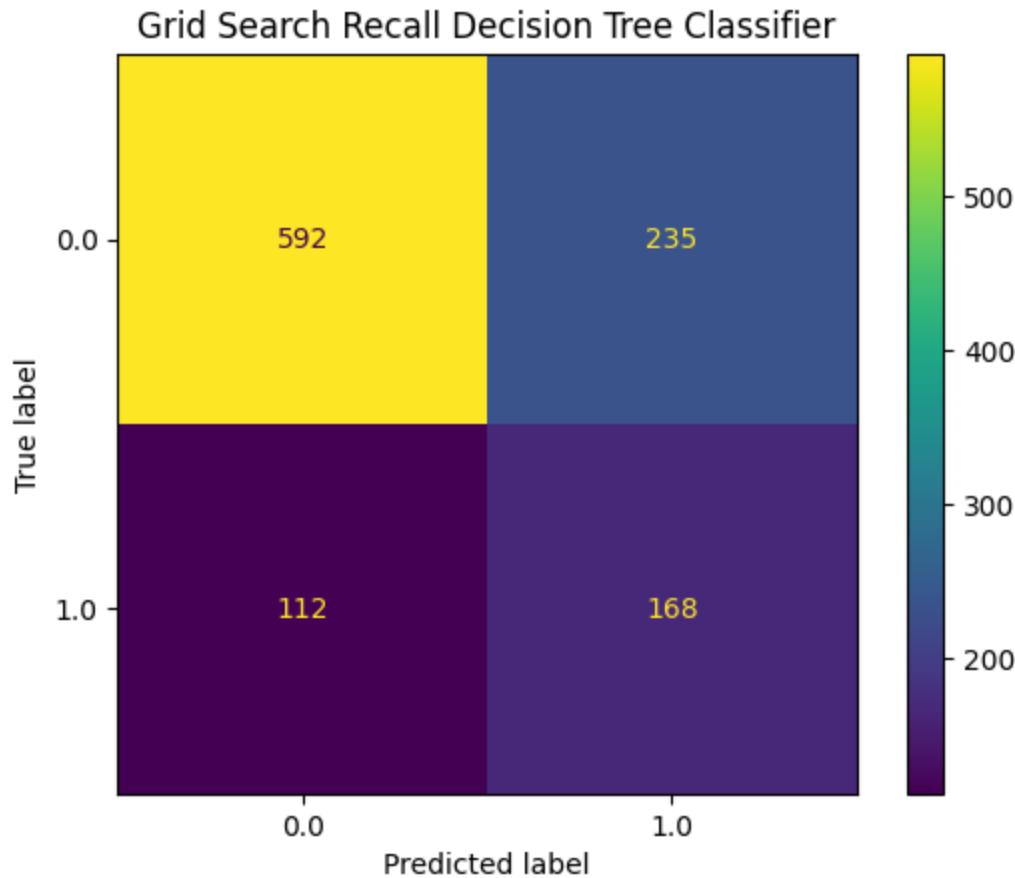
Out[187]: DecisionTreeClassifier(class\_weight='balanced', max\_depth=10)

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 [188]: ❶ pred = tree_clf.predict(X_test)
```

```
In [189]: 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))
```



	precision	recall	f1-score	support
0.0	0.84	0.72	0.77	827
1.0	0.42	0.60	0.49	280
accuracy			0.69	1107
macro avg	0.63	0.66	0.63	1107
weighted avg	0.73	0.69	0.70	1107

```
In [143]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [144]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [153]: 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     'ccp_alpha': [0.0, 0.01, 0.1],  
7     'class_weight': [None, 'balanced', {0:1, 1:2}, {0:1, 1:3}]  
8 }  
9  
10 tree_clf = DecisionTreeClassifier()  
11 scorer = make_scorer(recall_score)  
12 grid_search = GridSearchCV(estimator=tree_clf, param_grid=param_grid,  
13 grid_search.fit(X_train, y_train)  
14  
15 print("Best parameters:", grid_search.best_params_)  
16 print("Best score:", grid_search.best_score_)  
17  
18 best_tree = grid_search.best_estimator_  
19 y_pred = best_tree.predict(X_test)  
20 print("Test recall score:", recall_score(y_test, y_pred))
```

Best parameters: {'ccp\_alpha': 0.01, 'class\_weight': 'balanced', 'criterion': 'entropy', 'max\_depth': 10, 'min\_samples\_leaf': 1, 'min\_samples\_split': 2}  
Best score: 0.7420398622047244  
Test recall score: 0.7285714285714285

```
In [154]: 1 best_tree.fit(X_train, y_train)
```

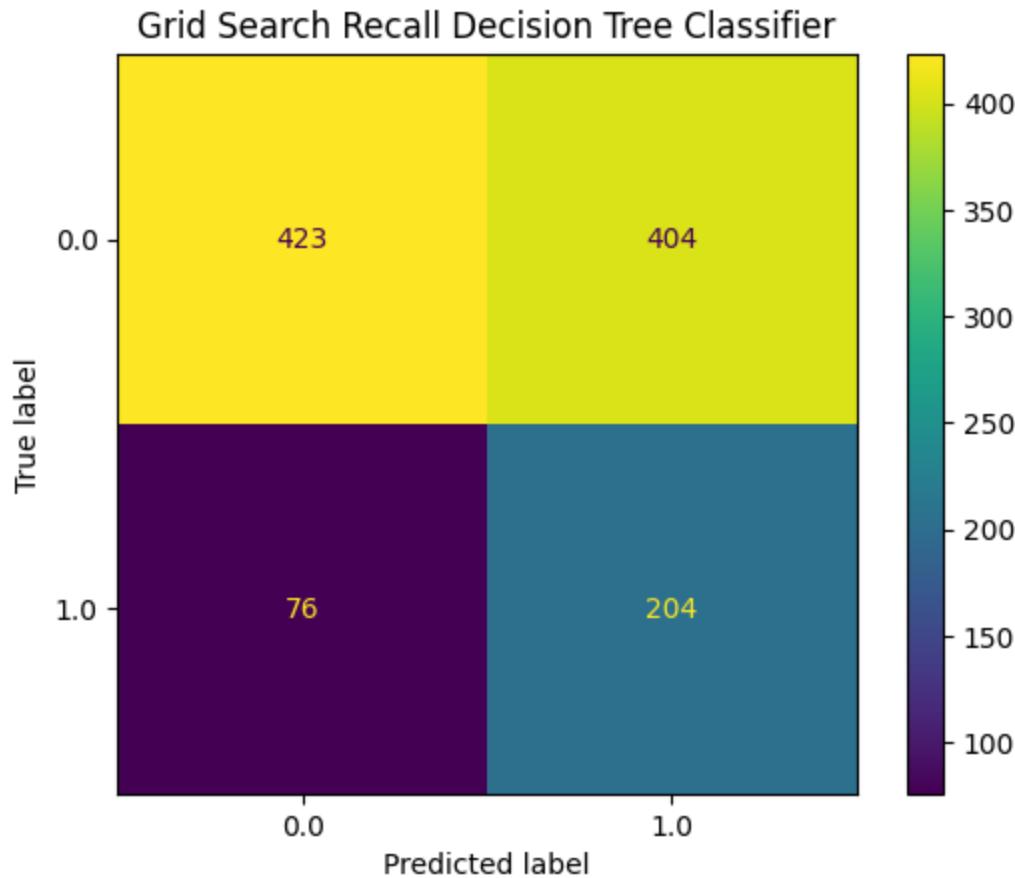
Out[154]: DecisionTreeClassifier(ccp\_alpha=0.01, class\_weight='balanced', criterion='entropy', max\_depth=10)

**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 [155]: 1 pred = best_tree.predict(X_test)
```

```
In [156]: 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))
```



	precision	recall	f1-score	support
0.0	0.85	0.51	0.64	827
1.0	0.34	0.73	0.46	280
accuracy			0.57	1107
macro avg	0.59	0.62	0.55	1107
weighted avg	0.72	0.57	0.59	1107

```
In [117]: 1 tree_clf = DecisionTreeClassifier(criterion='gini', class_weight='balanced')
2 tree_clf.fit(X_train, y_train)
```

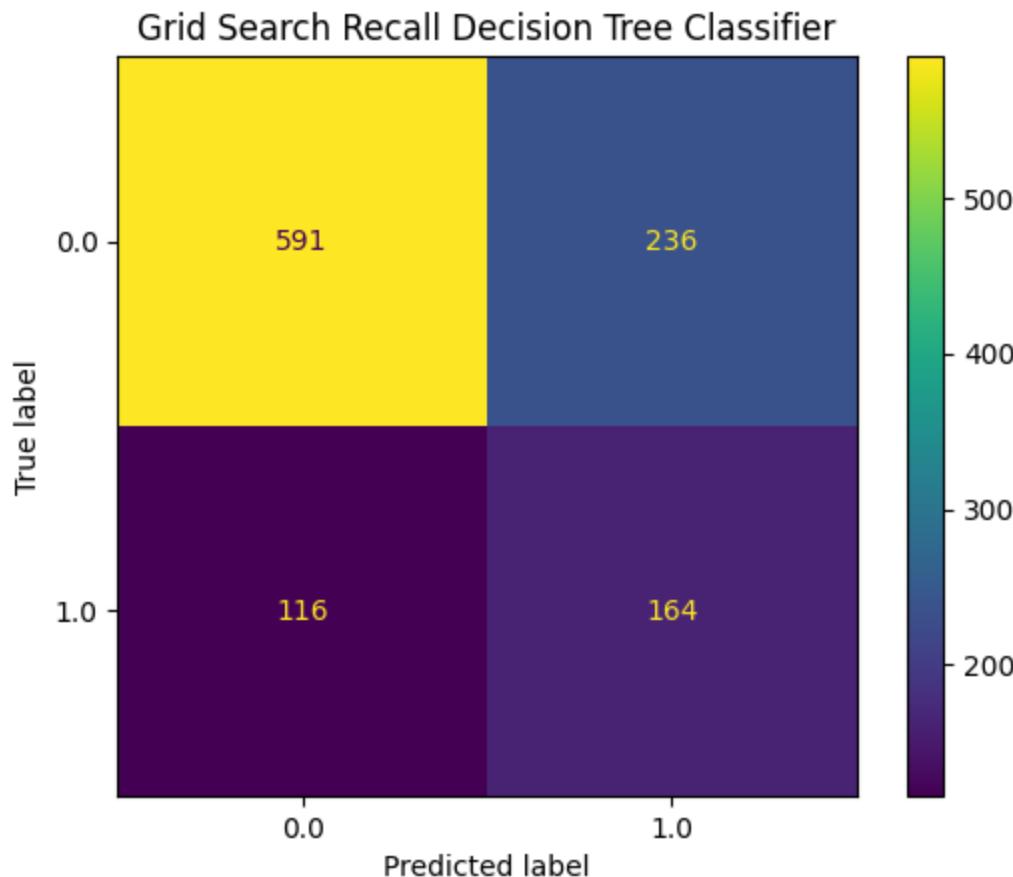
Out[117]: DecisionTreeClassifier(class\_weight='balanced', max\_depth=10)

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 [118]: 1 pred = tree_clf.predict(X_test)
```

```
In [119]: 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))
```



	precision	recall	f1-score	support
0.0	0.84	0.71	0.77	827
1.0	0.41	0.59	0.48	280
accuracy			0.68	1107
macro avg	0.62	0.65	0.63	1107
weighted avg	0.73	0.68	0.70	1107

Try to improve on this model, attempt PCA

```
In [25]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [8]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [19]: 1 tree_pipeline = ImbPipeline([
2     ('pca', PCA()),
3     ('smote', SMOTE()),
4     ('tree_clf', DecisionTreeClassifier())
5 ])
6
7 param_grid = {
8     'pca_n_components': [5, 10, 15],
9     'smote_sampling_strategy': [0.55, 0.6, 0.65, 0.7],
10    'tree_clf_criterion': ['gini', 'entropy'],
11    'tree_clf_max_depth': [3, 4, 5, 10],
12    'tree_clf_min_samples_split': [5, 10],
13    'tree_clf_min_samples_leaf': [1, 2, 3],
14    'tree_clf_class_weight': [None, 'balanced', {0:1, 1:1}, {0:1, 1:2}]
15 }
16
17 scorer = make_scorer(recall_score)
18 grid_search = GridSearchCV(tree_pipeline, param_grid, scoring=scorer,
19 grid_search.fit(X_train, y_train)
20
21 print("Best parameters:", grid_search.best_params_)
22 print("Best score:", grid_search.best_score_)
23
24 best_tree = grid_search.best_estimator_
25 y_pred = best_tree.predict(X_test)
26 print("Test set recall:", recall_score(y_test, y_pred))
```

Best parameters: {'pca\_n\_components': 15, 'smote\_sampling\_strategy': 0.7, 'tree\_clf\_class\_weight': {0: 1, 1: 2}, 'tree\_clf\_criterion': 'gini', 'tree\_clf\_max\_depth': 3, 'tree\_clf\_min\_samples\_leaf': 2, 'tree\_clf\_min\_samples\_split': 10}  
 Best score: 0.8226747047244094  
 Test set recall: 0.6535714285714286

```
In [20]: 1 best_tree
```

```
Out[20]: Pipeline(steps=[('pca', PCA(n_components=15)),
                         ('smote', SMOTE(sampling_strategy=0.7)),
                         ('tree_clf',
                          DecisionTreeClassifier(class_weight={0: 1, 1: 2}, max_depth=3,
                                                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 [21]: 1 best\_tree.fit(X\_train, y\_train)

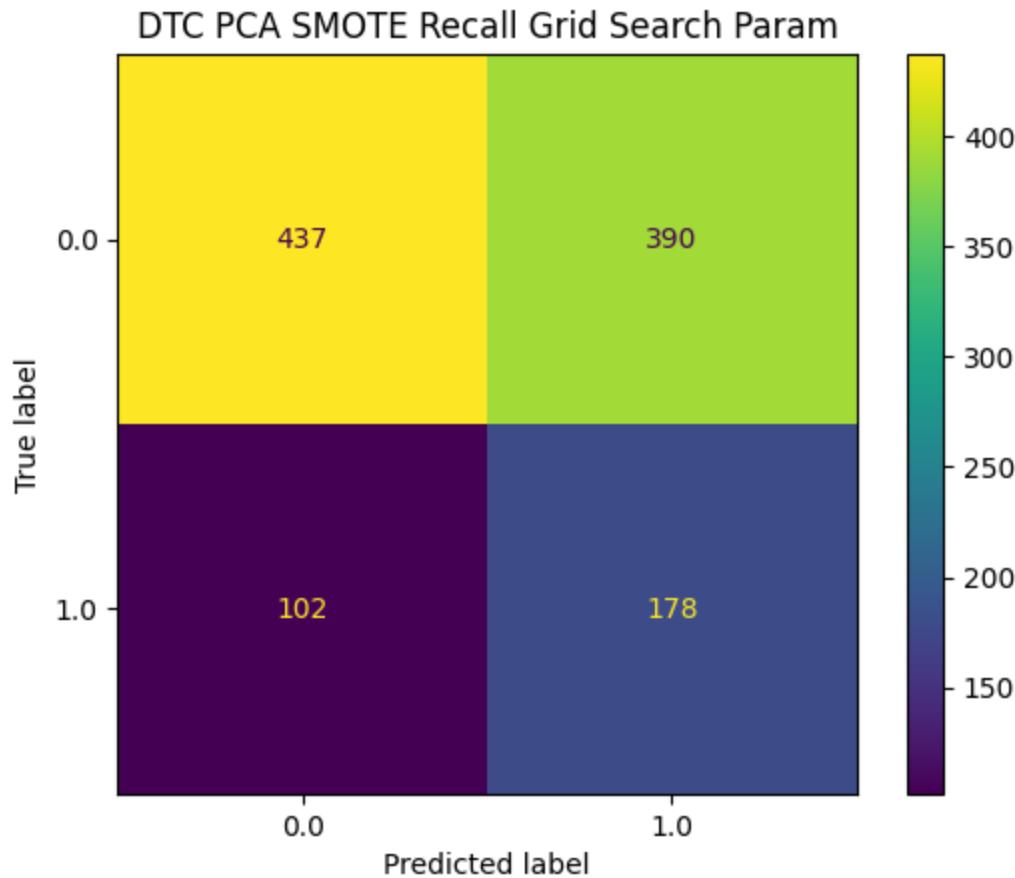
```
Out[21]: Pipeline(steps=[('pca', PCA(n_components=15)),
                         ('smote', SMOTE(sampling_strategy=0.7)),
                         ('tree_clf',
                          DecisionTreeClassifier(class_weight={0: 1, 1: 2}, max_depth=3,
                                                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 [22]: 1 y\_pred = best\_tree.predict(X\_test)

