

# Logistic Regression Baseline Model Development

The purpose of this notebook is to explore preprocessing techniques beyond the simple data cleaning formalized in the notebook "EDA" (!!!hyperlink this!!!).

My main goal is to discover the best sampling method to address the class imbalance present in the target variable.

I will also prototype various hyper parameter settings.

## Load and clean data


```
In [1]: 1 from preprocessor import data_cleaner
        2
        3 import pandas as pd
        4 import numpy as np
        5
        6 from sklearn.linear_model import LogisticRegression, LogisticRegressionCV
        7 from sklearn.metrics import plot_confusion_matrix, classification_report
        8 from sklearn.model_selection import cross_val_score, cross_val_predict
        9
```

```
In [2]: 1 # import and define training data. See preprocessor.py file for details
        2 X_train, y_train = data_cleaner("../data/train.csv.zip")
```

In [3]:  1 X\_train.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 103904 entries, 0 to 103903
Data columns (total 23 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   Age                                         103904 non-null  int64
1   Flight Distance                           103904 non-null  int64
2   Inflight wifi service                     103904 non-null  int64
3   Departure/Arrival time convenient         103904 non-null  int64
4   Ease of Online booking                    103904 non-null  int64
5   Gate location                             103904 non-null  int64
6   Food and drink                            103904 non-null  int64
7   Online boarding                           103904 non-null  int64
8   Seat comfort                              103904 non-null  int64
9   Inflight entertainment                    103904 non-null  int64
10  On-board service                          103904 non-null  int64
11  Leg room service                          103904 non-null  int64
12  Baggage handling                          103904 non-null  int64
13  Checkin service                           103904 non-null  int64
14  Inflight service                          103904 non-null  int64
15  Cleanliness                               103904 non-null  int64
16  Departure Delay in Minutes                103904 non-null  int64
17  Arrival Delay in Minutes                  103904 non-null  float64
18  Gender_Male                               103904 non-null  uint8
19  Type of Travel_Personal Travel            103904 non-null  uint8
20  Class_Eco                                 103904 non-null  uint8
21  Class_Eco Plus                            103904 non-null  uint8
22  satisfaction_satisfied                     103904 non-null  uint8
dtypes: float64(1), int64(17), uint8(5)
memory usage: 15.6 MB
```

## Develop baseline model

In [4]:  1 **import** matplotlib.pyplot **as** plt  
 2 **import** seaborn **as** sns  
 3 %matplotlib inline  
 4  
 5 **from** sklearn.model\_selection **import** cross\_val\_predict, cross\_val\_score  
 6 **from** sklearn.metrics **import** plot\_confusion\_matrix, confusion\_matrix, class:

## fit estimator

```
In [5]: 1 logreg = LogisticRegression(solver='liblinear')
2 logreg.fit(X_train,y_train)
3
4 y_train_pred = cross_val_predict(logreg,X_train,y_train)
5
6 train_score = cross_val_score(logreg,X_train,y_train_pred,scoring='f1')
```

```
In [6]: 1 train_score.mean()
```

Out[6]: 0.9848444640204921

## print classification report

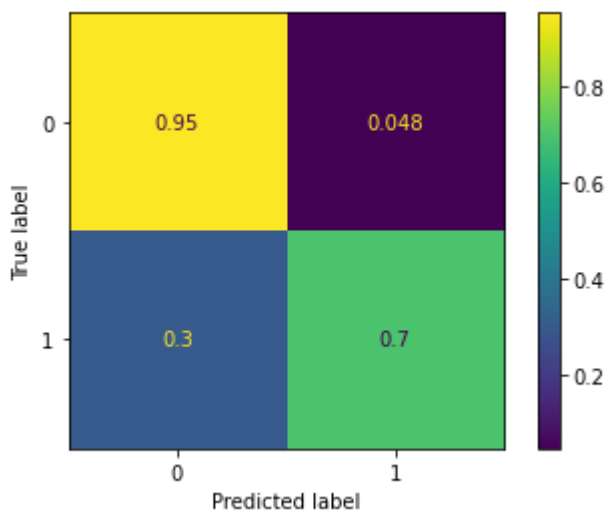
```
In [7]: 1 baseline_report = classification_report(y_train,y_train_pred,output_dict=True)
2 baseline_report = pd.DataFrame(baseline_report).iloc[:,0:3]
3 baseline_report
```

Out[7]:

	0	1	accuracy
<b>precision</b>	0.932692	0.765552	0.905076
<b>recall</b>	0.952604	0.692429	0.905076
<b>f1-score</b>	0.942543	0.727157	0.905076
<b>support</b>	84923.000000	18981.000000	0.905076

## plot confusion matrix

```
In [8]: 1 plot_confusion_matrix(logreg,X_train,y_train,normalize='true')
2 plt.show()
```



```
In [9]: ▶ 1 # Calculate the probability scores of each point in the training set
          2 y_train_score = logreg.decision_function(X_train)
          3
          4 # Calculate the fpr, tpr, and thresholds for the training set
          5 train_fpr, train_tpr, thresholds = roc_curve(y_train, y_train_score)
```

