

## Part III\_SVM + Logit, no\_norm

March 11, 2021

### 0.1 Part III. SVM Logit

#### 0.1.1 Importing Packages

```
[1]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings("ignore")
from imblearn.over_sampling import SMOTE
from imblearn.under_sampling import NearMiss
from imblearn.pipeline import make_pipeline
from pylab import rcParams
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score, recall_score, confusion_matrix
from sklearn.metrics import f1_score, roc_auc_score, roc_curve
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from conf_matrix import func_confusion_matrix
```

#### 0.1.2 Functions for Model Evaluation

```
[2]: def generate_model_report(y_actual, y_predicted):
    print("Accuracy: " , accuracy_score(y_actual, y_predicted))
    print("Precision:" ,precision_score(y_actual, y_predicted))
    print("Recall:    " ,recall_score(y_actual, y_predicted))
    print("F1 Score: " ,f1_score(y_actual, y_predicted))
    print("AUC Score:" , roc_auc_score(y_actual, y_predicted))
    pass
```

```
[3]: def generate_auc_roc_curve(clf, X_test):
    y_pred_proba = clf.predict_proba(X_test)[: , 1]
    fpr, tpr, thresholds = roc_curve(Y_test, y_pred_proba)
    auc = roc_auc_score(Y_test, y_pred_proba)
    plt.plot(fpr,tpr,label="ROC Curve with AUC =" +str(auc))
    plt.legend(loc=4)
    plt.show()
```

```
pass
```

### 0.1.3 Importing and Cleaning Data

```
[4]: sdata = pd.read_csv('healthcare-dataset-stroke-data.csv')
```

```
[5]: sdata.shape
```

```
[5]: (5110, 12)
```

```
[6]: sdata = sdata[sdata.gender != 'Other']
sdata['gender'] = sdata['gender'].replace({'Male':0, 'Female':1}).astype(np.
    ↳uint8)
sdata['ever_married'] = sdata['ever_married'].replace({'No':0, 'Yes':1}).
    ↳astype(np.uint8)
sdata['Residence_type'] = sdata['Residence_type'].replace({'Rural':0, 'Urban':1}).
    ↳astype(np.uint8)
```

```
[7]: one_hot = pd.get_dummies(sdata['work_type'])
sdata = sdata.drop('work_type', axis = 1)
sdata = sdata.join(one_hot)
one_hot = pd.get_dummies(sdata['smoking_status'])
sdata = sdata.drop('smoking_status', axis = 1)
sdata = sdata.join(one_hot)
```

```
[8]: from sklearn.impute import KNNImputer
# there are 201 missing values in the BMI column, 40 of which experienced a
    ↳stroke event,
# so we impute using kNN:
imputer = KNNImputer(n_neighbors=70)
sdatakNN = imputer.fit_transform(sdata)
sdata = pd.DataFrame(sdatakNN, columns=sdata.columns)

#sdata = sdata.dropna()
```

```
[9]: from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = "all"
```

### 0.1.4 Splitting Data

```
[10]: X, Y = (sdata.drop(columns=['id', 'stroke']), sdata['stroke'])
```

```
[11]: from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
```

```

#X = MinMaxScaler().fit_transform(X)

#Splitting data into train, validation, test

print("Features data shape: {}".format(X.shape))
print("Target data shape: {}".format(Y.shape))

#Test train split for train vs test
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.
↪5,random_state=425)

#Test train split for train vs validation
X_train, X_val, Y_train, Y_val = train_test_split(X_train, Y_train, test_size=0.
↪2,random_state=425)

```

Features data shape: (5109, 17)  
Target data shape: (5109,)

# SVM

## 1.1 Baseline- SVM

```

[12]: #Importing necessary packages
import sklearn.svm as svm
from sklearn.metrics import confusion_matrix, accuracy_score

```

```

[13]: ## Baseline SVM
svm_base = svm.SVC(random_state=425).fit(X_train, Y_train)
#predicting the model using the test data
svm_base_predict = svm_base.predict(X_test)
#generating the model report
generate_model_report(Y_test, svm_base_predict)
#creating the matrix for actual vs predicted
pd.crosstab(svm_base_predict, Y_test, rownames=['Predicted'],
↪colnames=['Actual'])

```

Accuracy: 0.950293542074364  
Precision: 0.0  
Recall: 0.0  
F1 Score: 0.0  
AUC Score: 0.5

```

[13]: Actual      0.0   1.0
Predicted
0.0           2428   127

```

```
[14]: #predicting the svm base model and printing it
svm_base_predict = pd.Series(svm_base_predict)
svm_base_predict
```

```
[14]: 0      0.0
      1      0.0
      2      0.0
      3      0.0
      4      0.0
      ...
      2550    0.0
      2551    0.0
      2552    0.0
      2553    0.0
      2554    0.0
      Length: 2555, dtype: float64
```