```
In [23]: ⏎ 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('DTC PCA SMOTE Recall Grid Search Param')
3 plt.show()
4 print(classification_report(y_test, y_pred))
```



	precision	recall	f1-score	support
0.0	0.81	0.53	0.64	827
1.0	0.31	0.64	0.42	280
accuracy			0.56	1107
macro avg	0.56	0.58	0.53	1107
weighted avg	0.68	0.56	0.58	1107

Try to further improve. Will try Borderline SMOTE

```
In [28]: ⏎ 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [29]: ⏎ 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [40]: 1 tree_pipeline = ImbPipeline([
2     ('pca', PCA()),
3     ('sm', BorderlineSMOTE()),
4     ('tree_clf', DecisionTreeClassifier())
5 ])
6
7 param_grid = {
8     'pca_n_components': [5, 10, 15],
9     'sm_sampling_strategy': [0.55, 0.6, 0.65, 0.7],
10    'tree_clf_criterion': ['gini', 'entropy'],
11    'tree_clf_max_depth': [3, 4, 5, 10],
12    'tree_clf_min_samples_split': [5, 10],
13    'tree_clf_min_samples_leaf': [1, 2, 3],
14    'tree_clf_class_weight': [None, 'balanced', {0:1, 1:1}, {0:1, 1:2}]
15 }
16
17 scorer = make_scorer(recall_score)
18 grid_search = GridSearchCV(tree_pipeline, param_grid, scoring=scorer,
19 grid_search.fit(X_train, y_train)
20
21 print("Best parameters:", grid_search.best_params_)
22 print("Best score:", grid_search.best_score_)
23
24 best_tree = grid_search.best_estimator_
25 y_pred = best_tree.predict(X_test)
26 print("Test set recall score:", recall_score(y_test, y_pred))
```

Best parameters: {'pca\_n\_components': 10, 'sm\_sampling\_strategy': 0.7, 'tree\_clf\_class\_weight': {0: 1, 1: 2}, 'tree\_clf\_criterion': 'gini', 'tree\_clf\_max\_depth': 3, 'tree\_clf\_min\_samples\_leaf': 2, 'tree\_clf\_min\_samples\_split': 5}  
 Best score: 0.872613188976378  
 Test set recall score: 0.8642857142857143

```
In [41]: 1 best_tree.fit(X_train, y_train)
```

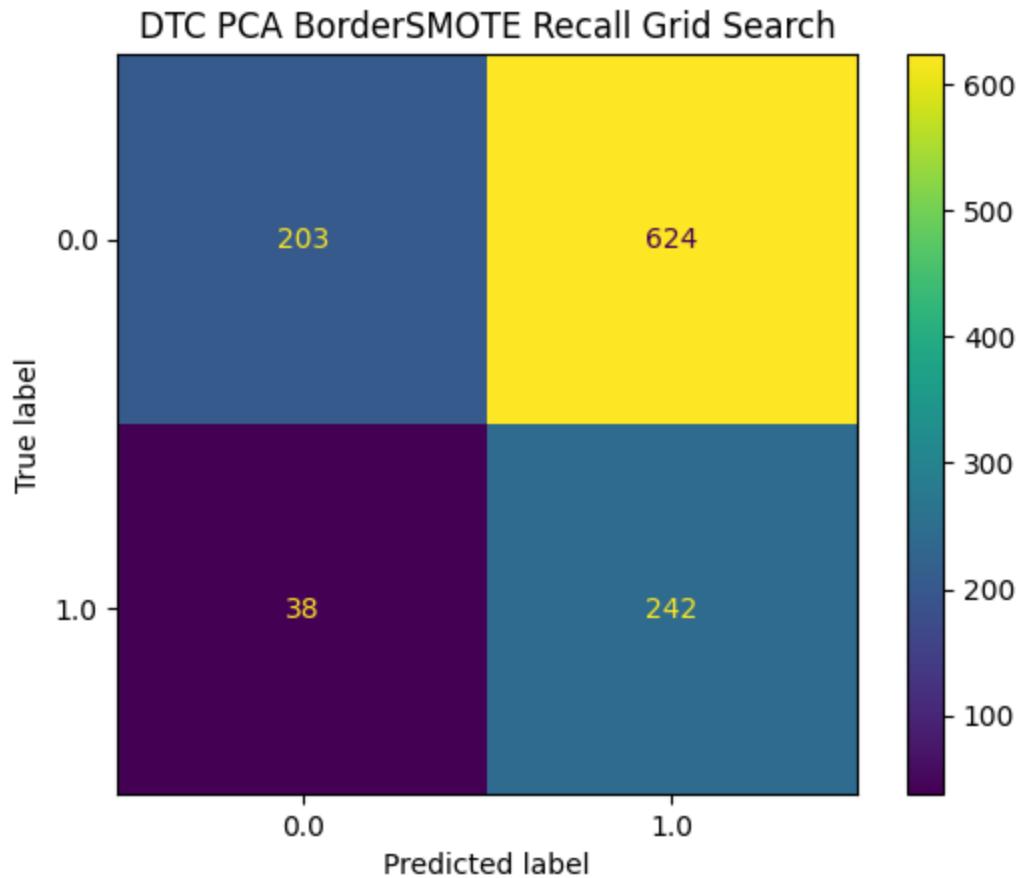
```
Out[41]: Pipeline(steps=[('pca', PCA(n_components=10)),
                         ('sm', BorderlineSMOTE(sampling_strategy=0.7)),
                         ('tree_clf',
                           DecisionTreeClassifier(class_weight={0: 1, 1: 2}, max_depth=3,
                                                 min_samples_leaf=2,
                                                 min_samples_split=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 [38]: 1 y_pred = best_tree.predict(X_test)
```

```
In [42]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('DTC PCA BorderSMOTE Recall Grid Search')
3 plt.show()
4 print(classification_report(y_test, y_pred))
```



	precision	recall	f1-score	support
0.0	0.84	0.25	0.38	827
1.0	0.28	0.86	0.42	280
accuracy			0.40	1107
macro avg	0.56	0.55	0.40	1107
weighted avg	0.70	0.40	0.39	1107

cond\_pivoted\_df, Decision Tree Classifier

In [89]: 1 cond\_pivoted\_df.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3688 entries, 0 to 3687
Data columns (total 39 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   season          3688 non-null    int64  
 1   Age              3688 non-null    float64 
 2   Throws           3688 non-null    int64  
 3   Surgery          3688 non-null    float64 
 4   BB_release_speed_weighted_avg  3688 non-null    float64 
 5   FB_release_speed_weighted_avg  3688 non-null    float64 
 6   OS_release_speed_weighted_avg  3688 non-null    float64 
 7   OT_release_speed_weighted_avg  3688 non-null    float64 
 8   SB_release_speed_weighted_avg  3688 non-null    float64 
 9   BB_release_pos_x_weighted_avg 3688 non-null    float64 
 10  FB_release_pos_x_weighted_avg 3688 non-null    float64 
 11  OS_release_pos_x_weighted_avg 3688 non-null    float64 
 12  OT_release_pos_x_weighted_avg 3688 non-null    float64 
 13  SB_release_pos_x_weighted_avg 3688 non-null    float64 
 14  BB_release_pos_y_weighted_avg 3688 non-null    float64 
 15  FB_release_pos_y_weighted_avg 3688 non-null    float64 
 16  OS_release_pos_y_weighted_avg 3688 non-null    float64 
 17  OT_release_pos_y_weighted_avg 3688 non-null    float64 
 18  SB_release_pos_y_weighted_avg 3688 non-null    float64 
 19  BB_release_pos_z_weighted_avg 3688 non-null    float64 
 20  FB_release_pos_z_weighted_avg 3688 non-null    float64 
 21  OS_release_pos_z_weighted_avg 3688 non-null    float64 
 22  OT_release_pos_z_weighted_avg 3688 non-null    float64 
 23  SB_release_pos_z_weighted_avg 3688 non-null    float64 
 24  BB_vx0_weighted_avg         3688 non-null    float64 
 25  FB_vx0_weighted_avg         3688 non-null    float64 
 26  OS_vx0_weighted_avg         3688 non-null    float64 
 27  OT_vx0_weighted_avg         3688 non-null    float64 
 28  SB_vx0_weighted_avg         3688 non-null    float64 
 29  BB_vy0_weighted_avg         3688 non-null    float64 
 30  FB_vy0_weighted_avg         3688 non-null    float64 
 31  OS_vy0_weighted_avg         3688 non-null    float64 
 32  OT_vy0_weighted_avg         3688 non-null    float64 
 33  SB_vy0_weighted_avg         3688 non-null    float64 
 34  BB_vz0_weighted_avg         3688 non-null    float64 
 35  FB_vz0_weighted_avg         3688 non-null    float64 
 36  OS_vz0_weighted_avg         3688 non-null    float64 
 37  OT_vz0_weighted_avg         3688 non-null    float64 
 38  SB_vz0_weighted_avg         3688 non-null    float64 
dtypes: float64(37), int64(2)
memory usage: 1.1 MB
```

In [95]: 1 y = cond\_pivoted\_df['Surgery']

2 X = cond\_pivoted\_df.drop('Surgery', axis=1)

```
In [96]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

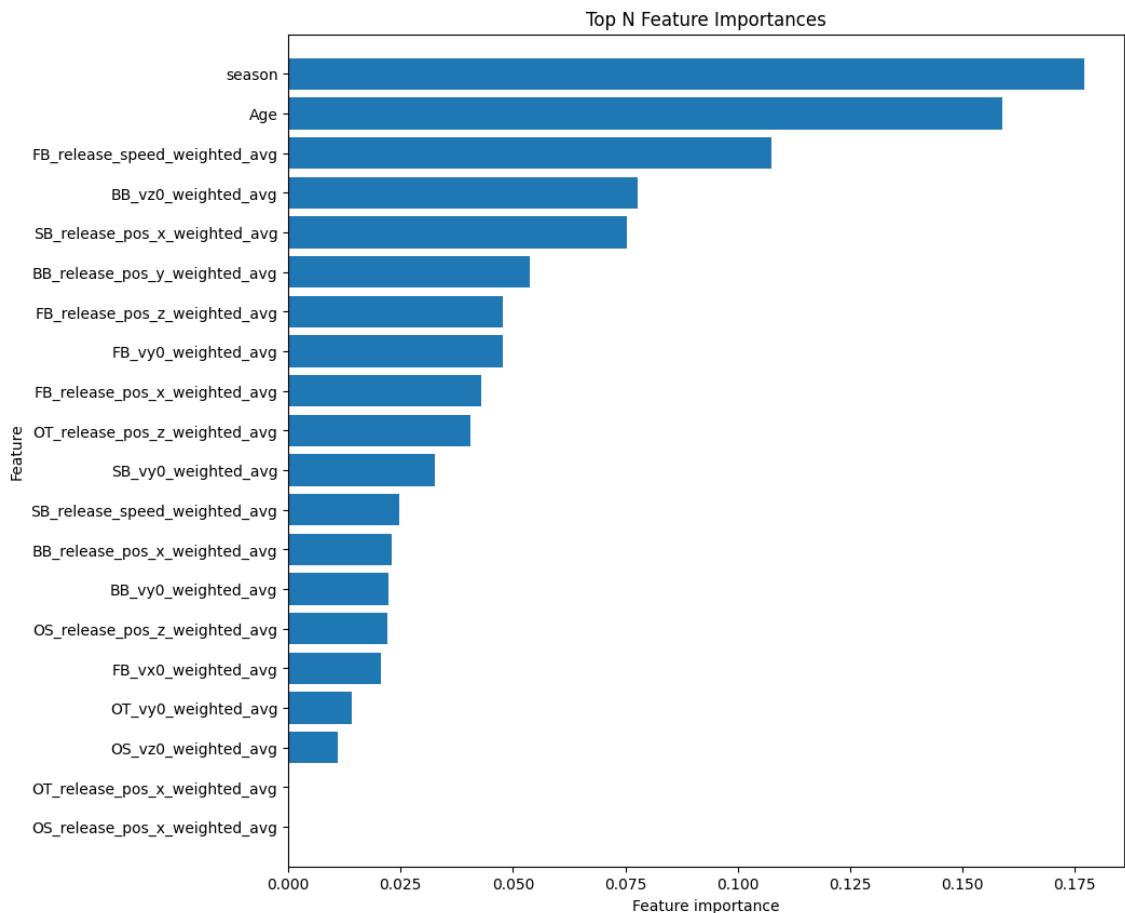
```
In [92]: 1 tree_clf = DecisionTreeClassifier(criterion='gini', max_depth=5)
2 tree_clf.fit(X_train, y_train)
```

Out[92]: `DecisionTreeClassifier(max_depth=5)`

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

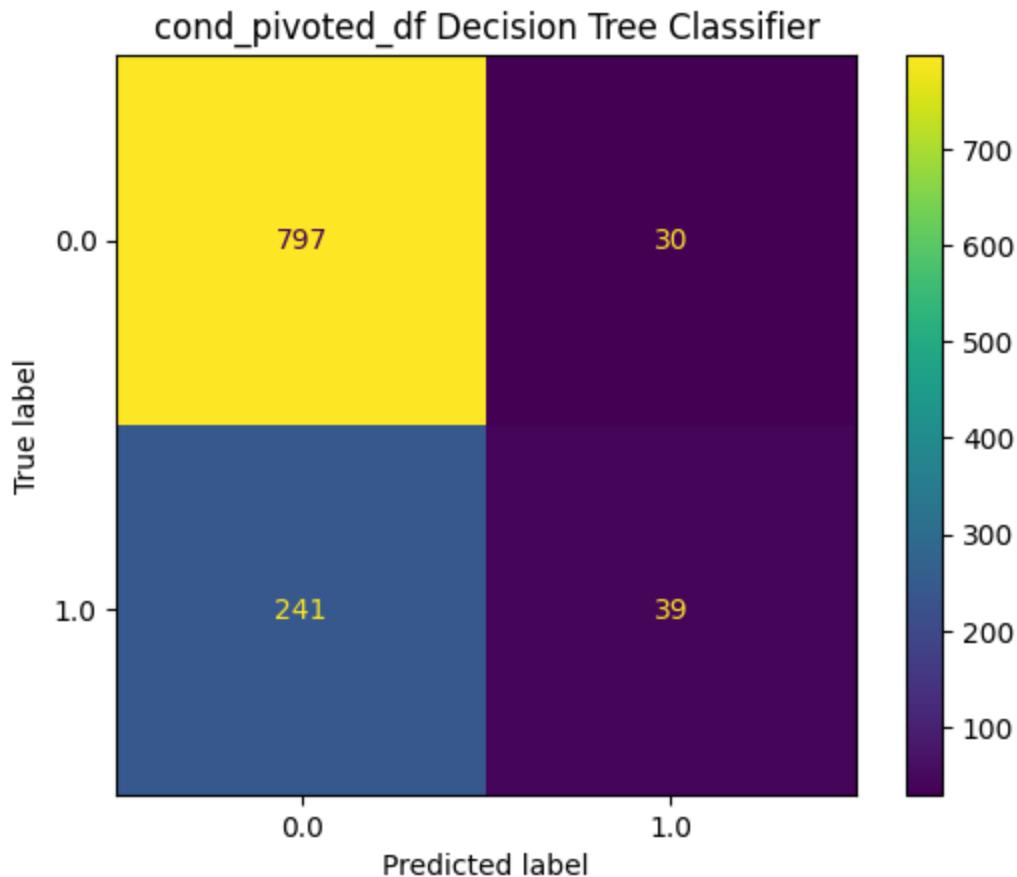
On GitHub, the HTML representation is unable to render, please try loading this page with [nbviewer.org](http://nbviewer.org).