```
In [10]: ▶ 1 baseline_auc = auc(train_fpr, train_tpr)
```

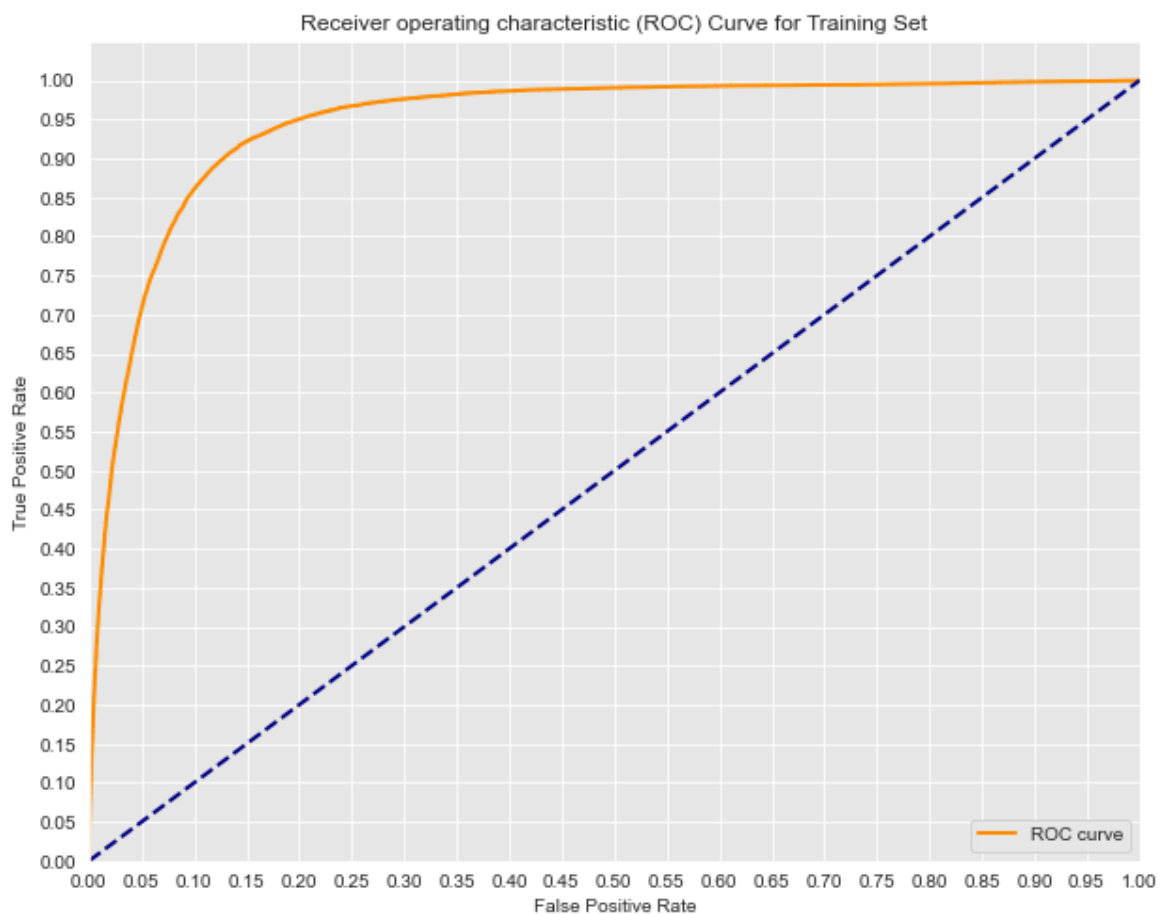
**calculate ROC AUC and plot curve**

```

In [11]: 1 # Seaborn's beautiful styling
2 sns.set_style('darkgrid', {'axes.facecolor': '0.9'})
3
4 # ROC curve for training set
5 plt.figure(figsize=(10, 8))
6 lw = 2
7 plt.plot(train_fpr, train_tpr, color='darkorange',
8          lw=lw, label='ROC curve')
9 plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
10 plt.xlim([0.0, 1.0])
11 plt.ylim([0.0, 1.05])
12 plt.yticks([i/20.0 for i in range(21)])
13 plt.xticks([i/20.0 for i in range(21)])
14 plt.xlabel('False Positive Rate')
15 plt.ylabel('True Positive Rate')
16 plt.title('Receiver operating characteristic (ROC) Curve for Training Set')
17 plt.legend(loc='lower right')
18 print('Training AUC: {}'.format(auc(train_fpr, train_tpr)))
19 plt.show()

```

Training AUC: 0.9471988366981169



## Prototype sampling methods

- SMOTE (synthetic over sampling)
- Tomek Links (under sampling against decision boundary)

- Near Miss (distance based under sampling)
- Edited Nearest Neighbors (under sampling against decision boundary)
- SMOTETomek (SMOTE/Tomek Link ensemble)
- SMOETENN (SMOTE/Edited Nearest Neighbors Ensemble)

## SMOTE

### *resample data*

```
In [12]: 1 from imblearn.over_sampling import SMOTE
          2 smote = SMOTE()
          3 X_smote, y_smote = smote.fit_resample(X_train, y_train)
          4
          5 print('Original dataset weights:', y_train.value_counts(normalize=True))
          6 print('Original dataset size:', len(y_train))
          7 print('\nResample dataset weights', y_smote.value_counts(normalize=True))
          8 print('Resample dataset size:', len(y_smote))
```

```
Original dataset weights: 0    0.817322
                          1    0.182678
Name: disloyal Customer, dtype: float64
Original dataset size: 103904
```