## 1.2 Over Sampling- SVM

```
[15]: from imblearn.over_sampling import SMOTE

## Initialize dataframe of all parameters to be validated
results = pd.DataFrame(columns = ['Weight',
                                  'Kernel',
                                  'C Value',
                                  'Accuracy',
                                  'F1 Score',
                                  'AUC Score'])

## Loop over all parameters
for weight in [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0]: # SMOTE Ratio
    for k in ['linear', 'poly', 'rbf']:
        for c in np.linspace(0.01, 10, 5): # with Normalization increase last_
            ↪ number

            ## Oversampling technique
            sm = SMOTE(random_state=425, sampling_strategy=weight)
            X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)

            ## Model definition
            model = svm.SVC(kernel=k,
                             C=c,
                             random_state=425)

            ## Model fitting
            SVM = model.fit(X_train_res, Y_train_res)
```

```

Y_val_pred = model.predict(X_val) # predicted values

## Model evaluation
accuracy = accuracy_score(Y_val, Y_val_pred)
f1 = f1_score(Y_val, Y_val_pred)
auc = roc_auc_score(Y_val, Y_val_pred)

## Append results to dataframe
results = results.append({'Weight':weight,
                        'Kernel':k,
                        'C Value': c,
                        'Accuracy':accuracy,
                        'F1 Score':f1,
                        'AUC Score': auc}, ignore_index = True)

```

```

[16]: ## Retrieve result with highest F1 Score
results.loc[results['F1 Score'].idxmax()]

```

```

[16]: Weight      0.300000
      Kernel      poly
      C Value    10.000000
      Accuracy    0.900196
      F1 Score    0.301370
      AUC Score   0.681934
      Name: 39, dtype: object

```

```

[17]: ## Another loop to get more granularity of SMOTE
      # Before I tested 10 different smote weights, now I test 3000
results_sm = pd.DataFrame(columns = ['Weight', 'Accuracy', 'F1 Score',
      ↪ 'MA_Acc', 'MA_F1'])
weights = np.linspace(0.1, 1, 3000)

for weight in weights:
    sm = SMOTE(random_state=425, sampling_strategy=weight)
    X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
    model = svm.SVC(kernel='poly',
                    C= 10,
                    random_state=425).fit(X_train_res, Y_train_res)
    Y_val_pred = model.predict(X_val)
    accuracy = accuracy_score(Y_val, Y_val_pred)
    f1 = f1_score(Y_val, Y_val_pred)
    results_sm = results_sm.append({'Weight':weight, 'Accuracy':accuracy, 'F1_
    ↪ Score':f1}, ignore_index = True)

```

```

[18]: ## Plotting moving average of SMOTE weights and scor
for i in range(15,len(results)):
    results_sm['MA_F1'][i] = (results_sm['F1 Score'][i-15:i]).mean()

```

```

    results_sm['MA_Acc'][i] = (results_sm['Accuracy'][i-15:i]).mean()

## Plot
fig, ax = plt.subplots(figsize=(12,6))

ax.plot(results_sm['Weight'], results_sm['MA_F1'], color='tab:red', label="F1_
↪Score")
ax.plot(results_sm['Weight'], results_sm['F1 Score'], color='tab:red', ↪
↪label="F1 Score", alpha=0.15)
ax.set_ylabel('F1 Score', color='tab:red')
ax.set_xlabel("SMOTE Weight")
ax2=ax.twinx()
ax2.plot(results_sm['Weight'], results_sm['MA_Acc'], color='tab:blue', ↪
↪label="Accuracy")
ax2.set_ylabel('Accuracy', color='tab:blue')
plt.show()

```

[18]: [matplotlib.lines.Line2D at 0x7fcb98613df0]