```
In [93]: 1 plot_feature_importances(tree_clf, n_top_features=20)
2 plt.show()
```



In [94]:

```
1 pred = tree_clf.predict(X_test)
2
3 ConfusionMatrixDisplay.from_predictions(y_test, pred)
4 plt.title('cond_pivoted_df Decision Tree Classifier')
5 plt.show()
6 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.77	0.96	0.85	827
1.0	0.57	0.14	0.22	280
accuracy			0.76	1107
macro avg	0.67	0.55	0.54	1107
weighted avg	0.72	0.76	0.70	1107

cond\_pivoted\_df not predicting TP & FP well.

```
In [97]: 1 param_grid = {
2     'criterion': ['gini', 'entropy'],
3     'max_depth': [None, 5, 10, 15, 20],
4     'min_samples_split': [2, 5, 10],
5     'min_samples_leaf': [1, 2, 4],
6     'class_weight': [None, '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
13 # Fit the grid search to the data
14 grid_search.fit(X_train, y_train)
15
16 # Print the best parameters and the best score
17 print("Best parameters:", grid_search.best_params_)
18 print("Best score:", grid_search.best_score_)
19
20 # Evaluate the best model found by GridSearchCV on the test set
21 best_tree = grid_search.best_estimator_
22 y_pred = best_tree.predict(X_test)
23 print("Test set accuracy:", accuracy_score(y_test, y_pred))
24
```

```
Best parameters: {'class_weight': 'balanced', 'criterion': 'entropy', 'max_depth': 5, 'min_samples_leaf': 4, 'min_samples_split': 10}
Best score: 0.7031373031496063
Test set accuracy: 0.6621499548328816
```

```
In [98]: 1 y = cond_pivoted_df['Surgery']
2 X = cond_pivoted_df.drop('Surgery', axis=1)
```

```
In [99]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

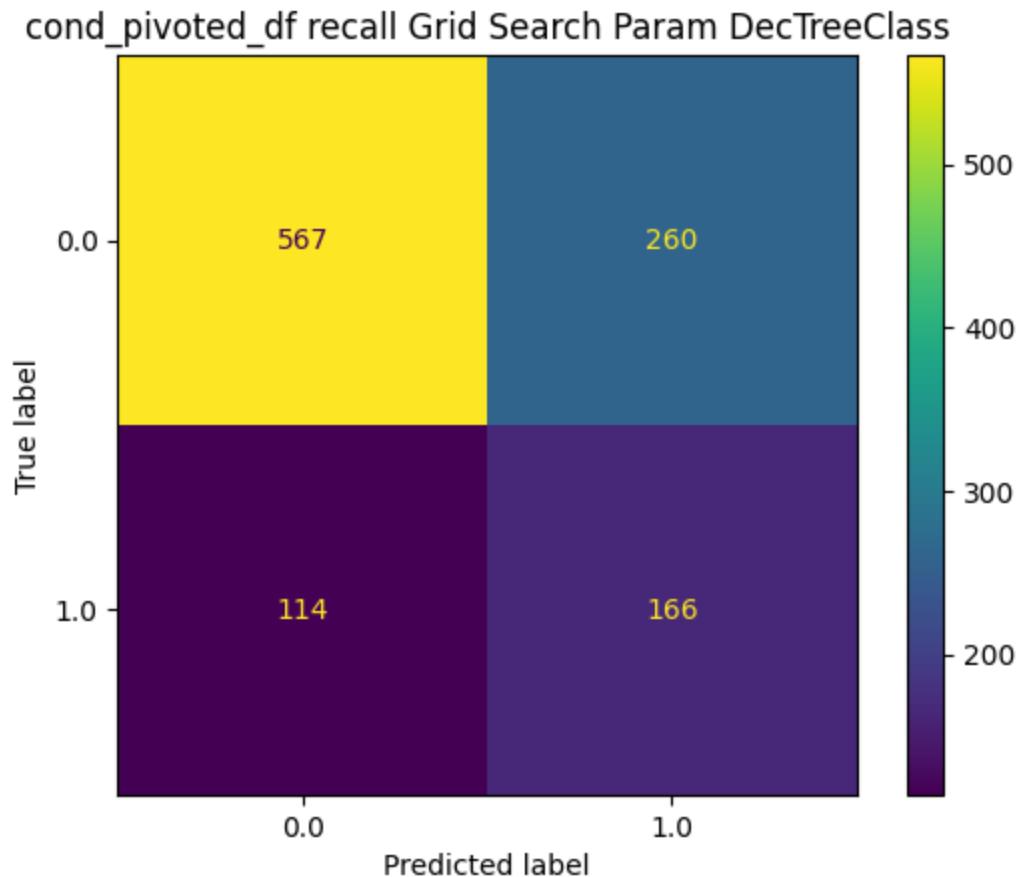
```
In [100]: 1 tree_clf = DecisionTreeClassifier(class_weight='balanced', criterion='entropy')
2 tree_clf.fit(X_train, y_train)
```

```
Out[100]: DecisionTreeClassifier(class_weight='balanced', criterion='entropy',
                                  max_depth=5, min_samples_leaf=4, 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 [101]: pred = tree_clf.predict(X_test)
2
3 ConfusionMatrixDisplay.from_predictions(y_test, pred)
4 plt.title('cond_pivoted_df recall Grid Search Param DecTreeClass')
5 plt.show()
6 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.83	0.69	0.75	827
1.0	0.39	0.59	0.47	280
accuracy			0.66	1107
macro avg	0.61	0.64	0.61	1107
weighted avg	0.72	0.66	0.68	1107

Better, but FN is high compared to other models.

pivoted\_df, Random Forest

```
In [102]: y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [103]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [104]: 1 forest_clf = RandomForestClassifier(n_estimators=100, max_depth=5)
2 forest_clf.fit(X_train, y_train)
```

Out[104]: RandomForestClassifier(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](#).

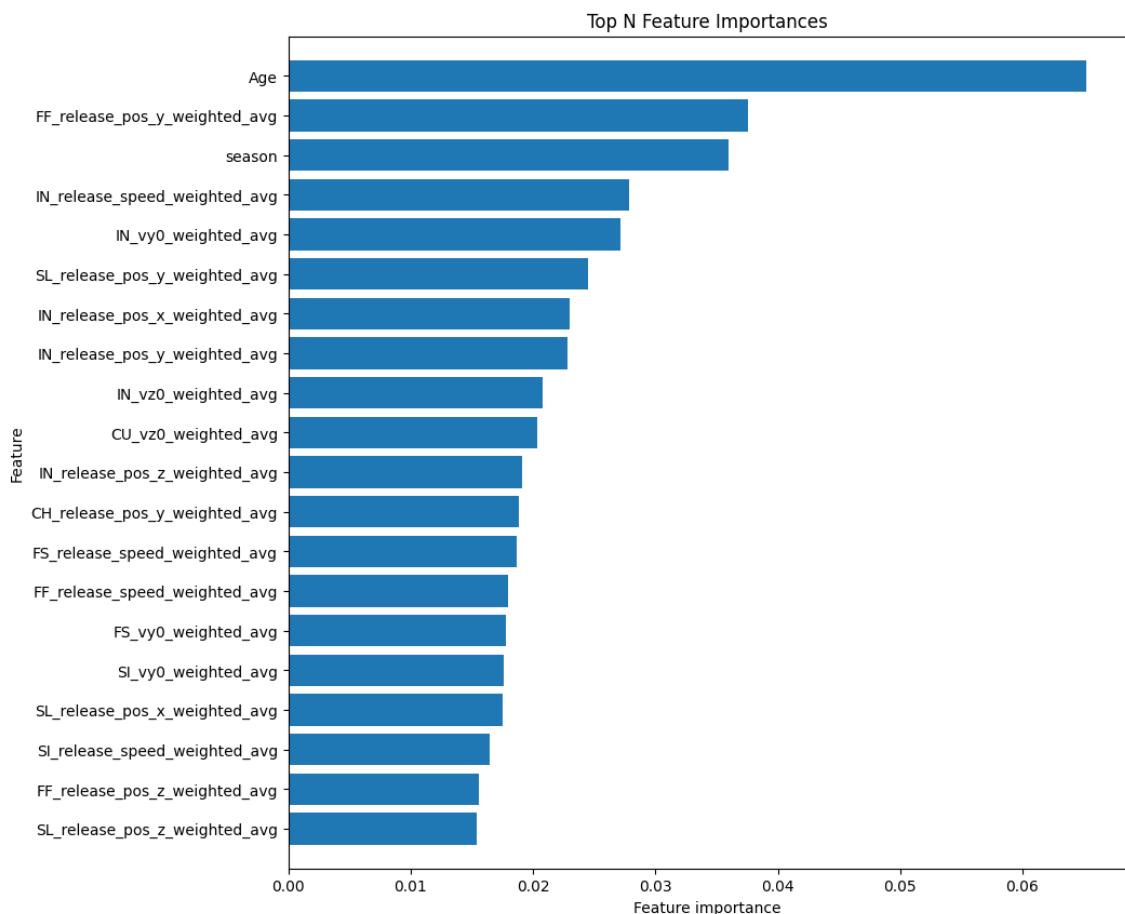
```
In [105]: 1 forest_clf.score(X_train, y_train)
```

Out[105]: 0.7667570709027509

```
In [106]: 1 forest_clf.score(X_test, y_test)
```

Out[106]: 0.7515808491418248

```
In [107]: 1 plot_feature_importances(forest_clf, n_top_features=20)
2 plt.show()
```



```
In [108]: 1 forest2_clf = RandomForestClassifier(n_estimators=5, max_features=10,  
2 forest2_clf.fit(X_train, y_train)
```

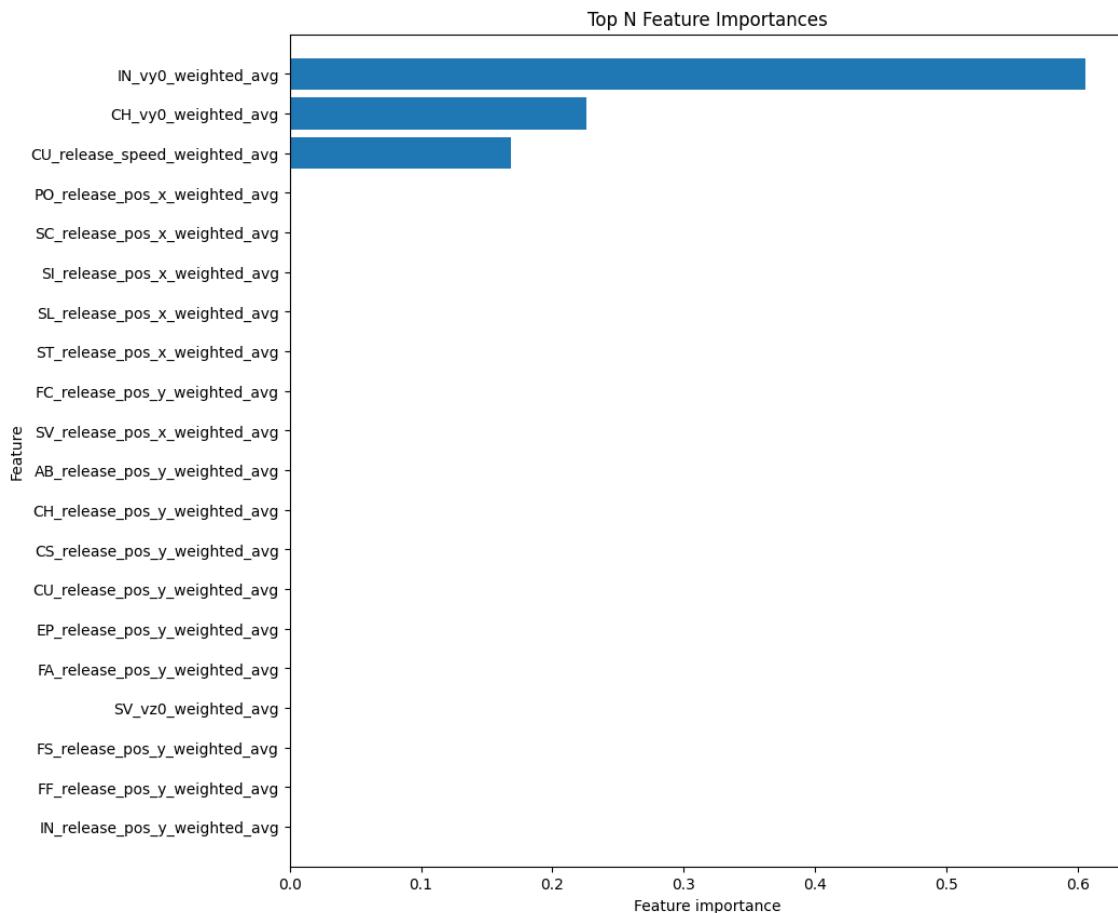
Out[108]: RandomForestClassifier(max\_depth=2, max\_features=10, n\_estimators=5)

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](#).

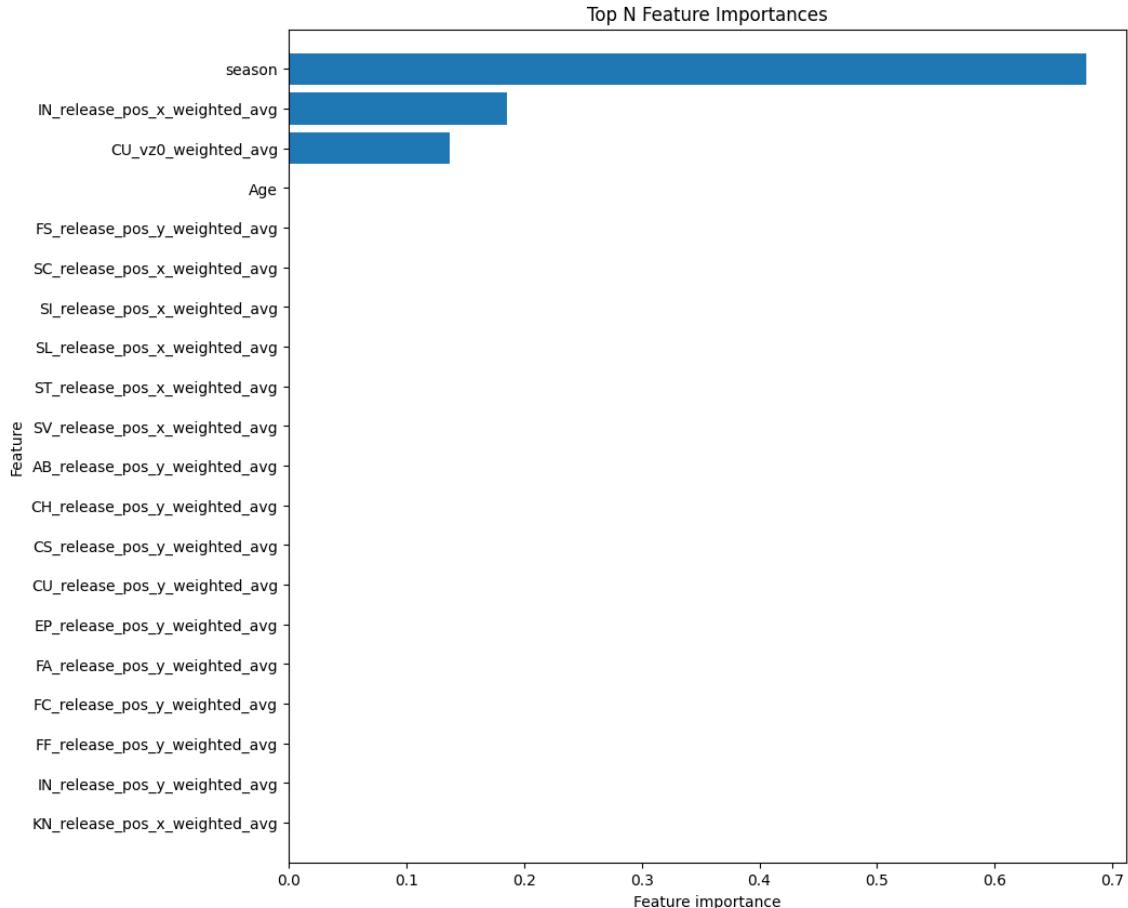
```
In [109]: 1 rf_tree_1 = forest2_clf.estimators_[0]
```

```
In [110]: 1 plot_feature_importances(rf_tree_1, n_top_features=20)  
2 plt.show()
```



```
In [111]: 1 rf_tree_2 = forest2_clf.estimators_[1]
```

```
In [112]: 1 plot_feature_importances(rf_tree_2, n_top_features=20)
2 plt.show()
```



```
In [43]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [44]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [46]: 1 param_grid_rf = {  
2     'n_estimators': [100, 200, 300],  
3     'criterion': ['gini', 'entropy'],  
4     'max_depth': [None, 5, 10, 15],  
5     'min_samples_split': [2, 5, 10],  
6     'min_samples_leaf': [1, 2, 4],  
7     'class_weight': [None, 'balanced', {0:1, 1:2}, {0:1, 1:3}]  
8 }  
9  
10 forest_clf = RandomForestClassifier()  
11 accuracy_scorer = make_scorer(recall_score)  
12 grid_search_rf = GridSearchCV(estimator=forest_clf, param_grid=param_g  
13 grid_search_rf.fit(X_train, y_train)  
14  
15 print("Best parameters:", grid_search_rf.best_params_)  
16 print("Best accuracy score:", grid_search_rf.best_score_)  
17  
18 best_rf = grid_search_rf.best_estimator_  
19 y_pred_rf = best_rf.predict(X_test)  
20  
21 print("Test set recall score:", recall_score(y_test, y_pred_rf))
```

```
Best parameters: {'class_weight': 'balanced', 'criterion': 'entropy', 'ma  
x_depth': 5, 'min_samples_leaf': 4, 'min_samples_split': 2, 'n_estimators': 100}  
Best accuracy score: 0.6069635826771653  
Test set recall score: 0.6071428571428571
```

```
In [47]: 1 #forest_clf = RandomForestClassifier(class_weight='balanced', criterio  
2 best_rf.fit(X_train, y_train)
```