```
Resample dataset weights 1    0.5
                          0    0.5
Name: disloyal Customer, dtype: float64
Resample dataset size: 169846
```

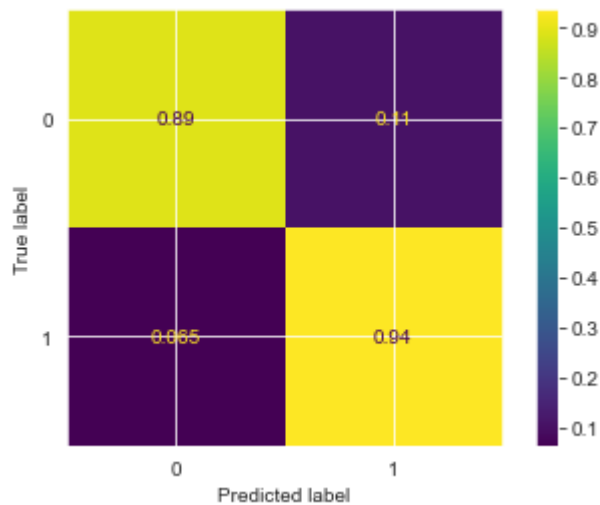
### **fit estimator**

```
In [13]: 1 logreg_smote = LogisticRegression(solver='liblinear')
          2 logreg_smote.fit(X_smote, y_smote)
          3
          4 smote_train_pred = logreg_smote.predict(X_smote)
          5
          6 smote_train_score = cross_val_score(logreg_smote, X_smote, smote_train_pred)
          7 print(smote_train_score.mean())
```

```
0.993825113914695
```

### **plot confusion matrix**

```
In [14]: 1 plot_confusion_matrix(logreg_smote,X_smote,y_smote,normalize='true')
2 plt.show()
```



print classification report

```
In [15]: 1 smote_report = classification_report(y_smote,smote_train_pred,output_dic
2 smote_report = pd.DataFrame(smote_report).iloc[:,0:3]
3 smote_report
```

Out[15]:

	0	1	accuracy
<b>precision</b>	0.932420	0.896163	0.913498
<b>recall</b>	0.891619	0.935377	0.913498
<b>f1-score</b>	0.911563	0.915350	0.913498
<b>support</b>	84923.000000	84923.000000	0.913498

calculate ROC AUC

```
In [16]: 1 # Calculate the probability scores of each point in the training set
2 y_smote_score = logreg_smote.decision_function(X_smote)
3 # Calculate the fpr, tpr, and thresholds for the training set
4 smote_fpr, smote_tpr, smote_thresholds = roc_curve(y_smote, y_smote_score)
5 smote_auc = auc(smote_fpr, smote_tpr)
6 print('Training AUC: {}'.format(smote_auc))
```

Training AUC: 0.9645822106754068

## Tomek Links

### resample data

```
In [17]: 1 from imblearn.under_sampling import TomekLinks
2 tl = TomekLinks(sampling_strategy='majority',n_jobs=3)
3 X_tl, y_tl = tl.fit_resample(X_train, y_train)
4
5 print('Original dataset weights:', y_train.value_counts(normalize=True))
6 print('Original dataset size:',len(y_train))
7 print('\nResample dataset weights', y_tl.value_counts(normalize=True))
8 print('Resample dataset size:', len(y_tl))
```

```
Original dataset weights: 0    0.817322
1    0.182678
Name: disloyal Customer, dtype: float64
Original dataset size: 103904
```

```
Resample dataset weights 0    0.809029
1    0.190971
Name: disloyal Customer, dtype: float64
Resample dataset size: 99392
```

### fit estimator

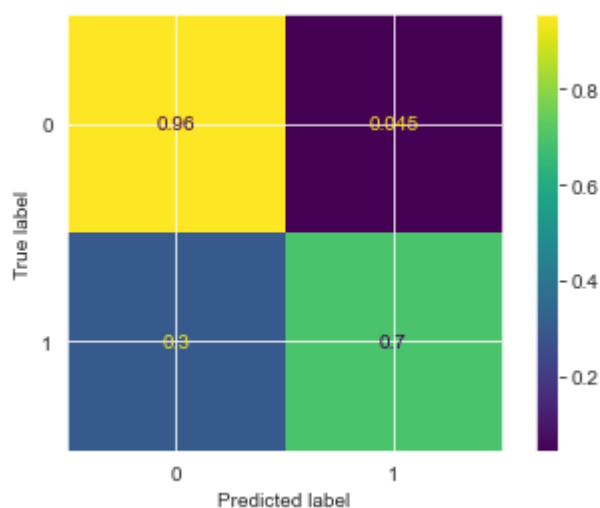
```
In [18]: 1 logreg_tl = LogisticRegression(solver='liblinear')
2 logreg_tl.fit(X_tl,y_tl)
3
4 tl_train_pred = logreg_tl.predict(X_tl)
5
6 tl_train_score = cross_val_score(logreg_tl,X_tl,y_tl,n_jobs=3,scoring='f1')
7 print(tl_train_score.mean())
```

```
0.7469474086997104
```

### plot confusion matrix



```
In [19]: 1 plot_confusion_matrix(logreg_tl,X_tl,y_tl,normalize='true')
        2 plt.show()
```



**print classification report**