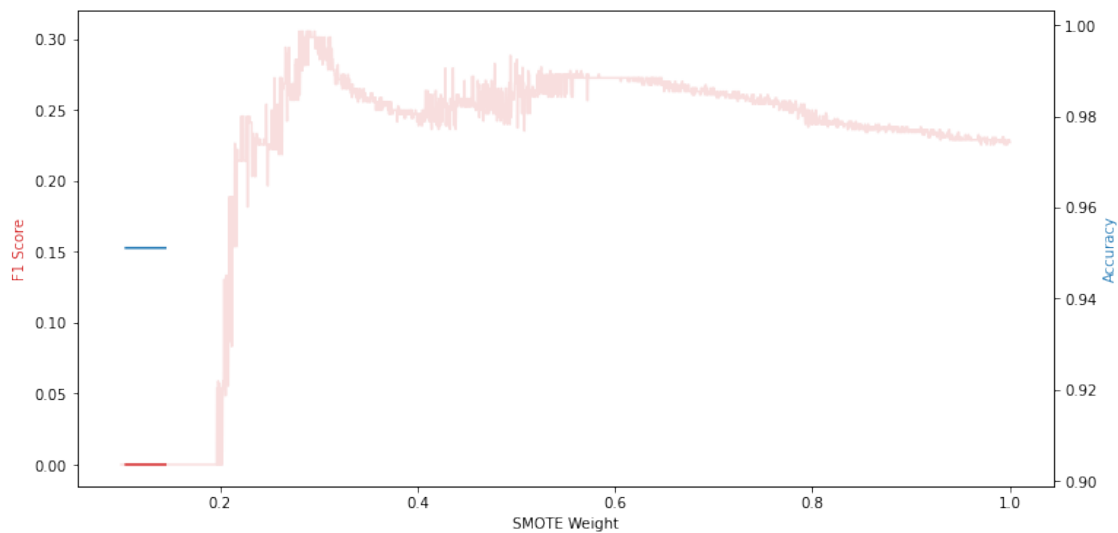
[18]: [matplotlib.lines.Line2D at 0x7fcb986291c0]

[18]: Text(0, 0.5, 'F1 Score')

[18]: Text(0.5, 0, 'SMOTE Weight')

[18]: [matplotlib.lines.Line2D at 0x7fcb98874610]

[18]: Text(0, 0.5, 'Accuracy')



```
[19]: ## Retrieve result with highest F1 Score
results_sm.loc[results_sm['F1 Score'].idxmax()]
```

```
[19]: Weight      0.280360
Accuracy    0.902153
F1 Score    0.305556
MA_Acc      NaN
MA_F1       NaN
Name: 601, dtype: float64
```

```
[20]: ## Run best model on testing data, only trying to optimize weight, keep the
      ↪ kernel and C parameter from before
sm = SMOTE(random_state=425, sampling_strategy=0.280360)
X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
model = svm.SVC(kernel='poly',
                  C= 10,
                  random_state=425).fit(X_train_res, Y_train_res)
#testing the best model on the test data
Y_test_pred = model.predict(X_test)
#generating the model report
generate_model_report(Y_test, Y_test_pred)
#creating the confusion matrix
conf_mat = confusion_matrix(Y_test, Y_test_pred)
pd.crosstab(Y_test, Y_test_pred, rownames=['Predicted'], colnames=['Actual'])
```

```
Accuracy:  0.8759295499021527
Precision: 0.1895424836601307
Recall:    0.4566929133858268
F1 Score:  0.2678983833718244
AUC Score: 0.6772756164952198
```

```
[20]: Actual      0.0   1.0
Predicted
0.0          2180   248
1.0           69    58
```

### 1.3 Under Sampling- SVM

```
[21]: from imblearn.under_sampling import NearMiss

## Initialize dataframe of all parameters to be validated
results = pd.DataFrame(columns = ['Neighbor',
                                  'Kernel',
                                  'C Value',
                                  'Accuracy',
                                  'F1 Score',
                                  'AUC Score'])
```

```

## Loop over all parameters
for neighbor in [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]:
    for k in ['linear', 'poly', 'rbf']:
        for c in np.linspace(0.01, 10, 5): # with Normalization increase last_
            ↪ number

            ## Underampling technique
            nm = NearMiss(version=1, n_neighbors=neighbor)
            X_train_res, Y_train_res = nm.fit_resample(X_train, Y_train)

            ## Model definition
            model = svm.SVC(kernel=k,
                             C=c,
                             random_state=425)

            ## Model fitting
            SVM = model.fit(X_train_res, Y_train_res)
            Y_val_pred = model.predict(X_val) # predicted values

            ## Model evaluation
            accuracy = accuracy_score(Y_val, Y_val_pred)
            f1 = f1_score(Y_val, Y_val_pred)
            auc = roc_auc_score(Y_val, Y_val_pred)

            ## Append results to dataframe
            results = results.append({'Neighbor': neighbor,
                                     'Kernel': k,
                                     'C Value': c,
                                     'Accuracy': accuracy,
                                     'F1 Score': f1,
                                     'AUC Score': auc}, ignore_index = True)

```

```

[22]: ## Retrieve result with highest F1 Score
      results.loc[results['F1 Score'].idxmax()]

```

```

[22]: Neighbor      1
      Kernel      linear
      C Value     5.005000
      Accuracy    0.831703
      F1 Score     0.232143
      AUC Score    0.683868
      Name: 2, dtype: object

```

```

[23]: ## Another loop to get more granularity of SMOTE
      # Before I tested 10 different neighbor, now I test 3000