Out[47]: RandomForestClassifier(class\_weight='balanced', criterion='entropy',  
max\_depth=5, min\_samples\_leaf=4)

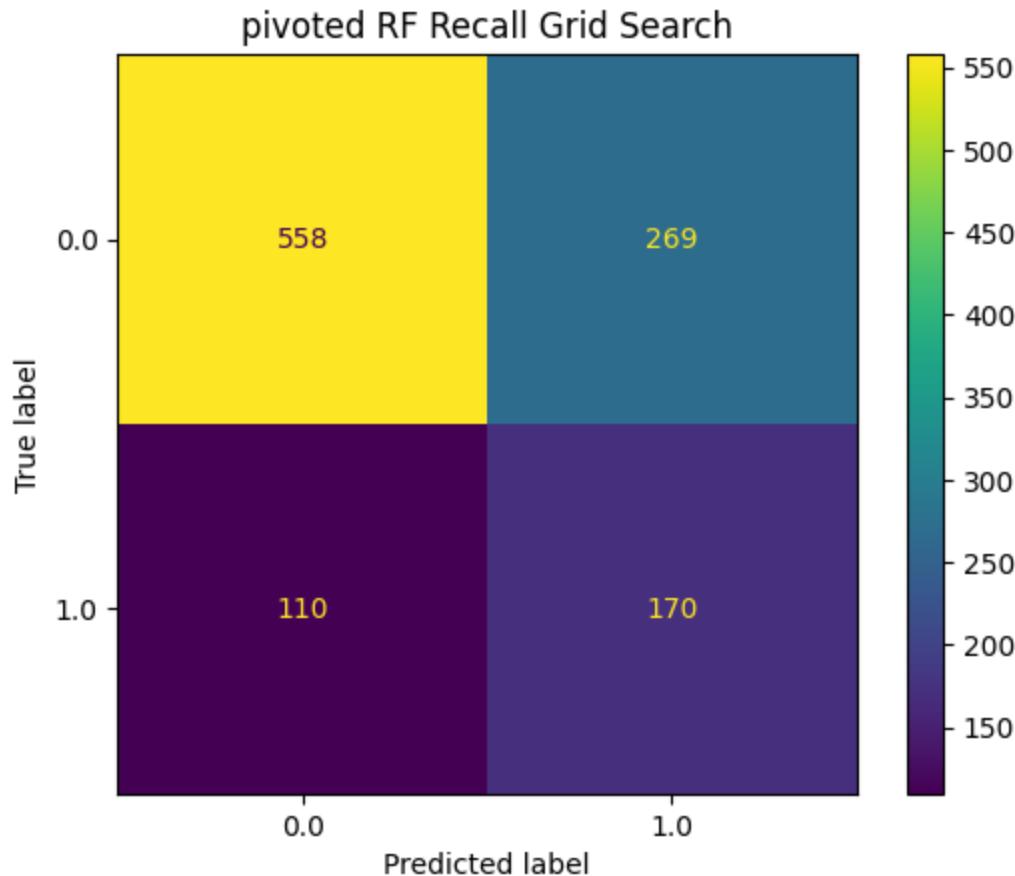
**In a Jupyter environment, please rerun this cell to show the HTML representation or  
trust the notebook.**

**On GitHub, the HTML representation is unable to render, please try loading this page  
with nbviewer.org.**

```
In [48]: 1 pred = best_rf.predict(X_test)
```

In [49]:

```
1 ConfusionMatrixDisplay.from_predictions(y_test, pred)
2 plt.title('pivoted RF Recall Grid Search')
3 plt.show()
4 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.84	0.67	0.75	827
1.0	0.39	0.61	0.47	280
accuracy			0.66	1107
macro avg	0.61	0.64	0.61	1107
weighted avg	0.72	0.66	0.68	1107

Again, TN & TP are good, but FP is too small compared to FN. Will try this with recall score.

```
In [132]: 1 param_grid_rf = {
2     'n_estimators': [100, 200, 300],
3     'criterion': ['gini', 'entropy'],
4     'max_depth': [None, 5, 10, 15],
5     'min_samples_split': [2, 5, 10],
6     'min_samples_leaf': [1, 2, 4],
7     'class_weight': [None, 'balanced', {0:1, 1:2}, {0:1, 1:3}]
8 }
9
10 forest_clf = RandomForestClassifier()
11 recall_scorer = make_scorer(recall_score)
12 grid_search_rf = GridSearchCV(estimator=forest_clf, param_grid=param_g
13 grid_search_rf.fit(X_train, y_train)
14
15 print("Best parameters:", grid_search_rf.best_params_)
16 print("Best accuracy score:", grid_search_rf.best_score_)
17
18 best_rf = grid_search_rf.best_estimator_
19 y_pred_rf = best_rf.predict(X_test)
20
21 print("Test set accuracy score:", recall_score(y_test, y_pred_rf))
```

Best parameters: {'class\_weight': 'balanced', 'criterion': 'entropy', 'max\_depth': None, 'min\_samples\_leaf': 2, 'min\_samples\_split': 10, 'n\_estimators': 100}  
 Best accuracy score: 0.810923935045657  
 Test set accuracy score: 0.34285714285714286

```
In [133]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [134]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
```

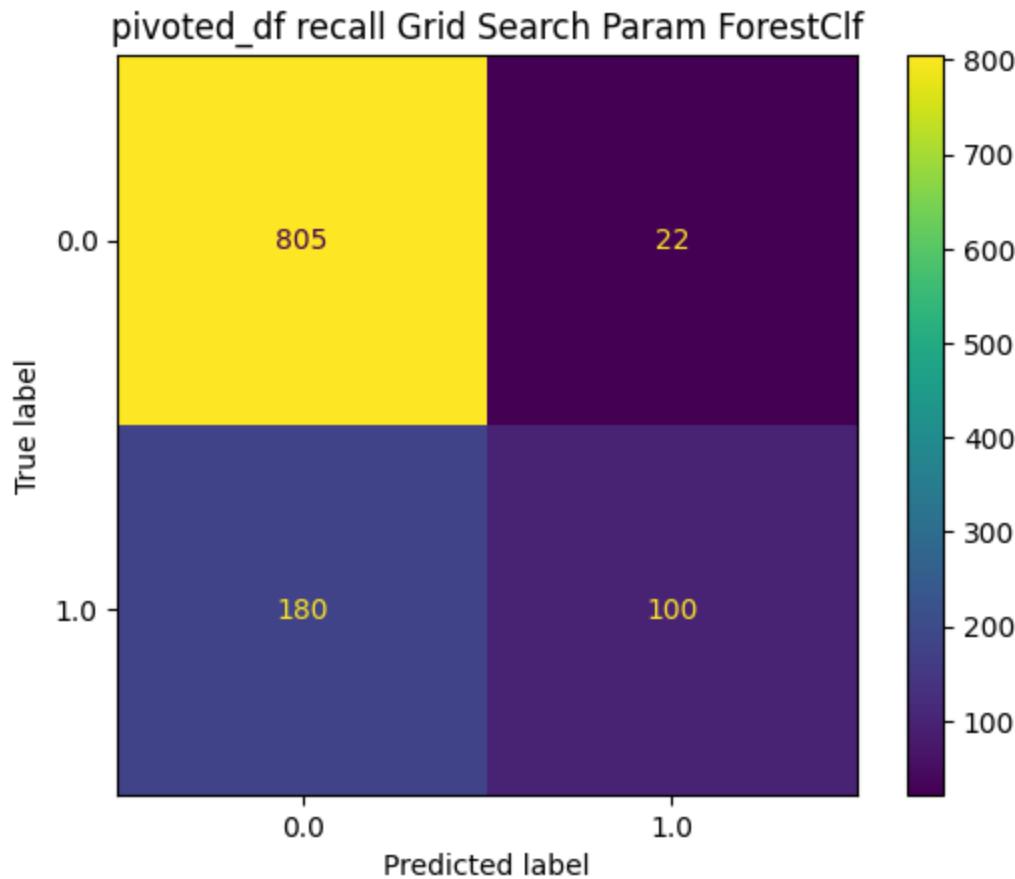
```
In [136]: 1 forest_clf = RandomForestClassifier(class_weight='balanced', criterion='entropy',
2 forest_clf.fit(X_train, y_train)
```

**Out[136]:** RandomForestClassifier(class\_weight='balanced', criterion='entropy',  
 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 [137]: pred = forest_clf.predict(X_test)
2
3 ConfusionMatrixDisplay.from_predictions(y_test, pred)
4 plt.title('pivoted_df recall Grid Search Param ForestClf')
5 plt.show()
6 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.82	0.97	0.89	827
1.0	0.82	0.36	0.50	280
accuracy			0.82	1107
macro avg	0.82	0.67	0.69	1107
weighted avg	0.82	0.82	0.79	1107

Not great. Would like to see lower FN and higher TP & FP.

cond\_pivoted\_df, Random Forest

```
In [138]: y = cond_pivoted_df['Surgery']
2 X = cond_pivoted_df.drop('Surgery', axis=1)
```

```
In [139]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [140]: 1 forest_clf = RandomForestClassifier(n_estimators=100, max_depth=5)
2 forest_clf.fit(X_train, y_train)
```

Out[140]: RandomForestClassifier(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](http://nbviewer.org).

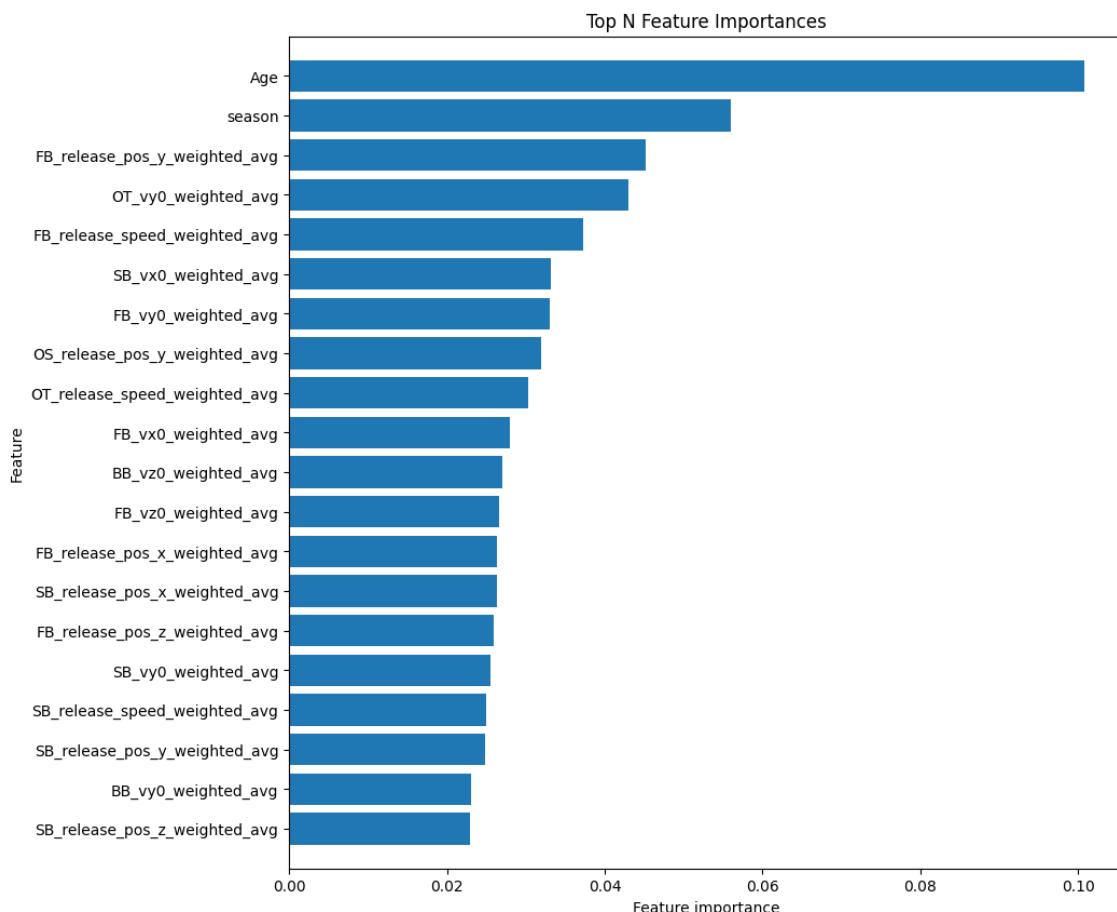
```
In [141]: 1 forest_clf.score(X_train, y_train)
```

Out[141]: 0.7729562185199536

```
In [142]: 1 forest_clf.score(X_test, y_test)
```

Out[142]: 0.7497741644083108

```
In [143]: 1 plot_feature_importances(forest_clf, n_top_features=20)
2 plt.show()
```



```
In [146]: 1 param_grid_rf = {
2     'n_estimators': [100, 200, 300],
3     'criterion': ['gini', 'entropy'],
4     'max_depth': [None, 5, 10, 15],
5     'min_samples_split': [2, 5, 10],
6     'min_samples_leaf': [1, 2, 4],
7     'class_weight': [None, 'balanced', {0:1, 1:2}, {0:1, 1:3}]
8 }
9
10 forest_clf = RandomForestClassifier()
11 accuracy_scorer = make_scorer(accuracy_score)
12 grid_search_rf = GridSearchCV(estimator=forest_clf, param_grid=param_g
13 grid_search_rf.fit(X_train, y_train)
14
15 print("Best parameters:", grid_search_rf.best_params_)
16 print("Best accuracy score:", grid_search_rf.best_score_)
17
18 best_rf = grid_search_rf.best_estimator_
19 y_pred_rf = best_rf.predict(X_test)
20
21 print("Test set accuracy score:", accuracy_score(y_test, y_pred_rf))
```

```
Best parameters: {'class_weight': 'balanced', 'criterion': 'entropy', 'ma
x_depth': 15, 'min_samples_leaf': 2, 'min_samples_split': 5, 'n_estimators': 200}
Best accuracy score: 0.8085953548348404
Test set accuracy score: 0.8121047877145439
```

```
In [147]: 1 y = cond_pivoted_df['Surgery']
2 X = cond_pivoted_df.drop('Surgery', axis=1)
```

```
In [148]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
```

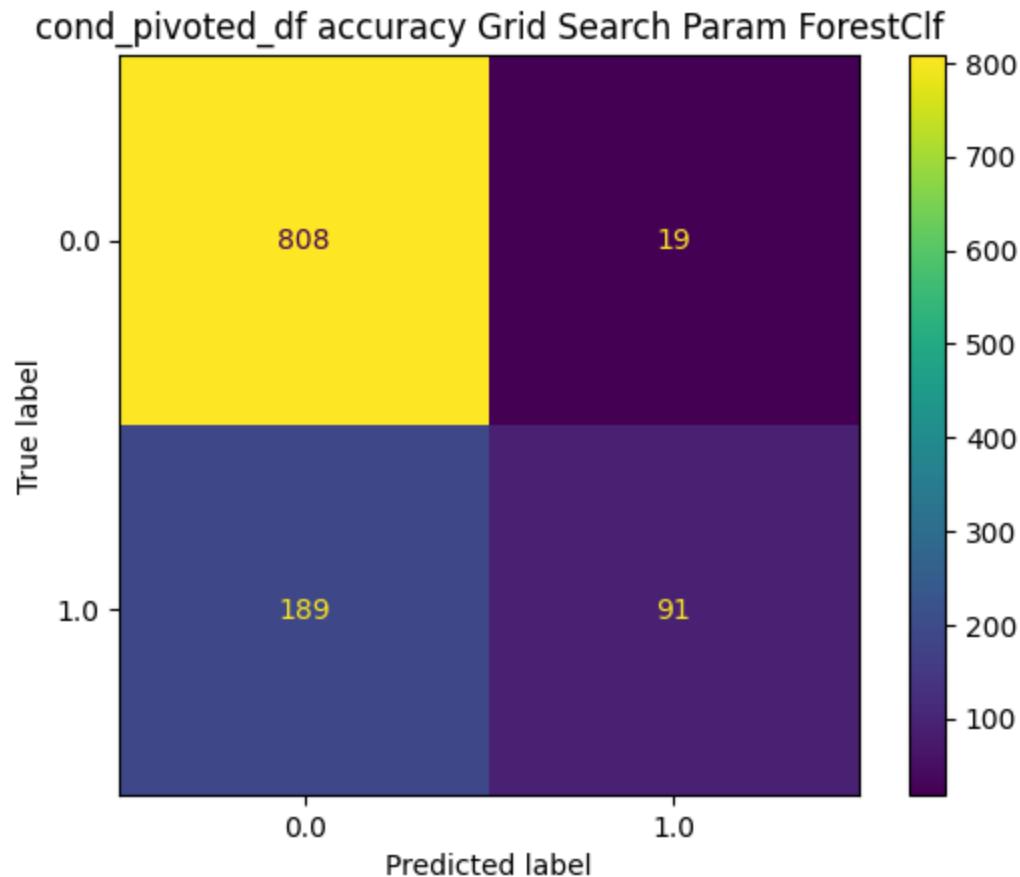
```
In [149]: 1 forest_clf = RandomForestClassifier(class_weight='balanced', criterio
2 forest_clf.fit(X_train, y_train)
```

```
Out[149]: RandomForestClassifier(class_weight='balanced', criterion='entropy',
                                 max_depth=15, min_samples_leaf=2, min_samples_spli
t=5,
                                 n_estimators=200)
```

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](http://nbviewer.org).