```
In [20]: 1 tokek_report = classification_report(y_tl,tl_train_pred,output_dict=True)
        2 tokek_report = pd.DataFrame(tokek_report).iloc[:,0:3]
        3 tokek_report
```

Out[20]:

	0	1	accuracy
<b>precision</b>	0.930535	0.785719	0.905968
<b>recall</b>	0.955068	0.697961	0.905968
<b>f1-score</b>	0.942642	0.739244	0.905968
<b>support</b>	80411.000000	18981.000000	0.905968

**calculate ROC AUC**

```
In [21]: 1 # Calculate the probability scores of each point in the training set
2 y_tomek_score = logreg_smote.decision_function(X_t1)
3 # Calculate the fpr, tpr, and thresholds for the training set
4 tl_fpr , tl_tpr, tomek_thresholds = roc_curve(y_t1, y_tomek_score)
5 tl_auc = auc(tl_fpr, tl_tpr)
6 print('Training AUC: {}'.format(tl_auc))
```

Training AUC: 0.9434908347763292

## Near Miss

### resample the data

```
In [22]: 1 from imblearn.under_sampling import NearMiss
2 nm = NearMiss(sampling_strategy='all',n_jobs=3)
3 X_nm, y_nm = nm.fit_resample(X_train,y_train)
4
5 print('Original dataset weights:', y_train.value_counts(normalize=True))
6 print('Original dataset size:',len(y_train))
7 print('\nResample dataset weights', y_nm.value_counts(normalize=True))
8 print('Resample dataset size:', len(y_nm))
```

Original dataset weights: 0 0.817322  
1 0.182678  
Name: disloyal Customer, dtype: float64  
Original dataset size: 103904

Resample dataset weights 1 0.5  
0 0.5  
Name: disloyal Customer, dtype: float64  
Resample dataset size: 37962

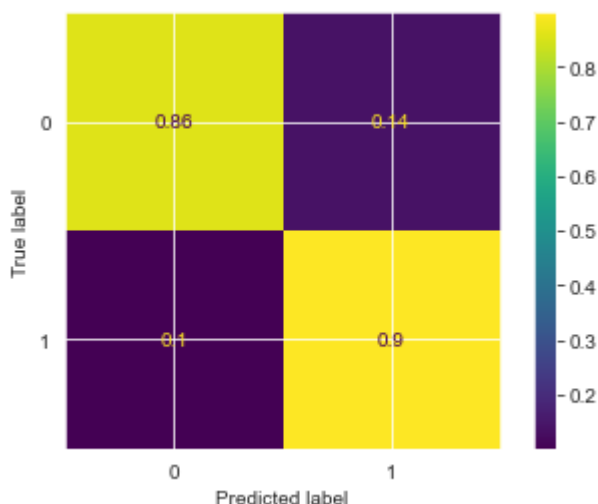
### fit estimator

```
In [23]: 1 logreg_nm = LogisticRegression(solver='liblinear')
2 logreg_nm.fit(X_nm,y_nm)
3
4 nm_train_pred = logreg_nm.predict(X_nm)
5
6 nm_train_score = cross_val_score(logreg_nm,X_nm,y_nm,n_jobs=3,scoring='f1')
7 print(nm_train_score.mean())
```

0.8708729312696878

### plot confusion matrix

```
In [24]: 1 plot_confusion_matrix(logreg_nm,X_nm,y_nm,normalize='true')
        2 plt.show()
```



### print classification report

```
In [25]: 1 nm_report = classification_report(y_nm,nm_train_pred,output_dict=True)
        2 nm_report = pd.DataFrame(nm_report).iloc[:,0:3]
        3 nm_report
```

```
Out[25]:
```

	0	1	accuracy
<b>precision</b>	0.894884	0.865828	0.879801
<b>recall</b>	0.860703	0.898899	0.879801
<b>f1-score</b>	0.877461	0.882053	0.879801
<b>support</b>	18981.000000	18981.000000	0.879801

### calculate ROC AUC

```
In [26]: 1 # Calculate the probability scores of each point in the training set
        2 y_nm_score = logreg_nm.decision_function(X_nm)
        3 # Calculate the fpr, tpr, and thresholds for the training set
        4 nm_fpr , nm_tpr, nm_thresholds = roc_curve(y_nm, y_nm_score )
        5 nm_auc = auc(nm_fpr, nm_tpr)
        6 print('Training AUC: {}'.format(nm_auc))
```

Training AUC: 0.9499421309957607

## Edited Nearest Neighbors

## resample the data

```
In [27]: 1 from imblearn.under_sampling import EditedNearestNeighbours
2 ENN = EditedNearestNeighbours(sampling_strategy='majority')
3 X_enh, y_enh = ENN.fit_resample(X_train,y_train)
4
5 print('Original dataset weights:', y_train.value_counts(normalize=True))
6 print('Original dataset size:',len(y_train))
7 print('\nResample dataset weights', y_enh.value_counts(normalize=True))
8 print('Resample dataset size:', len(y_enh))
```

```
Original dataset weights: 0    0.817322
1    0.182678
Name: disloyal Customer, dtype: float64
Original dataset size: 103904
```