```



```

results_sm = pd.DataFrame(columns = ['Neighbor', 'Accuracy', 'F1 Score',
    ↪ 'MA_Acc', 'MA_F1'])
neighbors = np.linspace(0.10, 1, 100)#had to make this smaller in order to get
    ↪it to run, ran it over night before and still didnt work

for neighbor in neighbors:
    sm = NearMiss(version=1,sampling_strategy=neighbor)
    X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
    model = svm.SVC(kernel='linear',
        C= 5.005000,
        random_state=425).fit(X_train_res, Y_train_res)
    Y_val_pred = model.predict(X_val)
    accuracy = accuracy_score(Y_val, Y_val_pred)
    f1 = f1_score(Y_val, Y_val_pred)
    results_sm = results_sm.append({'Neighbor':neighbor, 'Accuracy':
    ↪accuracy, 'F1 Score':f1}, ignore_index = True)

```

```

[24]: ## Plotting moving average of Best near miss model neighbor and scores
for i in range(15,len(results)):
    results_sm['MA_F1'][i] = (results_sm['F1 Score'][i-15:i]).mean()
    results_sm['MA_Acc'][i] = (results_sm['Accuracy'][i-15:i]).mean()

## Plot
fig, ax = plt.subplots(figsize=(12,6))

ax.plot(results_sm['Neighbor'], results_sm['MA_F1'], color='tab:red', label="F1
    ↪Score")
ax.plot(results_sm['Neighbor'], results_sm['F1 Score'], color='tab:red',
    ↪label="F1 Score", alpha=0.15)
ax.set_ylabel('F1 Score', color='tab:red')
ax.set_xlabel("NearMiss Neighbor")
ax2=ax.twinx()
ax2.plot(results_sm['Neighbor'], results_sm['MA_Acc'], color='tab:blue',
    ↪label="Accuracy")
ax2.set_ylabel('Accuracy', color='tab:blue')
plt.show()

```

[24]: [`<matplotlib.lines.Line2D at 0x7fcb988b0850>`]

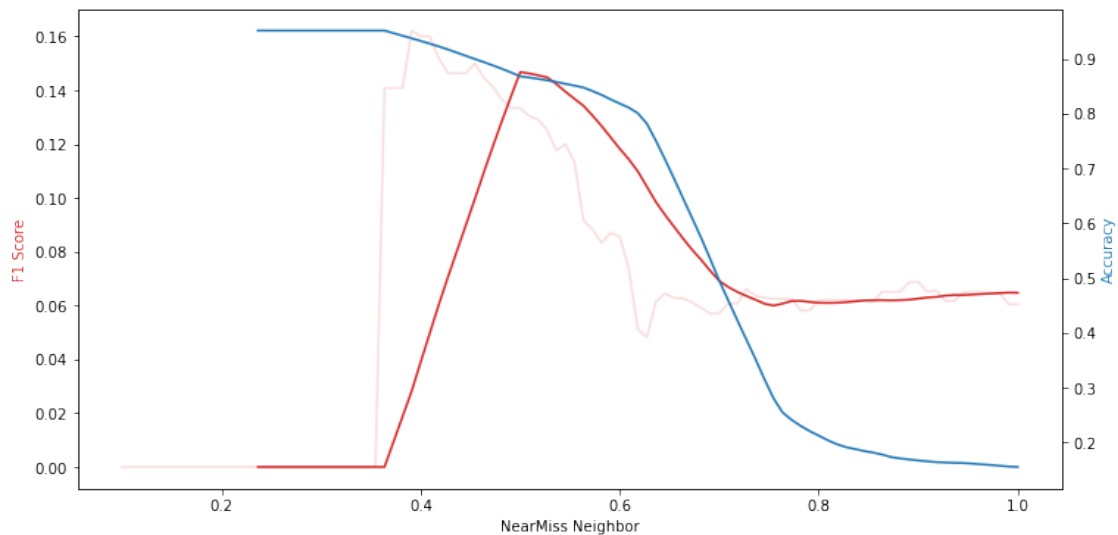
[24]: [`<matplotlib.lines.Line2D at 0x7fcb988b0ca0>`]

[24]: `Text(0, 0.5, 'F1 Score')`

[24]: `Text(0.5, 0, 'NearMiss Neighbor')`

[24]: [`<matplotlib.lines.Line2D at 0x7fcb98bd8ca0>`]

```
[24]: Text(0, 0.5, 'Accuracy')
```



```
[25]: ## Retrieve result with highest F1 Score
results_sm.loc[results_sm['F1 Score'].idxmax()]
```

```
[25]: Neighbor      0.390909
Accuracy      0.878669
F1 Score      0.162162
MA_Acc        0.936986
MA_F1         0.028169
Name: 32, dtype: float64
```

```
[26]: ## Run best model on testing data, only trying to optimize neighbor, keep the
      ↪kernel and C parameter from before
sm = NearMiss(version=1, sampling_strategy=0.390909)
X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
model = svm.SVC(kernel='linear',
                  C= 5.005000,
                  random_state=425).fit(X_train_res, Y_train_res)
#using the best model on the test data
Y_test_pred = model.predict(X_test)
#generate the model report
generate_model_report(Y_test, Y_test_pred)
#creating the confusion matrix
conf_mat = confusion_matrix(Y_test, Y_test_pred)
pd.crosstab(Y_test, Y_test_pred, rownames=['Predicted'], colnames=['Actual'])
```

```
Accuracy:  0.8634050880626223
Precision: 0.1396103896103896
```