```
In [151]: pred = forest_clf.predict(X_test)
2
3 ConfusionMatrixDisplay.from_predictions(y_test, pred)
4 plt.title('cond_pivoted_df accuracy Grid Search Param ForestClf')
5 plt.show()
6 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.81	0.98	0.89	827
1.0	0.83	0.33	0.47	280
accuracy			0.81	1107
macro avg	0.82	0.65	0.68	1107
weighted avg	0.81	0.81	0.78	1107

Try recall score.

```
In [154]: 1 param_grid_rf = {
2     'n_estimators': [100, 200, 300],
3     'criterion': ['gini', 'entropy'],
4     'max_depth': [None, 5, 10, 15],
5     'min_samples_split': [2, 5, 10],
6     'min_samples_leaf': [1, 2, 4],
7     'class_weight': [None, 'balanced', {0:1, 1:2}, {0:1, 1:3}]
8 }
9
10 forest_clf = RandomForestClassifier()
11 recall_scorer = make_scorer(recall_score)
12 grid_search_rf = GridSearchCV(estimator=forest_clf, param_grid=param_g
13 grid_search_rf.fit(X_train, y_train)
14
15 print("Best parameters:", grid_search_rf.best_params_)
16 print("Best accuracy score:", grid_search_rf.best_score_)
17
18 best_rf = grid_search_rf.best_estimator_
19 y_pred_rf = best_rf.predict(X_test)
20
21 print("Test set accuracy score:", recall_score(y_test, y_pred_rf))
```

Best parameters: {'class\_weight': 'balanced', 'criterion': 'entropy', 'max\_depth': 15, 'min\_samples\_leaf': 2, 'min\_samples\_split': 5, 'n\_estimators': 100}  
 Best accuracy score: 0.8074348132487668  
 Test set accuracy score: 0.33214285714285713

```
In [155]: 1 y = cond_pivoted_df['Surgery']
2 X = cond_pivoted_df.drop('Surgery', axis=1)
```

```
In [156]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
```

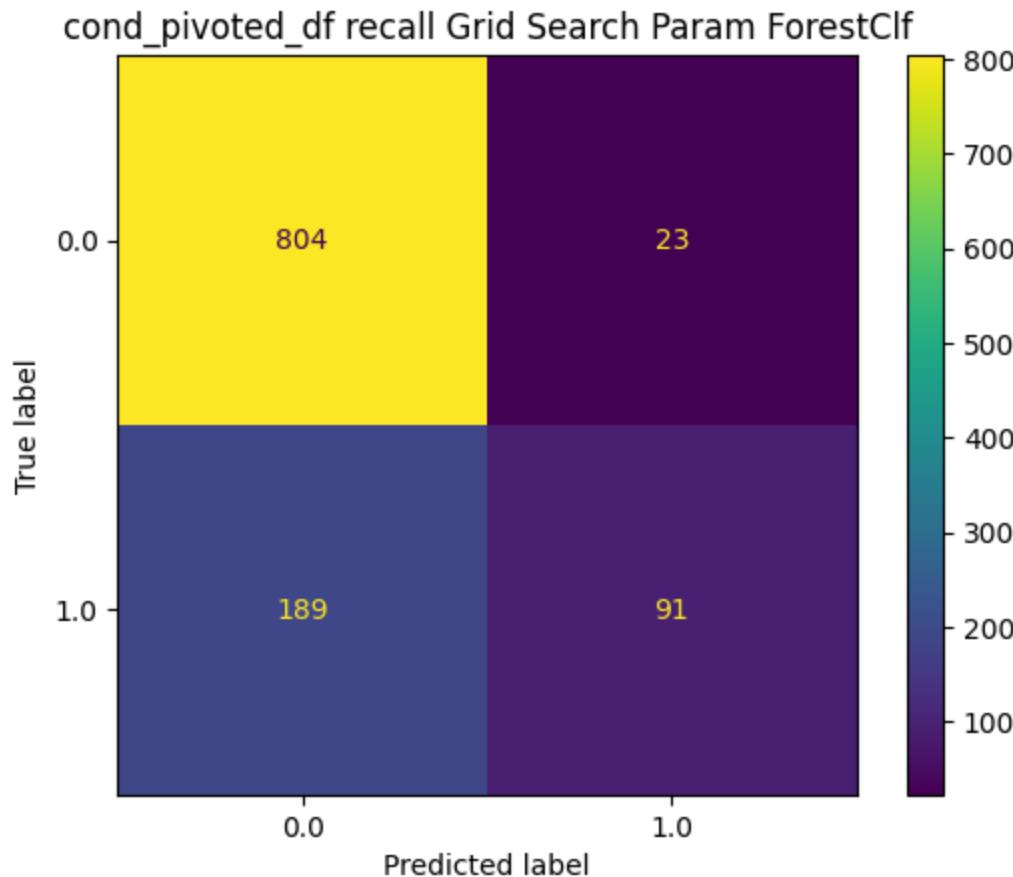
```
In [157]: 1 forest_clf = RandomForestClassifier(class_weight='balanced', criterion='
2 forest_clf.fit(X_train, y_train)
```

**Out[157]:** RandomForestClassifier(class\_weight='balanced', criterion='entropy',
max\_depth=15, min\_samples\_leaf=2, min\_samples\_split=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 [158]: pred = forest_clf.predict(X_test)
2
3 ConfusionMatrixDisplay.from_predictions(y_test, pred)
4 plt.title('cond_pivoted_df recall Grid Search Param ForestClf')
5 plt.show()
6 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.81	0.97	0.88	827
1.0	0.80	0.33	0.46	280
accuracy			0.81	1107
macro avg	0.80	0.65	0.67	1107
weighted avg	0.81	0.81	0.78	1107

Odd. Scores are very similar. Is there a way to better these scores? May just have to move onto another model. XG Boost?

pivoted\_df XG Boost

```
In [171]: y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [172]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [167]: 1 xgb_clf = XGBClassifier(use_label_encoder=False, eval_metric='logloss')
```

```
In [168]: 1 xgb_clf.fit(X_train, y_train)
```

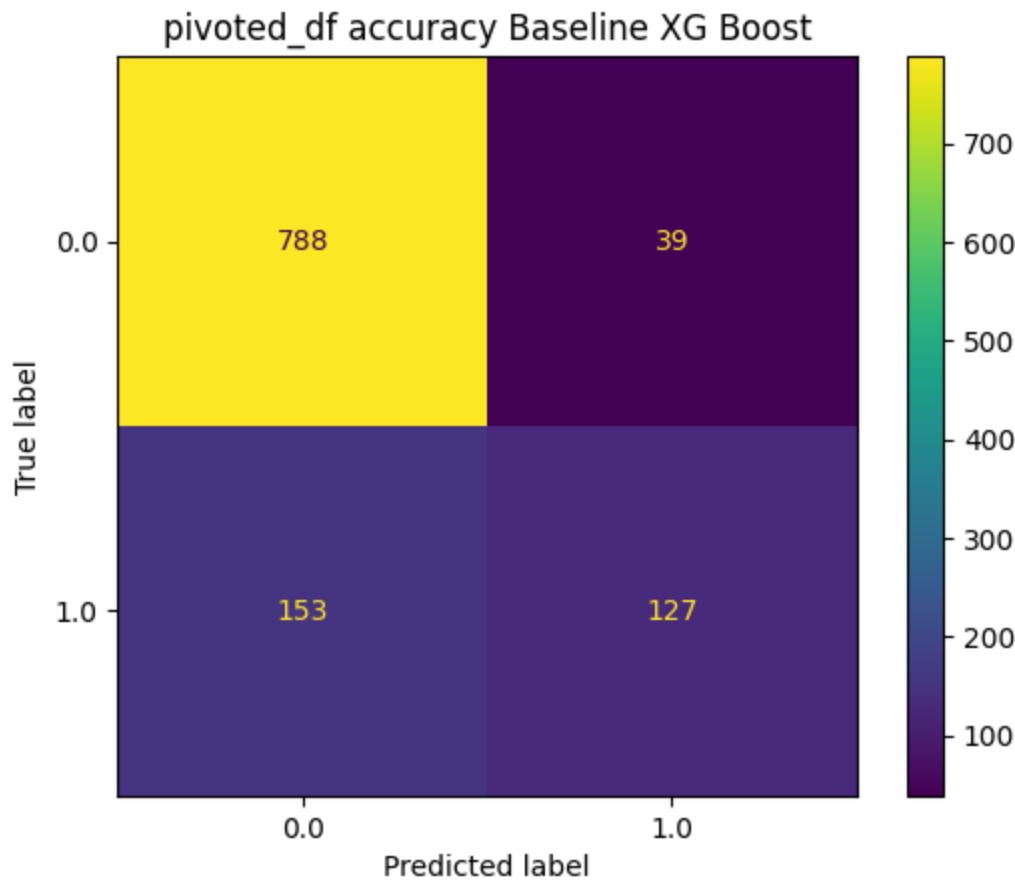
```
Out[168]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                        colsample_bylevel=None, colsample_bynode=None,
                        colsample_bytree=None, device=None, early_stopping_rounds=N
                        one,
                        enable_categorical=False, eval_metric='logloss',
                        feature_types=None, gamma=None, grow_policy=None,
                        importance_type=None, interaction_constraints=None,
                        learning_rate=None, max_bin=None, max_cat_threshold=None,
                        max_cat_to_onehot=None, max_delta_step=None, max_depth=None,
                        max_leaves=None, min_child_weight=None, missing=nan,
                        monotone_constraints=None, multi_strategy=None, n_estimators=None,
                        n_jobs=None, num_parallel_tree=None, random_state=None,
                        ...)
```

**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 [169]: 1 y_pred = xgb_clf.predict(X_test)
```

```
In [170]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('pivoted_df accuracy Baseline XG Boost')
3 plt.show()
4 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.81	0.97	0.88	827
1.0	0.80	0.33	0.46	280
accuracy			0.81	1107
macro avg	0.80	0.65	0.67	1107
weighted avg	0.81	0.81	0.78	1107

Pretty good. Would like to see lower FN and higher scores for everything else. Will try GridSearch.

```
In [183]: 1 param_grid_xgb = {
2     'n_estimators': [100, 200],
3     'max_depth': [3, 5, 7],
4     'learning_rate': [0.01, 0.1, 0.2],
5     'subsample': [0.8, 1],
6     'colsample_bytree': [0.8, 1]
7 }
8
9 xgb_clf = XGBClassifier(objective='binary:logistic', use_label_encoder=True)
10
11 grid_search_xgb = GridSearchCV(estimator=xgb_clf, param_grid=param_grid_xgb)
12 grid_search_xgb.fit(X_train, y_train)
13
14 print("Best parameters:", grid_search_xgb.best_params_)
15 print("Best recall score:", grid_search_xgb.best_score_)
16
17 best_xgb = grid_search_xgb.best_estimator_
18 y_pred_best_xgb = best_xgb.predict(X_test)
19 print("Test set recall score:", recall_score(y_test, y_pred_best_xgb))
20 print(classification_report(y_test, y_pred_best_xgb))
```

Best parameters: {'colsample\_bytree': 1, 'learning\_rate': 0.2, 'max\_depth': 7, 'n\_estimators': 200, 'subsample': 0.8}

Best recall score: 0.4323941929133858

Test set recall score: 0.45357142857142857

	precision	recall	f1-score	support
0.0	0.84	0.96	0.89	827
1.0	0.79	0.45	0.58	280
accuracy			0.83	1107
macro avg	0.81	0.71	0.74	1107
weighted avg	0.83	0.83	0.81	1107

```
In [184]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [185]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
In [186]: 1 xgb_clf = XGBClassifier(colsample_bytree=1, learning_rate=0.2, max_depth=7)
```

In [187]: 1 xgb\_clf.fit(X\_train, y\_train)

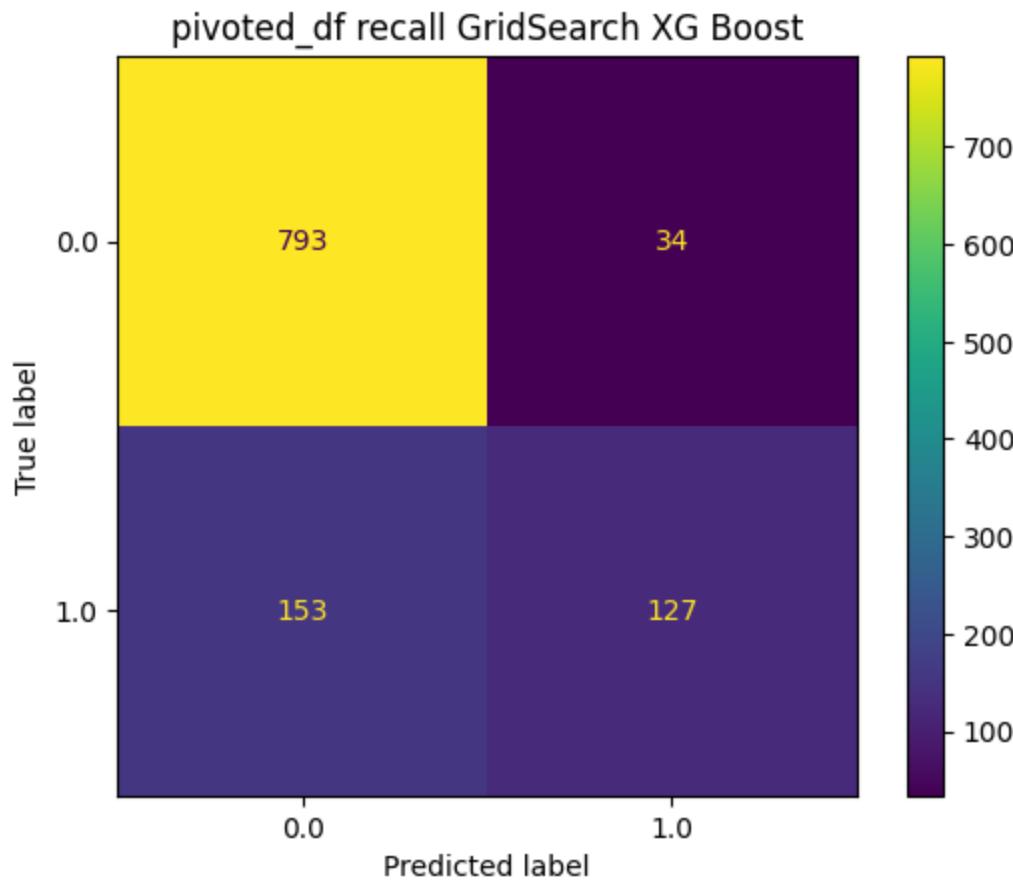
Out[187]: XGBClassifier(base\_score=None, booster=None, callbacks=None, colsample\_bylevel=None, colsample\_bynode=None, colsample\_bytree=1, device=None, early\_stopping\_rounds=None, enable\_categorical=False, eval\_metric='logloss', feature\_types=None, gamma=None, grow\_policy=None, importance\_type=None, interaction\_constraints=None, learning\_rate=0.2, max\_bin=None, max\_cat\_threshold=None, max\_cat\_to\_onehot=None, max\_delta\_step=None, max\_depth=7, max\_leaves=None, min\_child\_weight=None, missing=nan, monotone\_constraints=None, multi\_strategy=None, n\_estimators=200, n\_jobs=None, num\_parallel\_tree=None, random\_state=None, ...)