```
Resample dataset weights 0    0.760619
1    0.239381
Name: disloyal Customer, dtype: float64
Resample dataset size: 79292
```

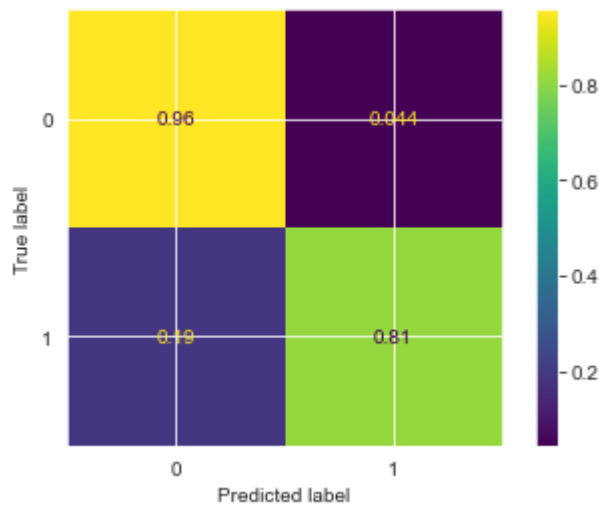
## fit estimator

```
In [28]: 1 logreg_ENN = LogisticRegression(solver='liblinear')
2 logreg_ENN.fit(X_enh,y_enh)
3
4 enn_train_pred = cross_val_predict(logreg_ENN,X_enh,y_enh,n_jobs=3)
5
6 enn_train_score = cross_val_score(logreg_ENN,X_enh,y_enh,n_jobs=3,scoring='accuracy')
7 print(enn_train_score.mean())
```

```
0.8329351715136429
```

## plot confusion matrix

```
In [29]: 1 plot_confusion_matrix(logreg_ENN,X_enn,y_enn,normalize='true')
          2 plt.show()
```



### print classification matrix

```
In [30]: 1 enn_report = classification_report(y_enn,enn_train_pred,output_dict=True)
          2 enn_report = pd.DataFrame(enn_report).iloc[:,0:3]
          3 enn_report
```

```
Out[30]:
```

	0	1	accuracy
<b>precision</b>	0.942915	0.850206	0.921606
<b>recall</b>	0.954735	0.816343	0.921606
<b>f1-score</b>	0.948788	0.832930	0.921606
<b>support</b>	60311.000000	18981.000000	0.921606

### calculate ROC AUC

```
In [31]: 1 # Calculate the probability scores of each point in the training set
          2 y_enn_score = logreg_nm.decision_function(X_enn)
          3 # Calculate the fpr, tpr, and thresholds for the training set
          4 enn_fpr , enn_tpr, enn_thresholds = roc_curve(y_enn, y_enn_score )
          5 enn_auc = auc(enn_fpr, enn_tpr)
          6 print('Training AUC: {}'.format(enn_auc))
```

Training AUC: 0.7817683396992051

## SMOTETomek

### resample the data

```
In [32]: 1 from imblearn.combine import SMOTETomek
2 SMOTek = SMOTETomek(sampling_strategy='all', smote=smote, totek=tl, n_jobs=
3 X_smotek, y_smotek = SMOTek.fit_resample(X_train, y_train)
4
5 print('Original dataset weights:', y_train.value_counts(normalize=True))
6 print('Original dataset size:', len(y_train))
7 print('\nResample dataset weights', y_smotek.value_counts(normalize=True))
8 print('Resample dataset size:', len(y_smotek))
```

```
Original dataset weights: 0    0.817322
1    0.182678
Name: disloyal Customer, dtype: float64
Original dataset size: 103904
```

```
Resample dataset weights 1    0.501423
0    0.498577
Name: disloyal Customer, dtype: float64
Resample dataset size: 169364
```

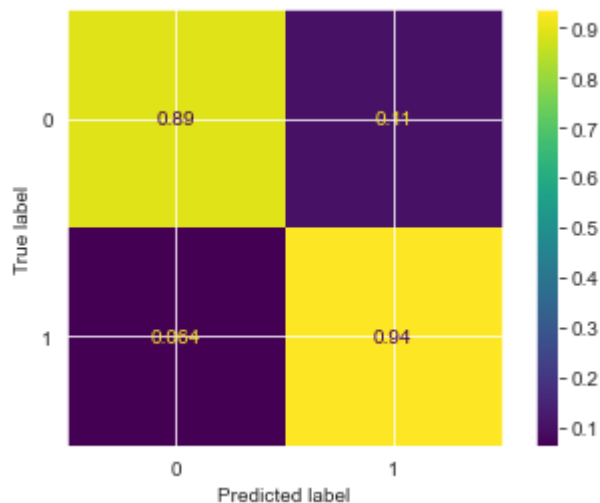
### fit estimator

```
In [33]: 1 logreg_SMOTek = LogisticRegression(solver='liblinear')
2 logreg_SMOTek.fit(X_smotek, y_smotek)
3
4 smotek_train_pred = logreg_SMOTek.predict(X_smotek)
5
6 smotek_train_score = cross_val_score(logreg_SMOTek, X_smotek, y_smotek, n_jo
7 print(smotek_train_score.mean())
```

```
0.9092644474215696
```

### plot confusion matrix

```
In [34]: 1 plot_confusion_matrix(logreg_SMOTek,X_smotek,y_smotek,normalize='true')
2 plt.show()
```



**print classification report**