Recall: 0.33858267716535434  
F1 Score: 0.19770114942528738  
AUC Score: 0.6147196746617546

```
[26]: Actual      0.0   1.0  
      Predicted  
      0.0      2163  265  
      1.0       84   43
```

## 1.4 Logit

```
[27]: #Import necessary packages  
      from sklearn.linear_model import LogisticRegression  
      from sklearn import metrics  
      import seaborn as sns
```

## 1.5 Logit- Base

```
[28]: #Fitting the model
lr= LogisticRegression()
lr.fit(X_train,Y_train)

#Predicting
Y_Test_Pred1=lr.predict(X_test)
print('Accuracy: ',metrics.accuracy_score(Y_test, Y_Test_Pred1))

#CF
import sklearn.metrics as sklmetrics
conf_mat = sklmetrics.confusion_matrix(Y_test, Y_Test_Pred1)
print(conf_mat)

# Confusion matrix
sns.heatmap(conf_mat, fmt='d',square=True, annot=True, cbar = False,
            xticklabels = ['Failure','Success'],
            yticklabels = ['Failure','Success'])
plt.xlabel("Predicted Value")
plt.ylabel("True Value")

#Metrics
print(metrics.classification_report(Y_test,Y_Test_Pred1))

auc = roc_auc_score(Y_test,Y_Test_Pred1)
auc

#generate the model report
generate_model_report(Y_test, Y_Test_Pred1)
```

```
[28]: LogisticRegression()

Accuracy: 0.950293542074364
[[2428  0]
 [ 127  0]]
```

```
[28]: <AxesSubplot:>
```

```
[28]: Text(0.5, 15.0, 'Predicted Value')
```

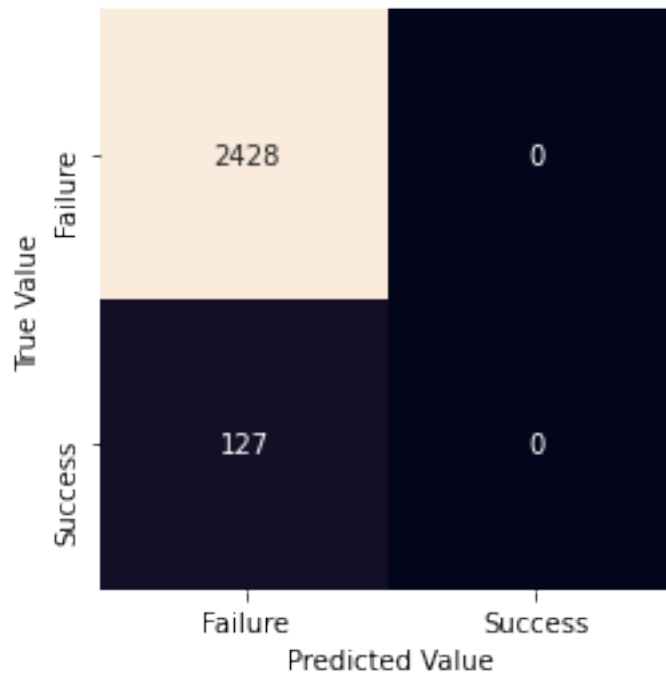
```
[28]: Text(91.68, 0.5, 'True Value')
```

	precision	recall	f1-score	support
0.0	0.95	1.00	0.97	2428
1.0	0.00	0.00	0.00	127
accuracy			0.95	2555

macro avg	0.48	0.50	0.49	2555
weighted avg	0.90	0.95	0.93	2555

[28]: 0.5

Accuracy: 0.950293542074364  
 Precision: 0.0  
 Recall: 0.0  
 F1 Score: 0.0  
 AUC Score: 0.5



[ ]:

## 1.6 Logit- Over Sampling

```
[29]: from imblearn.over_sampling import SMOTE

      ## Initialize dataframe of all parameters to be validated
      results = pd.DataFrame(columns = ['Weight',
                                       'C',
                                       'Solver',
                                       'Accuracy',
                                       'F1 Score',
                                       'AUC Score'])
```

```

## Loop over all parameters
for weight in np.linspace(0.1,1,100): # SMOTE Ratio
    for s in ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']:
        for c in range(1,10):

            ## Oversampling technique
            sm = SMOTE(random_state=425, sampling_strategy=weight)
            X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)

            ## Model definition
            model = LogisticRegression(
                                C=c,
                                solver=s,
                                random_state=425)

            ## Model fitting
            lr = model.fit(X_train_res, Y_train_res)
            Y_val_pred = model.predict(X_val) # predicted values

            ## Model evaluation
            accuracy = accuracy_score(Y_val, Y_val_pred)
            f1 = f1_score(Y_val, Y_val_pred)
            auc = roc_auc_score(Y_val, Y_val_pred)

            ## Append results to dataframe
            results = results.append({'Weight':weight,
                                    'C': c,
                                    'Solver':s,
                                    'Accuracy':accuracy,
                                    'F1 Score':f1,
                                    'AUC Score': auc}, ignore_index =_
↪True)

```

```

[30]: ## Retrieve result with highest F1 Score
      results.loc[results['F1 Score'].idxmax()]

```

```

[30]: Weight      0.218182
      C           7
      Solver      lbfgs
      Accuracy    0.925636
      F1 Score    0.366667
      AUC Score   0.695309
      Name: 600, dtype: object

```

```

[31]: ## Another loop to get more granularity of SMOTE
      # Before I tested 10 different smote weights, now I test 3000

```

```

results_sm = pd.DataFrame(columns = ['Weight', "C", "Solver", 'Accuracy', 'F1_
↳Score', 'MA_Acc', 'MA_F1'])
weights = np.linspace(0.07, 1, 3000)

for weight in weights:
    sm = SMOTE(random_state=425, sampling_strategy=weight)
    X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
    model = LogisticRegression(C=7,
                               solver="lbfgs",
                               random_state=425).fit(X_train_res, Y_train_res)
    Y_val_pred = model.predict(X_val)
    accuracy = accuracy_score(Y_val, Y_val_pred)
    f1 = f1_score(Y_val, Y_val_pred)
    results_sm = results_sm.append({'Weight':weight, "C":c, "Solver":s,
↳'Accuracy':accuracy, 'F1 Score':f1}, ignore_index = True)

```

```

[32]: ## Plotting moving average of SMOTE weights and scor
for i in range(15, len(results)):
    results_sm['MA_F1'][i] = (results_sm['F1 Score'][i-15:i]).mean()
    results_sm['MA_Acc'][i] = (results_sm['Accuracy'][i-15:i]).mean()

## Plot
fig, ax = plt.subplots(figsize=(12,6))

#ax.plot(results_sm['Weight'], results_sm['MA_F1'], color='tab:red', label="F1_
↳Score")
ax.plot(results_sm['Weight'], results_sm['F1 Score'], color='tab:red',
↳label="F1 Score", alpha=1)
ax.set_ylabel('F1 Score', color='tab:red')
ax.set_xlabel("SMOTE Weight")
ax2=ax.twinx()
ax2.plot(results_sm['Weight'], results_sm['Accuracy'], color='tab:blue',
↳label="Accuracy")
ax2.set_ylabel('Accuracy', color='tab:blue')
plt.show()

```

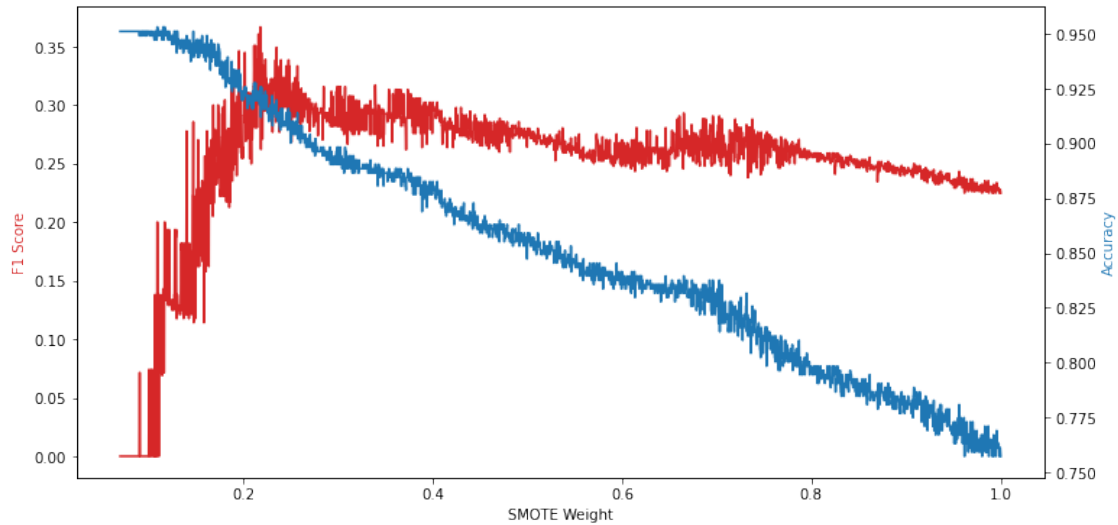
[32]: [

[32]: Text(0, 0.5, 'F1 Score')