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 [188]: 1 y\_pred = xgb\_clf.predict(X\_test)

```
In [190]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('pivoted_df recall GridSearch XG Boost')
3 plt.show()
4 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.81	0.97	0.88	827
1.0	0.80	0.33	0.46	280
accuracy			0.81	1107
macro avg	0.80	0.65	0.67	1107
weighted avg	0.81	0.81	0.78	1107

Essentially the same score as the previous model. Will try with cond\_pivoted\_df now.

cond\_pivoted\_df, XG Boost

```
In [191]: 1 y = cond_pivoted_df['Surgery']
2 X = cond_pivoted_df.drop('Surgery', axis=1)
```

```
In [192]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.)
```

```
In [193]: 1 xgb_clf = XGBClassifier(use_label_encoder=False, eval_metric='logloss')
```

```
In [194]: 1 xgb_clf.fit(X_train, y_train)
```

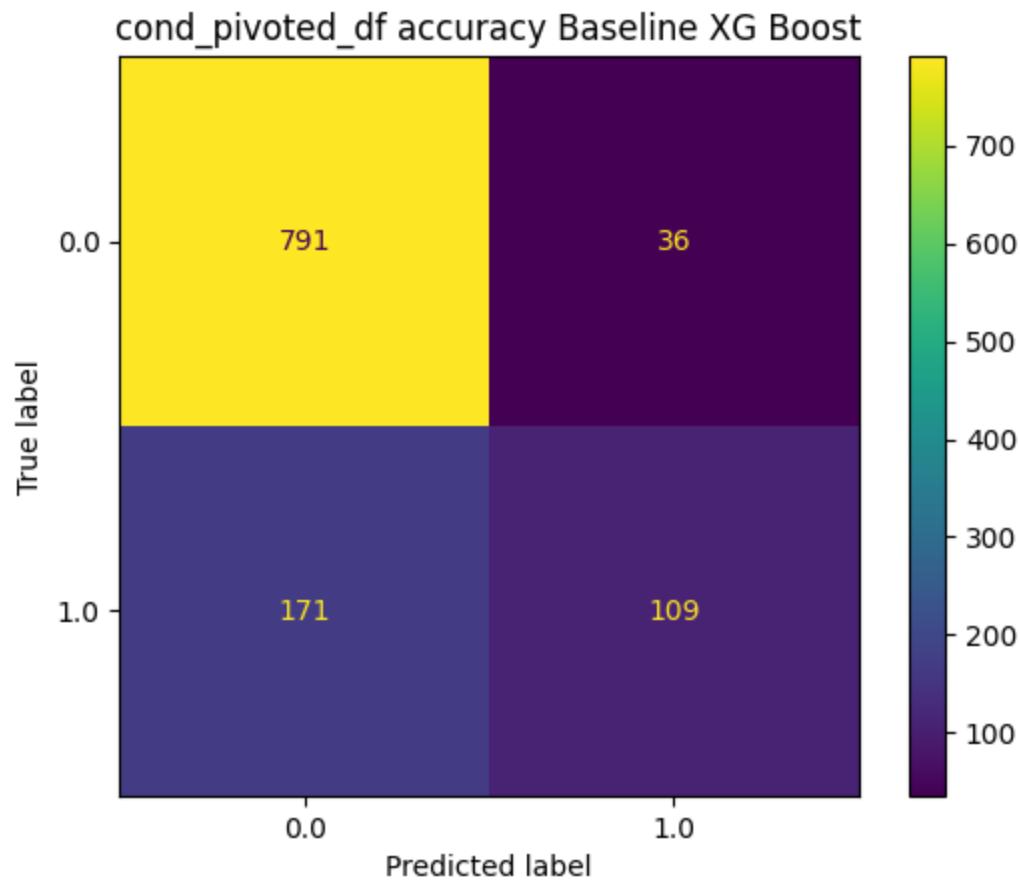
```
Out[194]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                        colsample_bylevel=None, colsample_bynode=None,
                        colsample_bytree=None, device=None, early_stopping_rounds=None,
                        enable_categorical=False, eval_metric='logloss',
                        feature_types=None, gamma=None, grow_policy=None,
                        importance_type=None, interaction_constraints=None,
                        learning_rate=None, max_bin=None, max_cat_threshold=None,
                        max_cat_to_onehot=None, max_delta_step=None, max_depth=None,
                        max_leaves=None, min_child_weight=None, missing=nan,
                        monotone_constraints=None, multi_strategy=None, n_estimators=None,
                        n_jobs=None, num_parallel_tree=None, random_state=None,
                        ...)
```

**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 [195]: 1 y_pred = xgb_clf.predict(X_test)
```

```
In [196]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('cond_pivoted_df accuracy Baseline XG Boost')
3 plt.show()
4 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.81	0.97	0.88	827
1.0	0.80	0.33	0.46	280
accuracy			0.81	1107
macro avg	0.80	0.65	0.67	1107
weighted avg	0.81	0.81	0.78	1107

```
In [202]: 1 param_grid_xgb = {
2     'n_estimators': [100, 200],
3     'max_depth': [3, 5, 7],
4     'learning_rate': [0.01, 0.1, 0.2],
5     'subsample': [0.8, 1],
6     'colsample_bytree': [0.8, 1]
7 }
8
9 cond_xgb_clf = XGBClassifier(objective='binary:logistic', use_label_encoder=True)
10
11 grid_search_xgb = GridSearchCV(estimator=cond_xgb_clf, param_grid=param_grid_xgb)
12 grid_search_xgb.fit(X_train, y_train)
13
14 print("Best parameters:", grid_search_xgb.best_params_)
15 print("Best recall score:", grid_search_xgb.best_score_)
16
17 best_xgb = grid_search_xgb.best_estimator_
18 y_pred_best_xgb = best_xgb.predict(X_test)
19 print("Test set recall score:", recall_score(y_test, y_pred_best_xgb))
20 print(classification_report(y_test, y_pred_best_xgb))
```

Best parameters: {'colsample\_bytree': 1, 'learning\_rate': 0.2, 'max\_depth': 7, 'n\_estimators': 200, 'subsample': 0.8}

Best recall score: 0.4323941929133858

Test set recall score: 0.45357142857142857

	precision	recall	f1-score	support
0.0	0.84	0.96	0.89	827
1.0	0.79	0.45	0.58	280
accuracy			0.83	1107
macro avg	0.81	0.71	0.74	1107
weighted avg	0.83	0.83	0.81	1107

```
In [203]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [204]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
In [205]: 1 cond_xgb_clf = XGBClassifier(colsample_bytree=1, learning_rate=0.2, max_depth=7)
```

In [206]: 1 cond\_xgb\_clf.fit(X\_train, y\_train)

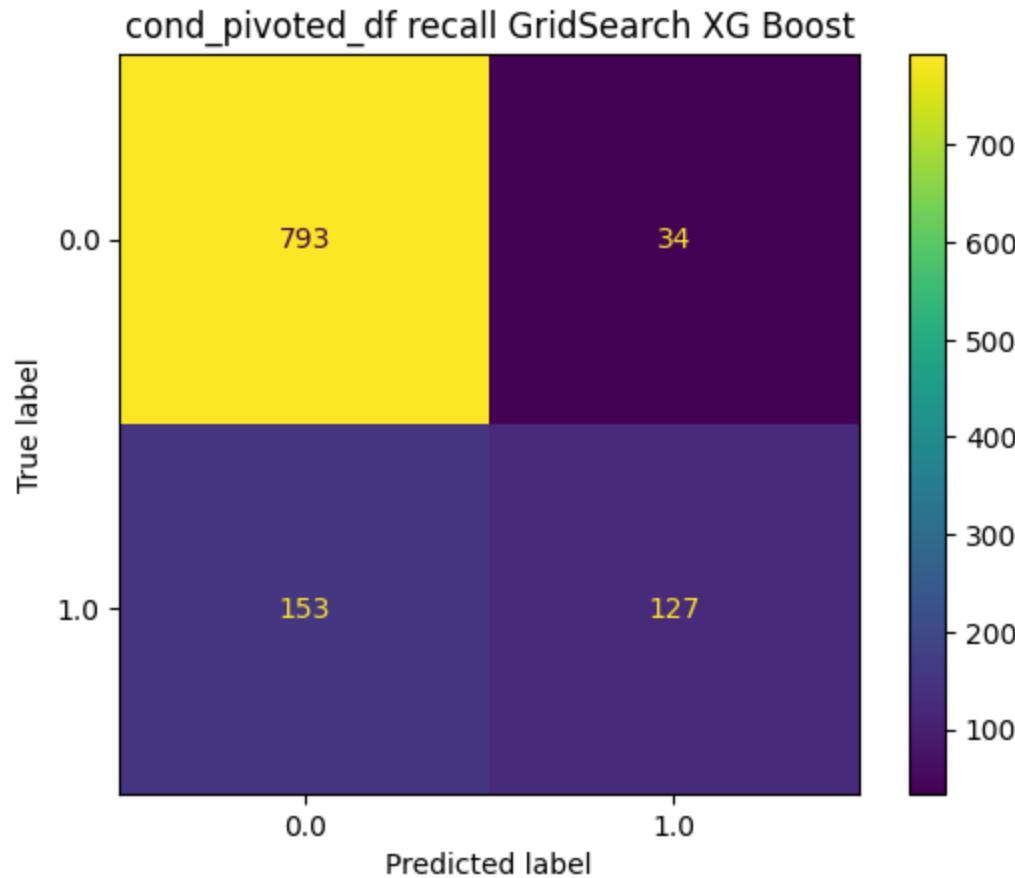
Out[206]: XGBClassifier(base\_score=None, booster=None, callbacks=None,  
colsample\_bylevel=None, colsample\_bynode=None, colsample\_bytree=1,  
device=None, early\_stopping\_rounds=None, enable\_categorical=False,  
eval\_metric='logloss', feature\_types=None, gamma=None,  
grow\_policy=None, importance\_type=None,  
interaction\_constraints=None, learning\_rate=0.2, max\_bin=None,  
max\_cat\_threshold=None, max\_cat\_to\_onehot=None,  
max\_delta\_step=None, max\_depth=7, max\_leaves=None,  
min\_child\_weight=None, missing=nan, monotone\_constraints=None,  
multi\_strategy=None, n\_estimators=200, n\_jobs=None,  
num\_parallel\_tree=None, random\_state=None, ...)

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 [207]: 1 y\_pred = cond\_xgb\_clf.predict(X\_test)

```
In [208]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('cond_pivoted_df recall GridSearch XG Boost')
3 plt.show()
4 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.81	0.97	0.88	827
1.0	0.80	0.33	0.46	280
accuracy			0.81	1107
macro avg	0.80	0.65	0.67	1107
weighted avg	0.81	0.81	0.78	1107

A little bit better. Will have to look into how to better tune all the various models parameters tomorrow.

Try to reduce dimensionality. Will work with Clustering and K Means on pivoted\_df.

```
In [77]: 1 y = pivoted_df['Surgery']
2 X = pivoted_df.drop('Surgery', axis=1)
```

```
In [78]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
```

```
In [79]: 1 def cross_validation(X_train, y_train, k, num_split = 5):
2
3     X_train = X_train.values
4     y_train = y_train.values
5
6     recall_train_list = []
7     recall_val_list = []
8
9     for train_index, valid_index in KFold(n_splits = num_split, shuffle=True):
10
11         # train and validation splitting
12         X_train_fold, X_val_fold = X_train[train_index], X_train[valid_index]
13         y_train_fold, y_val_fold = y_train[train_index], y_train[valid_index]
14
15         #create/fit the Standard scaler on the train fold
16         scaler = StandardScaler()
17         X_train_fold_scaled = scaler.fit_transform(X_train_fold)
18         # transform validation fold
19         X_val_fold_scaled = scaler.transform(X_val_fold)
20
21         # create/fit knearest neighbor
22         knn = KNeighborsClassifier(n_neighbors = k)
23         knn.fit(X_train_fold_scaled, y_train_fold)
24
25         # now how did we do?
26         y_train_pred = knn.predict(X_train_fold_scaled)
27         y_val_pred = knn.predict(X_val_fold_scaled)
28
29         recall_train = recall_score(y_train_fold, y_train_pred)
30         recall_val = recall_score(y_val_fold, y_val_pred)
31
32         recall_train_list.append(recall_train)
33         recall_val_list.append(recall_val)
34
35     return {'k': k, 'train_recall': np.mean(recall_train_list), 'validation_recall': np.mean(recall_val_list)}
```

```
In [80]: 1 results = []
2 for k in np.arange(1, 100):
3     results.append(cross_validation(X_train, y_train, k, 5))
4
5 # Create DataFrame from results
6 crossval_df = pd.DataFrame(results)
7
8 print(crossval_df.head())
```

	k	train_recall	validation_recall
0	1	1.000000	0.539758
1	2	0.522691	0.261689
2	3	0.665922	0.399245
3	4	0.380519	0.184001
4	5	0.487719	0.274063

In [81]:

```
1 # empty dataframe
2 crossval_df = pd.DataFrame(columns = ['k', 'train_recall', 'validation_
3 # append results for each value of k
4 for k in np.arange(1,100):
5     crossval_df = crossval_df.append(cross_validation(X_train, y_train, k))
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.
```

```
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)
```

```
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
```

is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
```

```
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.
```

```
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)
```

```
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
```

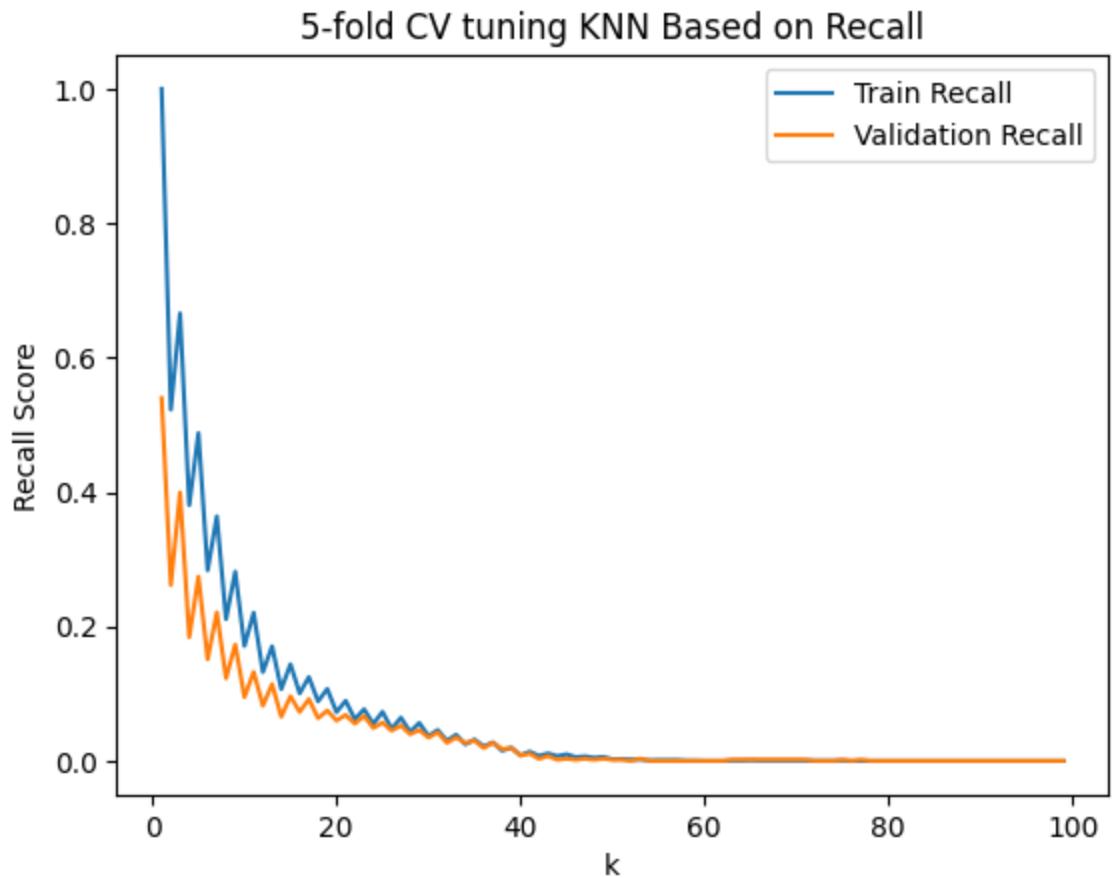
```
is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-81-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)
```

In [82]: 1 crossval\_df.info()

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 99 entries, 0 to 98  
Data columns (total 3 columns):  
 #   Column           Non-Null Count  Dtype     
 ---  --  
 0   k                99 non-null      float64  
 1   train_recall     99 non-null      float64  
 2   validation_recall 99 non-null      float64  
dtypes: float64(3)  
memory usage: 2.4 KB
```

In [83]:

```
1 fig, ax = plt.subplots()
2 sns.lineplot(x = 'k', y = 'train_recall',
3                 data = crossval_df,
4                 ax = ax, label = 'Train Recall')
5 sns.lineplot(x = 'k', y = 'validation_recall',
6                 data = crossval_df,
7                 ax = ax, label = 'Validation Recall')
8 ax.set_ylabel('Recall Score')
9 ax.set_title('5-fold CV tuning KNN Based on Recall')
10 plt.show()
```



This looks bizarre. Will try again without shuffle.