```
In [35]: 1 smotek_report = classification_report(y_nm,nm_train_pred,output_dict=True)
2 smotek_report = pd.DataFrame(smotek_report).iloc[:,0:3]
3 smotek_report
```

```
Out[35]:
```

	0	1	accuracy
<b>precision</b>	0.894884	0.865828	0.879801
<b>recall</b>	0.860703	0.898899	0.879801
<b>f1-score</b>	0.877461	0.882053	0.879801
<b>support</b>	18981.000000	18981.000000	0.879801

**calculate ROC AUC**

```
In [36]: 1 # Calculate the probability scores of each point in the training set
2 y_smotek_score = logreg_SMOTek.decision_function(X_smotek)
3 # Calculate the fpr, tpr, and thresholds for the training set
4 smotek_fpr , smotek_tpr, smotek_thresholds = roc_curve(y_smotek, y_smotek_score)
5 smotek_auc = auc(smotek_fpr, smotek_tpr)
6 print('Training AUC: {}'.format(smotek_auc))
```

Training AUC: 0.9647099965956508

## SMOTENN

**resmaple the data**

```
In [37]: 1 from imblearn.combine import SMOTEENN
2 SMN = SMOTEENN(sampling_strategy='all', smote=smote, enn=ENN, n_jobs=3)
3 X_smn, y_smn = SMN.fit_resample(X_train, y_train)
4
5 print('Original dataset weights:', y_train.value_counts(normalize=True))
6 print('Original dataset size:', len(y_train))
7 print('\nResample dataset weights', y_smn.value_counts(normalize=True))
8 print('Resample dataset size:', len(y_smn))
```

```
Original dataset weights: 0    0.817322
1    0.182678
Name: disloyal Customer, dtype: float64
Original dataset size: 103904
```

```
Resample dataset weights 1    0.618792
0    0.381208
Name: disloyal Customer, dtype: float64
Resample dataset size: 137240
```

### fit estimator

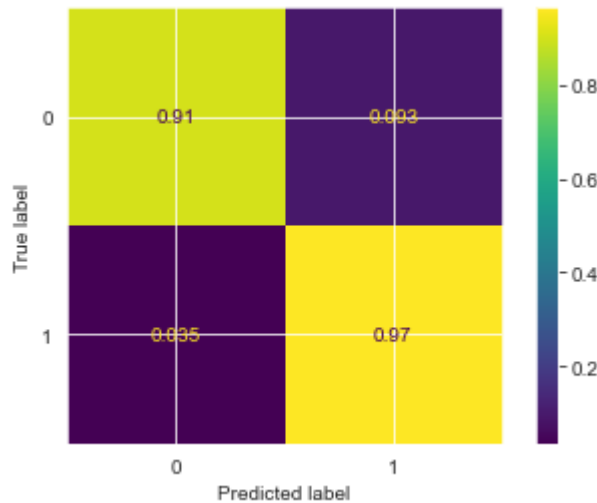
```
In [38]: 1 logreg_SMOTENN = LogisticRegression(solver='liblinear')
2 logreg_SMOTENN.fit(X_smn, y_smn)
3
4 smn_train_pred = logreg_SMOTENN.predict(X_smn)
5
6 smn_train_score = cross_val_score(logreg_SMOTENN, X_smn, y_smn, n_jobs=3, sc
7 print(smn_train_score.mean())
```

```
0.9504230038099541
```

### plot confusion matrix



```
In [39]: 1 plot_confusion_matrix(logreg_SMOTENN,X_smn,y_smn,normalize='true')
2 plt.show()
```



### print classification report

```
In [40]: 1 SMOTENN_report = classification_report(y_smn,smn_train_pred,output_dict=
2 SMOTENN_report = pd.DataFrame(SMOTENN_report).iloc[:,0:3]
3 SMOTENN_report
```

```
Out[40]:
```

	0	1	accuracy
<b>precision</b>	0.941517	0.943891	0.94302
<b>recall</b>	0.906856	0.965298	0.94302
<b>f1-score</b>	0.923862	0.954475	0.94302
<b>support</b>	52317.000000	84923.000000	0.94302

### calculate ROC AUC

```
In [41]: 1 # Calculate the probability scores of each point in the training set
2 SMOTENN_score = logreg_SMOTENN.decision_function(X_smn)
3 # Calculate the fpr, tpr, and thresholds for the training set
4 smotenn_fpr , smotenn_tpr, smotenn_thresholds = roc_curve(y_smn, SMOTENN_score)
5 smotenn_auc = auc(smotenn_fpr, smotenn_tpr)
6 print('Training AUC: {}'.format(smotenn_auc))
```

Training AUC: 0.9810652094683898

## Validate sample and select method(s)

- First tried SMOTENN because it had the highest F1, validation however had a ROC AUC of about .06, indicating the model is overfitting the data.

- Next I am trying Near Miss because it only undersamples the majority and had a training F1 score between SMOTE and Edited Nearest Neighbors, the two methods used in SMOTENN. Not synthesizing new data, and only undersampling I expect will reduce overfitting.