[32]: Text(0.5, 0, 'SMOTE Weight')

[32]: [

[32]: Text(0, 0.5, 'Accuracy')



```
[33]: ## Retrieve result with highest F1 Score
results_sm.loc[results_sm['F1 Score'].idxmax()]
```

```
[33]: Weight      0.217919
C              9
Solver         saga
Accuracy       0.925636
F1 Score       0.366667
MA_Acc         0.920939
MA_F1          0.321909
Name: 477, dtype: object
```

```
[34]: ## Run best model on testing data
sm = SMOTE(random_state=425, sampling_strategy= 0.217919)
X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
model = LogisticRegression(C=7,
                           solver='lbfgs',
                           random_state=425).fit(X_train_res, Y_train_res)
#using the model on test data
Y_test_pred = model.predict(X_test)
#generate the model report
generate_model_report(Y_test, Y_test_pred)
#confusion matrix
conf_mat = confusion_matrix(Y_test, Y_test_pred)
pd.crosstab(Y_test, Y_test_pred, rownames=['Predicted'], colnames=['Actual'])
```

```
Accuracy:  0.9048923679060665
Precision: 0.22380952380952382
Recall:    0.3700787401574803
```



F1 Score: 0.2789317507418398  
AUC Score: 0.6514726484971916

```
[34]: Actual      0.0  1.0  
      Predicted  
      0.0      2265  163  
      1.0       80   47
```

```
[ ]:
```

## 1.7 Logit- Under Sampling

```
[35]: from imblearn.under_sampling import NearMiss  
  
      ## Initialize dataframe of all parameters to be validated  
      results = pd.DataFrame(columns = ['Neighbor',  
                                       'C',  
                                       'Solver',  
                                       'Accuracy',  
                                       'F1 Score',  
                                       'AUC Score'])  
  
      ## Loop over all parameters  
      for neighbor in [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]:  
          for s in ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']:  
              for c in range(1,10):  
  
                  ## Underampling technique  
                  nm = NearMiss(version=1,n_neighbors=neighbor)  
                  X_train_res, Y_train_res = nm.fit_resample(X_train, Y_train)  
  
                  ## Model definition  
                  model = LogisticRegression(  
                                          C=c,  
                                          solver=s,  
                                          random_state=425)  
  
                  ## Model fitting  
                  lr = model.fit(X_train_res, Y_train_res)  
                  Y_val_pred = model.predict(X_val) # predicted values  
  
                  ## Model evaluation  
                  accuracy = accuracy_score(Y_val, Y_val_pred)  
                  f1 = f1_score(Y_val, Y_val_pred)  
                  auc = roc_auc_score(Y_val, Y_val_pred)  
  
                  ## Append results to dataframe
```

```

        results = results.append({'Neighbor':neighbor,
                                  'C': c,
                                  'Solver':s,
                                  'Accuracy':accuracy,
                                  'F1 Score':f1,
                                  'AUC Score': auc}, ignore_index =_
→True)

```

```

[36]: ## Retrieve result with highest F1 Score
results.loc[results['F1 Score'].idxmax()]

```

```

[36]: Neighbor          1
      C                4
      Solver          lbfgs
      Accuracy      0.818004
      F1 Score      0.243902
      AUC Score     0.714609
      Name: 12, dtype: object

```

```

[37]: ## Another loop to get more granularity of SMOTE
# Before I tested 10 different Neighbors, now I test 3000
results_sm = pd.DataFrame(columns = ['Neighbor','C',"Solver", 'Accuracy', 'F1_
→Score', 'MA_Acc', 'MA_F1'])
neighbors = np.linspace(0.07, 1, 3000)

for neighbor in neighbors:
    sm = NearMiss(sampling_strategy=neighbor)
    X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
    model = LogisticRegression(C=4,
                                solver="lbfgs",
                                random_state=425).fit(X_train_res, Y_train_res)
    Y_val_pred = model.predict(X_val)
    accuracy = accuracy_score(Y_val, Y_val_pred)
    f1 = f1_score(Y_val, Y_val_pred)
    results_sm = results_sm.append({'Neighbor':neighbor,"C":c,"Solver":s,_
→'Accuracy':accuracy,'F1 Score':f1}, ignore_index = True)