In [111]:

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

In [112]:

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

```
In [95]: ❶ def cross_validation(X_train, y_train, k, num_split = 5):
❷
❸     X_train = X_train.values
❹     y_train = y_train.values
❺
❻     recall_train_list = []
❼     recall_val_list = []
➋
⌋     for train_index, valid_index in KFold(n_splits = num_split).split():
⌌
⌍         # train and validation splitting
⌎         X_train_fold, X_val_fold = X_train[train_index], X_train[valid_
⌏         y_train_fold, y_val_fold = y_train[train_index], y_train[valid_
⌐
⌍         #create/fit the Standard scaler on the train fold
⌎         scaler = StandardScaler()
⌏         X_train_fold_scaled = scaler.fit_transform(X_train_fold)
⌍         # transform validation fold
⌎         X_val_fold_scaled = scaler.transform(X_val_fold)
⌐
⌍         smote=SMOTE(random_state=42)
⌎         X_train_fold_smote, y_train_fold_smote = smote.fit_resample(X_
⌐
⌍         # create/fit knearest neighbor
⌎         knn = KNeighborsClassifier(n_neighbors = k)
⌏         knn.fit(X_train_fold_smote, y_train_fold_smote)
⌐
⌍         # now how did we do?
⌎         y_train_pred = knn.predict(X_train_fold_scaled)
⌏         y_val_pred = knn.predict(X_val_fold_scaled)
⌐
⌍         recall_train = recall_score(y_train_fold, y_train_pred)
⌎         recall_val = recall_score(y_val_fold, y_val_pred)
⌐
⌍         recall_train_list.append(recall_train)
⌎         recall_val_list.append(recall_val)
⌐
⌋     return {'k': k, 'train_recall': np.mean(recall_train_list), 'valid_
```

In [96]:

```
1 # empty dataframe
2 crossval_df = pd.DataFrame(columns = ['k', 'train_recall', 'validation_
3 # append results for each value of k
4 for k in np.arange(1,100):
5     crossval_df = crossval_df.append(cross_validation(X_train, y_train, k))
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.
```

```
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)
```

```
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
```

is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
```

```
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.
```

```
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)  
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method  
is deprecated and will be removed from pandas in a future version. Use pa  
ndas.concat instead.  
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,  
5), ignore_index = True)
```

```
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
    crossval_df = crossval_df.append(cross_validation(X_train, y_train, k,
5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method
is deprecated and will be removed from pandas in a future version. Use pa
ndas.concat instead.
```

is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.

```
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
crossval_df = crossval_df.append(cross_validation(X_train, y_train, k, 5), ignore_index = True)
<ipython-input-96-9036620c5700>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
```

```
In [97]: 1 results = []
2 for k in np.arange(1, 100):
```

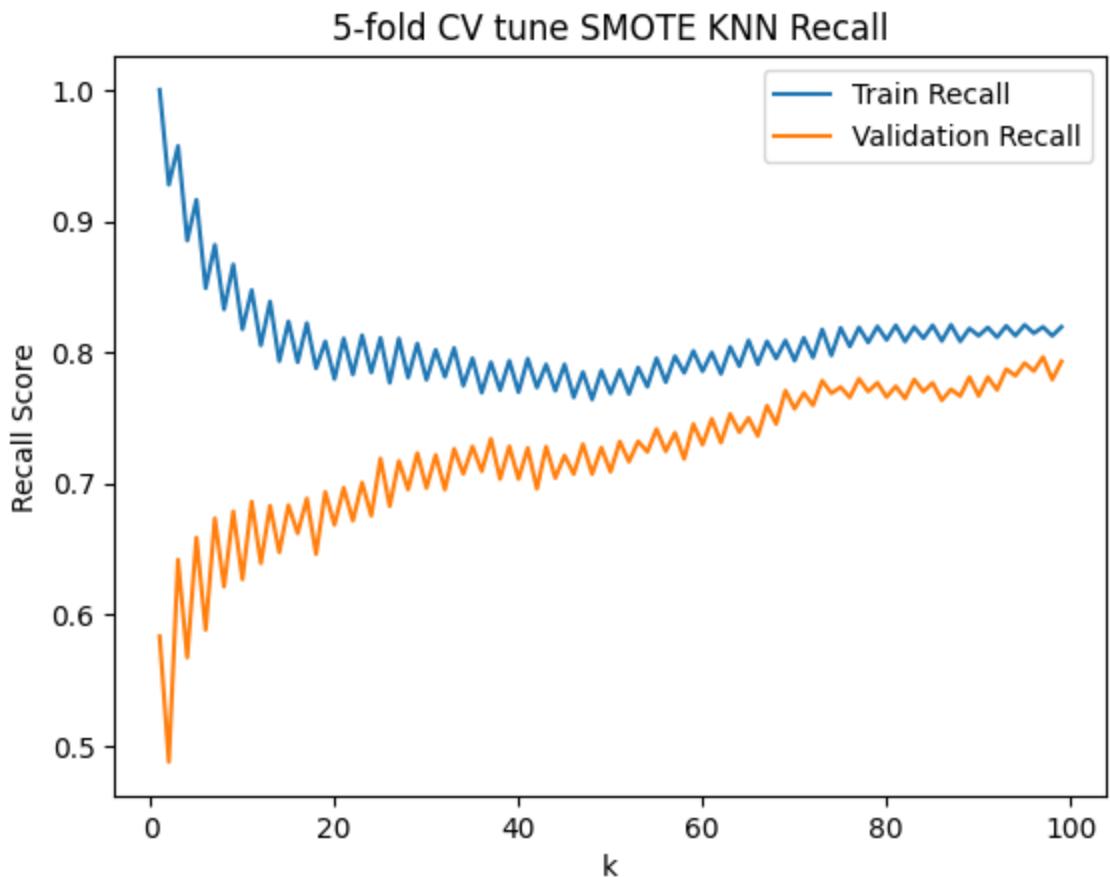
```
3     results.append(cross_validation)
4
5 # Create DataFrame from results
6 crossval_df = pd.DataFrame(results)
7
```

k	train_recall	validation_recall
0	1	1.000000
1	2	0.927756
2	3	0.957156
3	4	0.885197
4	5	0.916181

In [98]: 1 crossval\_df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 99 entries, 0 to 98
Data columns (total 3 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   k                99 non-null      int64  
 1   train_recall     99 non-null      float64 
 2   validation_recall 99 non-null      float64 
dtypes: float64(2), int64(1)
memory usage: 2.4 KB
```

In [100]: 1 fig, ax = plt.subplots()
2 sns.lineplot(x = 'k', y = 'train\_recall',
3 data = crossval\_df,
4 ax = ax, label = 'Train Recall')
5 sns.lineplot(x = 'k', y = 'validation\_recall',
6 data = crossval\_df,
7 ax = ax, label = 'Validation Recall')
8 ax.set\_ylabel('Recall Score')
9 ax.set\_title('5-fold CV tune SMOTE KNN Recall')
10 plt.show()



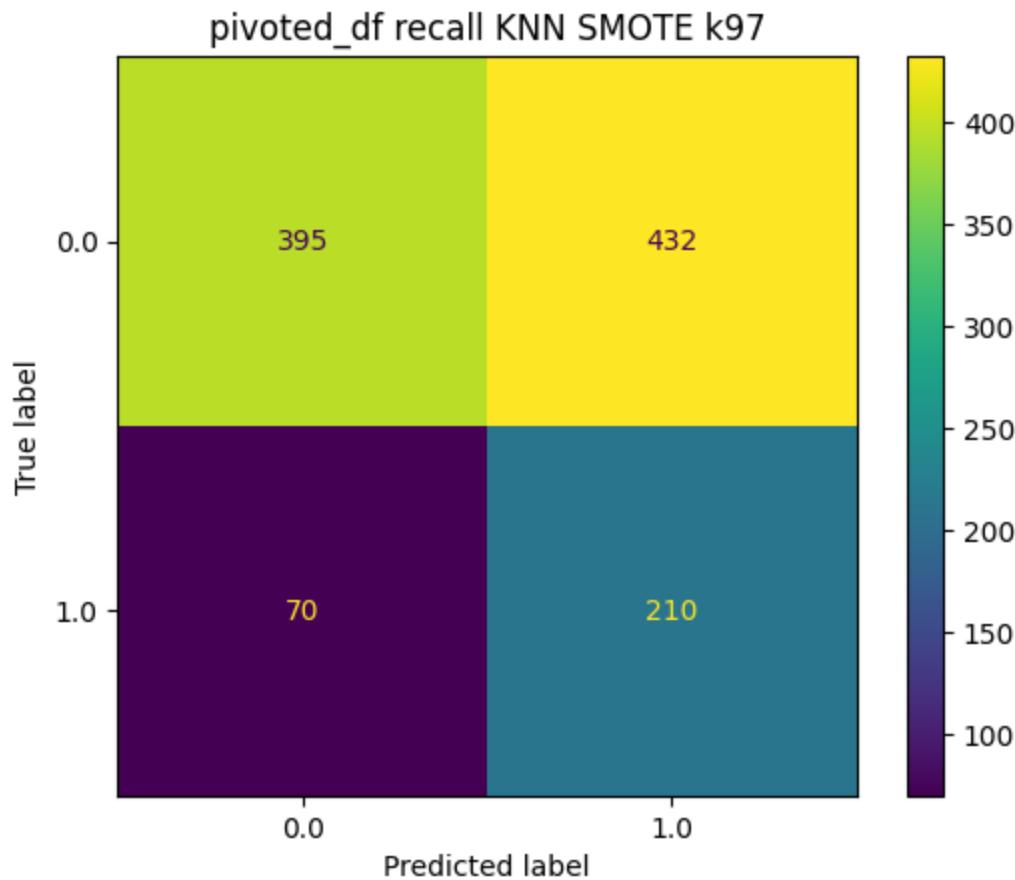
```
In [102]: 1 crossval_df.iloc[crossval_df['validation_recall'].idxmax()]
```

```
Out[102]: k          97.000000
train_recall      0.819350
validation_recall 0.796382
Name: 96, dtype: float64
```

```
In [107]: 1 fulltrain_scaler = StandardScaler()
2 X_train_sc = fulltrain_scaler.fit_transform(X_train)
3
4 smote=SMOTE(random_state=42)
5 X_train_smote, y_train_smote = smote.fit_resample(X_train_sc, y_train)
6
7 best_estimator = KNeighborsClassifier(n_neighbors = 97)
8 best_estimator.fit(X_train_smote, y_train_smote)
9
10 X_test_sc = fulltrain_scaler.transform(X_test)
11
12 # get predictions
13 y_pred = best_estimator.predict(X_test_sc)
14 test_recall = recall_score(y_test, y_pred)
15 print(f'Test Recall: {test_recall}')
```

Test Recall: 0.75

```
In [108]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('pivoted_df recall KNN SMOTE k97')
3 plt.show()
4 print(classification_report(y_test, pred))
```



	precision	recall	f1-score	support
0.0	0.84	0.67	0.75	827
1.0	0.39	0.61	0.47	280
accuracy			0.66	1107
macro avg	0.61	0.64	0.61	1107
weighted avg	0.72	0.66	0.68	1107

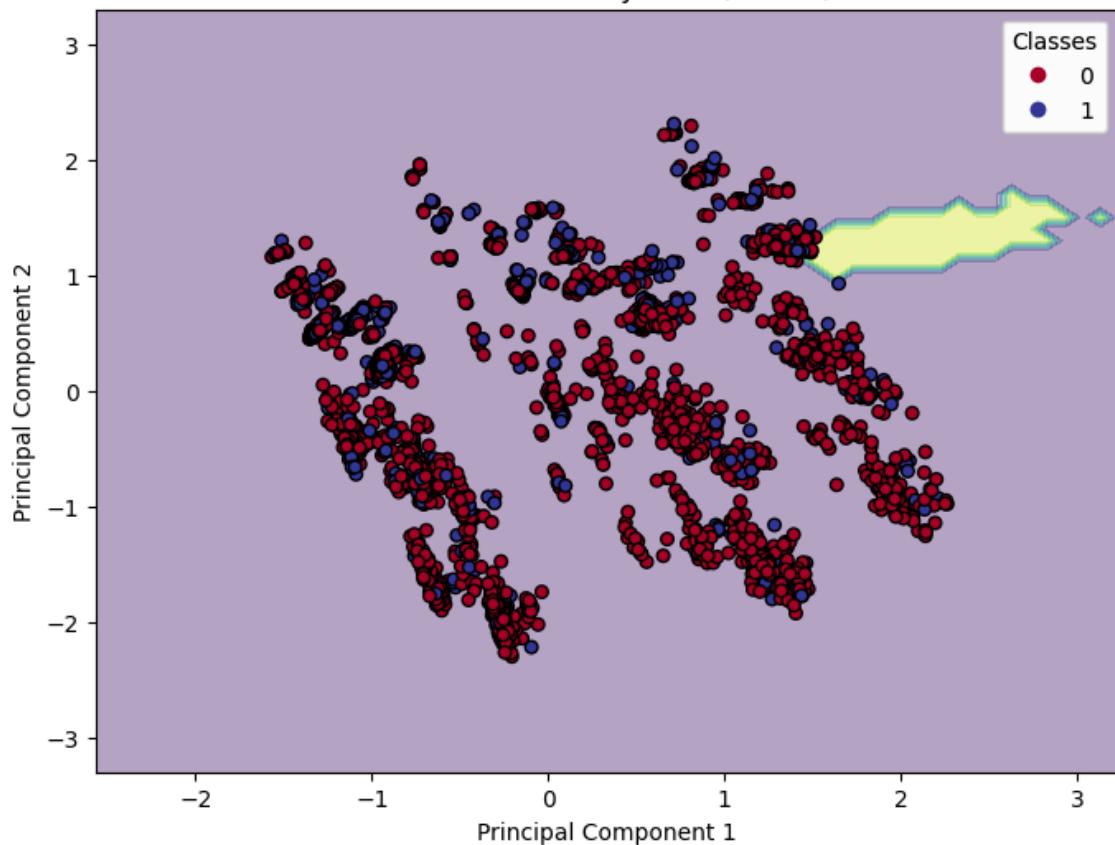
Want to pull more info from this. Cluster, analyze, label (power pitcher, etc.)? Hyperparameter tune?

To visualize the decision boundary, need to apply PCA (too many features)

In [109]:

```
1 # Reduce the data to two components for visualization
2 pca = PCA(n_components=2)
3 X_pca = pca.fit_transform(X)
4
5 # Scale the PCA-transformed data
6 fulltrain_scaler = StandardScaler()
7 X_full_sc = fulltrain_scaler.fit_transform(X_pca)
8
9 # Define the meshgrid for the contour plot based on the scaled data
10 x_min, x_max = X_full_sc[:, 0].min() - 1, X_full_sc[:, 0].max() + 1
11 y_min, y_max = X_full_sc[:, 1].min() - 1, X_full_sc[:, 1].max() + 1
12 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),
13                      np.arange(y_min, y_max, 0.1))
14
15 # Train the best estimator on the PCA-transformed and scaled data
16 best_estimator = KNeighborsClassifier(n_neighbors=97)
17 best_estimator.fit(X_full_sc, y) # Make sure y is your target array
18
19 # Perform the prediction on the meshgrid
20 Z = best_estimator.predict(np.c_[xx.ravel(), yy.ravel()])
21 Z = Z.reshape(xx.shape)
22
23 # Plot the decision boundary
24 f, ax = plt.subplots(figsize=(8, 6))
25 ax.contourf(xx, yy, Z, alpha=0.4)
26
27 # Scatter plot of the actual data points
28 scatter = ax.scatter(X_full_sc[:, 0], X_full_sc[:, 1], c=y, s=30, edge
29
30 # Create a Legend
31 legend1 = ax.legend(*scatter.legend_elements(), loc="upper right", tit
32 ax.add_artist(legend1)
33
34 # Label the axes and title
35 ax.set_xlabel('Principal Component 1')
36 ax.set_ylabel('Principal Component 2')
37 ax.set_title('Decision Boundary: KNN (k = 97)')
38
39 plt.show()
40
```