```
In [42]: 1 # load and clean validation set
        2 X_test, y_test = data_cleaner("../data/train.csv.zip")
```

## SMOTENN CV Test

### fit cross validated model

```
In [43]: 1 logreg_SMOTENN_CV = LogisticRegressionCV(solver='liblinear',n_jobs=3)
        2 logreg_SMOTENN_CV.fit(X_smn,y_smn)
```

Out[43]: LogisticRegressionCV(n\_jobs=3, solver='liblinear')

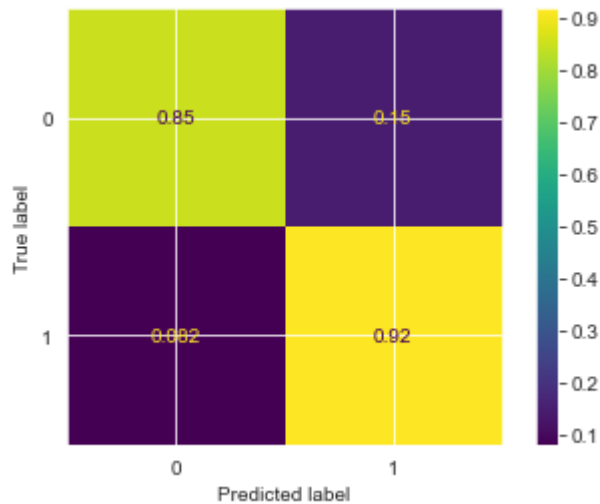
```
In [44]: 1 smn_test_pred = logreg_SMOTENN_CV.predict(X_test)
        2 smn_test_report = classification_report(y_test,smn_test_pred,output_dict=True)
        3 SMOTENN_CV_report = pd.DataFrame(smn_test_report).iloc[:,0:3]
        4 SMOTENN_CV_report
```

Out[44]:

	0	1	accuracy
<b>precision</b>	0.978902	0.577867	0.862517
<b>recall</b>	0.850111	0.918023	0.862517
<b>f1-score</b>	0.909972	0.709270	0.862517
<b>support</b>	84923.000000	18981.000000	0.862517

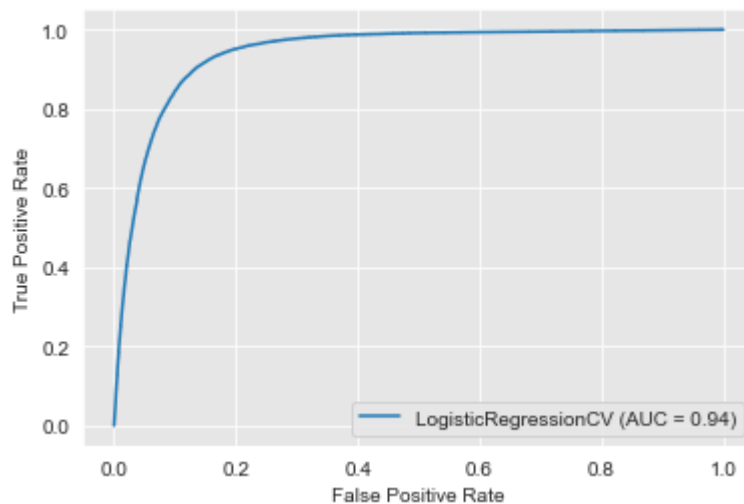
### plot confusion matrix

```
In [45]: 1 # plot confusion matrix
2 plot_confusion_matrix(logreg_SMOTENN_CV,X_test,y_test,normalize='true')
3 plt.show()
```



### plot ROC AUC

```
In [46]: 1 from sklearn.metrics import plot_roc_curve
2 plot_roc_curve(logreg_SMOTENN_CV,X_test,y_test)
3 plt.show()
```



```
In [47]: 1 # Calculate the probability scores of each point in the training set
2 SMOTENN_CV_score = logreg_SMOTENN_CV.decision_function(X_test)
3 # Calculate the fpr, tpr, and thresholds for the training set
4 smotenn_cv_fpr , smotenn_cv_tpr, smotenn_cv_thresholds = roc_curve(y_test, SMOTENN_CV_score)
5 smotenn_cv_auc = auc(smotenn_cv_fpr, smotenn_cv_tpr)
6 print('SMOTENN Test AUC: {}'.format(smotenn_cv_auc))
```

SMOTENN Test AUC: 0.9419915299042956

## Near Miss CV Test

### fit cross validated model

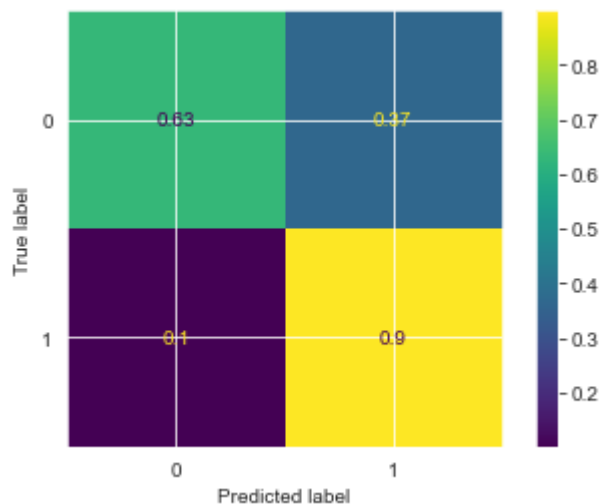
```
In [48]: 1 logreg_nm_cv = LogisticRegressionCV(solver='liblinear',n_jobs=3)
2 logreg_nm_cv.fit(X_nm,y_nm)
3
4 nm_test_pred = logreg_nm_cv.predict(X_test)
5 nm_test_report = classification_report(y_test,nm_test_pred,output_dict=True)
6 nm_test_report = pd.DataFrame(nm_test_report).iloc[:,0:3]
7 nm_test_report
```