```

```

[38]: ## Plotting moving average of Best near miss model Neighbor and scores
for i in range(15,len(results)):
    results_sm['MA_F1'][i] = (results_sm['F1 Score'][i-15:i]).mean()
    results_sm['MA_Acc'][i] = (results_sm['Accuracy'][i-15:i]).mean()

## Plot
fig, ax = plt.subplots(figsize=(12,6))

ax.plot(results_sm['Neighbor'], results_sm['MA_F1'], color='tab:red', label="F1_
→Score")

```

```

ax.plot(results_sm['Neighbor'], results_sm['F1 Score'], color='tab:red',
        label="F1 Score", alpha=0.15)
ax.set_ylabel('F1 Score', color='tab:red')
ax.set_xlabel("NearMiss Neighbor")
ax2=ax.twinx()
ax2.plot(results_sm['Neighbor'], results_sm['MA_Acc'], color='tab:blue',
        label="Accuracy")
ax2.set_ylabel('Accuracy', color='tab:blue')
plt.show()

```

[38]: [<matplotlib.lines.Line2D at 0x7fcb99c1fe50>]

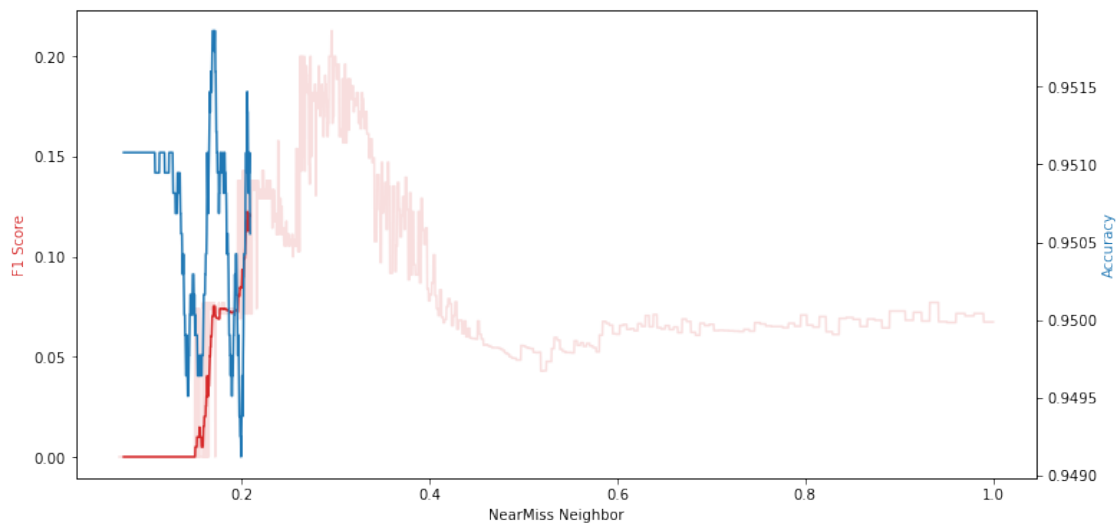
[38]: [<matplotlib.lines.Line2D at 0x7fcb99d3b2b0>]

[38]: Text(0, 0.5, 'F1 Score')

[38]: Text(0.5, 0, 'NearMiss Neighbor')

[38]: [<matplotlib.lines.Line2D at 0x7fcb98d5ca00>]

[38]: Text(0, 0.5, 'Accuracy')



```

[39]: ## Retrieve result with highest F1 Score
results_sm.loc[results_sm['F1 Score'].idxmax()]

```

```

[39]: Neighbor    0.295755
      C          9
      Solver      saga
      Accuracy    0.927593

```

F1 Score      0.212766  
MA\_Acc              NaN  
MA\_F1              NaN  
Name: 728, dtype: object

```
[40]: ## Run best model on testing data
sm = NearMiss(version=1,sampling_strategy= 0.295755)
X_train_res, Y_train_res = sm.fit_resample(X_train, Y_train)
model = LogisticRegression(C=4,
                           solver="lbfgs",
                           random_state=425).fit(X_train_res, Y_train_res)
#using the best model on the test data
Y_test_pred = model.predict(X_test)
#generating the model report
generate_model_report(Y_test, Y_test_pred)
#creating the confusion matrix
conf_mat = confusion_matrix(Y_test, Y_test_pred)
pd.crosstab(Y_test, Y_test_pred, rownames=['Predicted'], colnames=['Actual'])
```

Accuracy: 0.911545988258317  
Precision: 0.15384615384615385  
Recall: 0.1732283464566929  
F1 Score: 0.16296296296296295  
AUC Score: 0.5616965455512459

```
[40]: Actual      0.0   1.0
Predicted
0.0          2307   121
1.0          105   22
```

```
[ ]: 
```

```
[ ]: 
```

```
[ ]: 
```

```
[ ]: 
```

```
[ ]: 
```

```
[ ]: 
```