## Decision Boundary: KNN (k = 97)



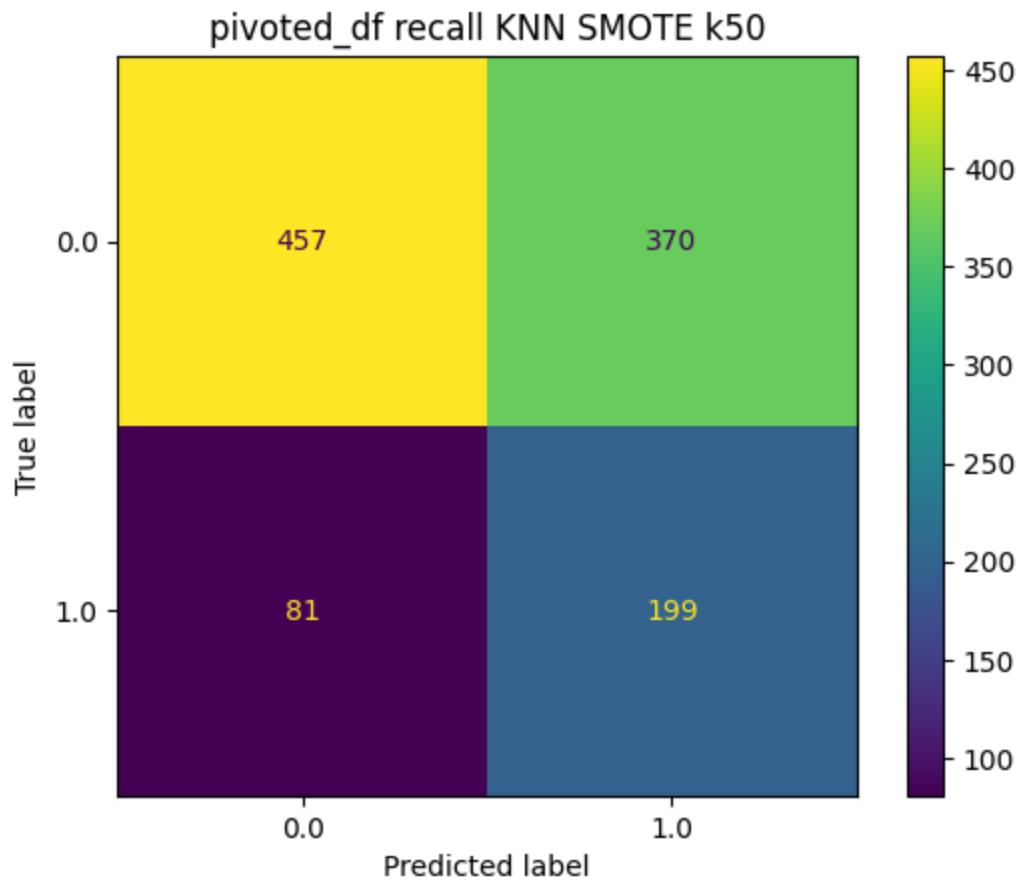
Don't think this makes sense. Will retry things with a smaller k value

In [113]:

```
1 fulltrain_scaler = StandardScaler()
2 X_train_sc = fulltrain_scaler.fit_transform(X_train)
3
4 smote=SMOTE(random_state=42)
5 X_train_smote, y_train_smote = smote.fit_resample(X_train_sc, y_train)
6
7 best_estimator = KNeighborsClassifier(n_neighbors = 50)
8 best_estimator.fit(X_train_smote, y_train_smote)
9
10 X_test_sc = fulltrain_scaler.transform(X_test)
11
12 # get predictions
13 y_pred = best_estimator.predict(X_test_sc)
14 test_recall = recall_score(y_test, y_pred)
15 print(f'Test Recall: {test_recall}')
```

Test Recall: 0.7107142857142857

```
In [114]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('pivoted_df recall KNN SMOTE k50')
3 plt.show()
4 print(classification_report(y_test, pred))
```

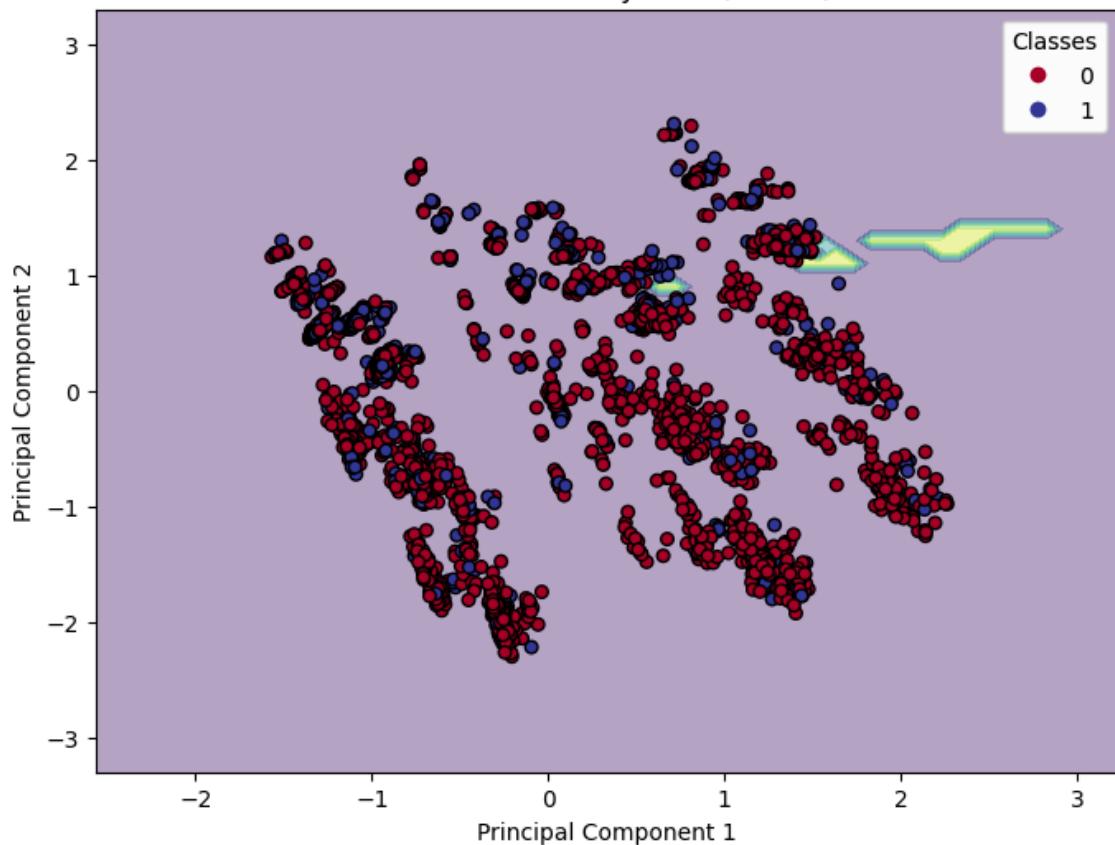


	precision	recall	f1-score	support
0.0	0.84	0.67	0.75	827
1.0	0.39	0.61	0.47	280
accuracy			0.66	1107
macro avg	0.61	0.64	0.61	1107
weighted avg	0.72	0.66	0.68	1107

In [116]:

```
1 # Reduce the data to two components for visualization
2 pca = PCA(n_components=2)
3 X_pca = pca.fit_transform(X)
4
5 # Scale the PCA-transformed data
6 fulltrain_scaler = StandardScaler()
7 X_full_sc = fulltrain_scaler.fit_transform(X_pca)
8
9 # Define the meshgrid for the contour plot based on the scaled data
10 x_min, x_max = X_full_sc[:, 0].min() - 1, X_full_sc[:, 0].max() + 1
11 y_min, y_max = X_full_sc[:, 1].min() - 1, X_full_sc[:, 1].max() + 1
12 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),
13                      np.arange(y_min, y_max, 0.1))
14
15 # Train the best estimator on the PCA-transformed and scaled data
16 best_estimator = KNeighborsClassifier(n_neighbors=50)
17 best_estimator.fit(X_full_sc, y) # Make sure y is your target array
18
19 # Perform the prediction on the meshgrid
20 Z = best_estimator.predict(np.c_[xx.ravel(), yy.ravel()])
21 Z = Z.reshape(xx.shape)
22
23 # Plot the decision boundary
24 f, ax = plt.subplots(figsize=(8, 6))
25 ax.contourf(xx, yy, Z, alpha=0.4)
26
27 # Scatter plot of the actual data points
28 scatter = ax.scatter(X_full_sc[:, 0], X_full_sc[:, 1], c=y, s=30, edge
29
30 # Create a Legend
31 legend1 = ax.legend(*scatter.legend_elements(), loc="upper right", tit
32 ax.add_artist(legend1)
33
34 # Label the axes and title
35 ax.set_xlabel('Principal Component 1')
36 ax.set_ylabel('Principal Component 2')
37 ax.set_title('Decision Boundary: KNN (k = 50)')
38
39 plt.show()
```

## Decision Boundary: KNN (k = 50)

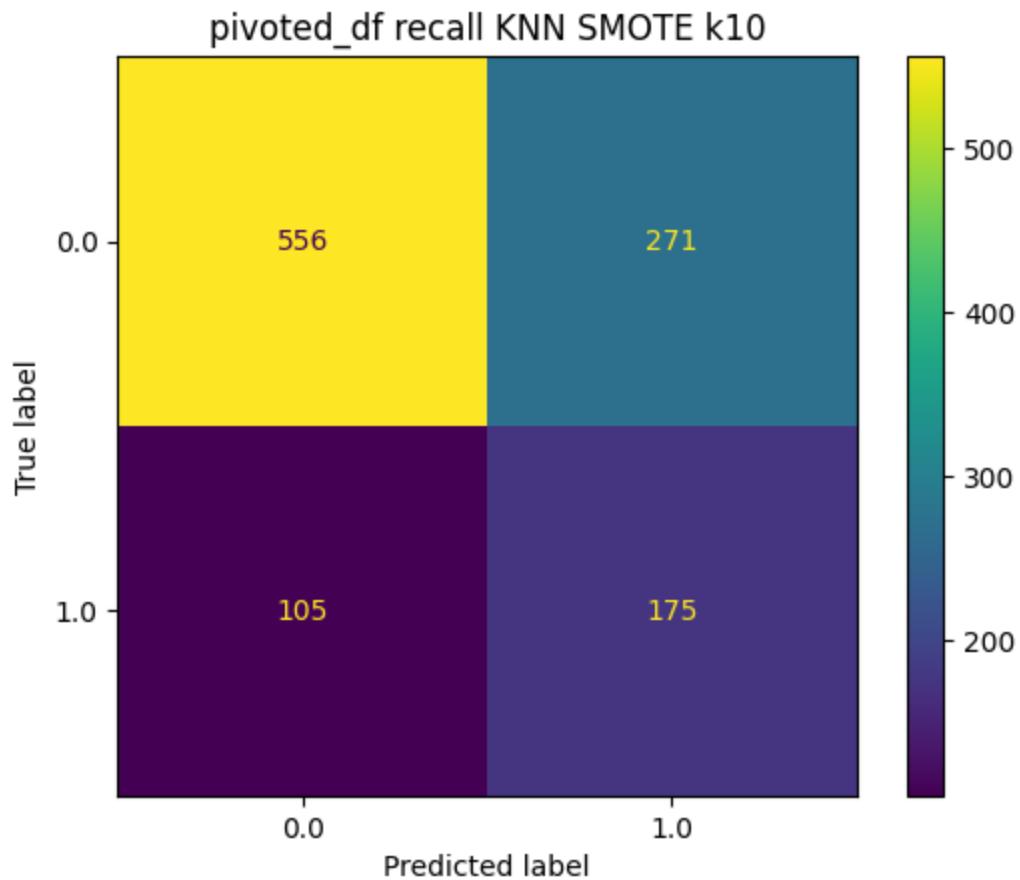


In [117]:

```
1 fulltrain_scaler = StandardScaler()
2 X_train_sc = fulltrain_scaler.fit_transform(X_train)
3
4 smote=SMOTE(random_state=42)
5 X_train_smote, y_train_smote = smote.fit_resample(X_train_sc, y_train)
6
7 best_estimator = KNeighborsClassifier(n_neighbors = 10)
8 best_estimator.fit(X_train_smote, y_train_smote)
9
10 X_test_sc = fulltrain_scaler.transform(X_test)
11
12 # get predictions
13 y_pred = best_estimator.predict(X_test_sc)
14 test_recall = recall_score(y_test, y_pred)
15 print(f'Test Recall: {test_recall}')
```

Test Recall: 0.625

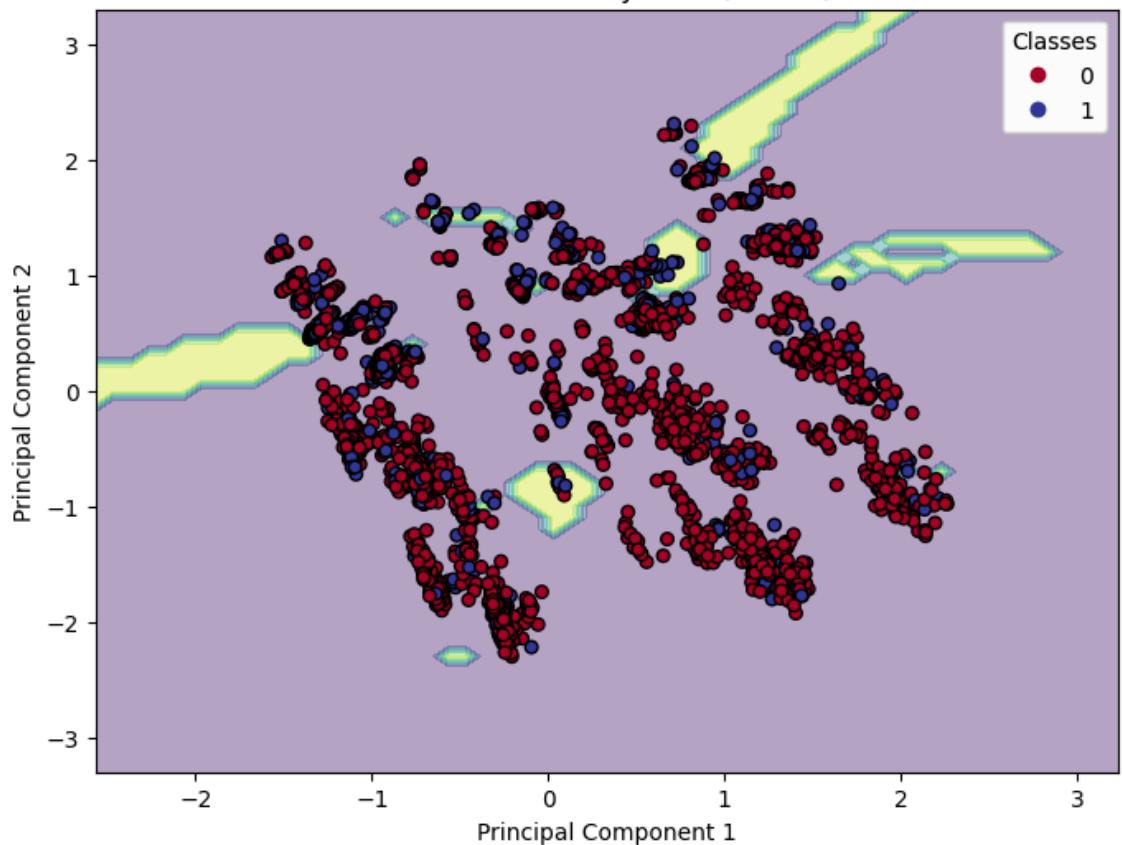
```
In [119]: 1 ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
2 plt.title('pivoted_df recall KNN SMOTE k10')
3 plt.show()
4 print(classification_report(y_test, pred))
```



In [120]:

```
1 # Reduce the data to two components for visualization
2 pca = PCA(n_components=2)
3 X_pca = pca.fit_transform(X)
4
5 # Scale the PCA-transformed data
6 fulltrain_scaler = StandardScaler()
7 X_full_sc = fulltrain_scaler.fit_transform(X_pca)
8
9 # Define the meshgrid for the contour plot based on the scaled data
10 x_min, x_max = X_full_sc[:, 0].min() - 1, X_full_sc[:, 0].max() + 1
11 y_min, y_max = X_full_sc[:, 1].min() - 1, X_full_sc[:, 1].max() + 1
12 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),
13                      np.arange(y_min, y_max, 0.1))
14
15 # Train the best estimator on the PCA-transformed and scaled data
16 best_estimator = KNeighborsClassifier(n_neighbors=10)
17 best_estimator.fit(X_full_sc, y) # Make sure y is your target array
18
19 # Perform the prediction on the meshgrid
20 Z = best_estimator.predict(np.c_[xx.ravel(), yy.ravel()])
21 Z = Z.reshape(xx.shape)
22
23 # Plot the decision boundary
24 f, ax = plt.subplots(figsize=(8, 6))
25 ax.contourf(xx, yy, Z, alpha=0.4)
26
27 # Scatter plot of the actual data points
28 scatter = ax.scatter(X_full_sc[:, 0], X_full_sc[:, 1], c=y, s=30, edge
29
30 # Create a Legend
31 legend1 = ax.legend(*scatter.legend_elements(), loc="upper right", tit
32 ax.add_artist(legend1)
33
34 # Label the axes and title
35 ax.set_xlabel('Principal Component 1')
36 ax.set_ylabel('Principal Component 2')
37 ax.set_title('Decision Boundary: KNN (k = 10)')
38
39 plt.show()
```

## Decision Boundary: KNN (k = 10)



This isn't working. KMeans and KNN may not be applicable to my dataset. Not enough data to differentiate groups of pitchers. Also dataset is very imbalanced, SMOTE doesn't work well for this.

In [ ]: 1