```
Out[48]:
```

	0	1	accuracy
<b>precision</b>	0.965475	0.354056	0.681985
<b>recall</b>	0.633562	0.898635	0.681985
<b>f1-score</b>	0.765071	0.507974	0.681985
<b>support</b>	84923.000000	18981.000000	0.681985

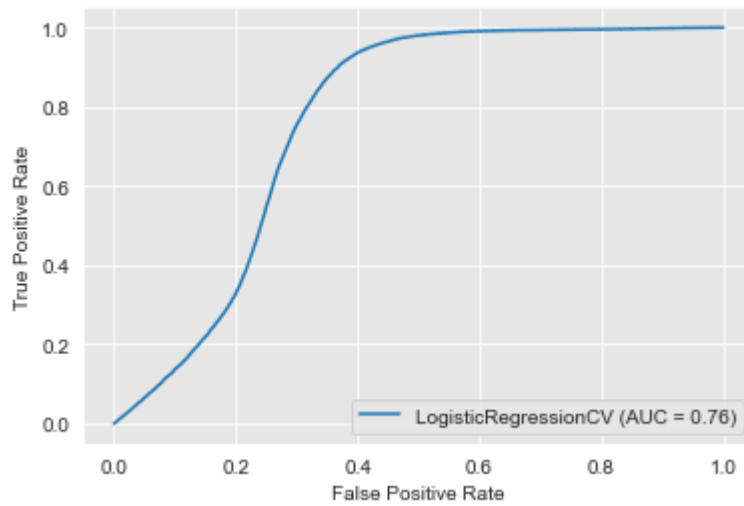
### plot confusion matrix

```
In [49]: 1 plot_confusion_matrix(logreg_nm_cv,X_test,y_test,normalize='true')
2 plt.show()
```



### plot ROC AUC

```
In [50]: 1 plot_roc_curve(logreg_nm_cv,X_test,y_test)
2 plt.show()
```



```
In [51]: 1 # Calculate the probability scores of each point in the training set
2 nm_cv_score = logreg_nm_cv.decision_function(X_test)
3 # Calculate the fpr, tpr, and thresholds for the training set
4 nm_cv_fpr , nm_cv_tpr, nm_cv_thresholds = roc_curve(y_test, nm_cv_score)
5 nm_cv_auc = auc(nm_cv_fpr, nm_cv_tpr)
6 print('Near Miss Test AUC: {}'.format(nm_cv_auc))
```

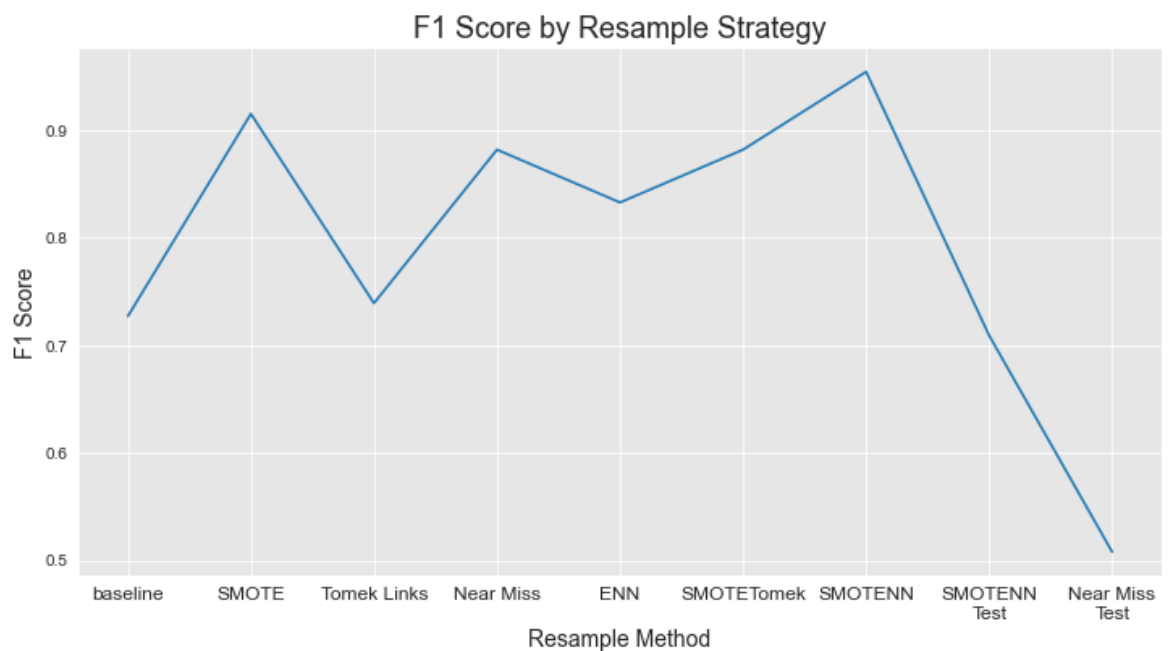
Near Miss Test AUC: 0.7630068283211038

## Visualize F1, ROC AUC iterative performance

```
In [52]: 1 f1_scores = [  
2     baseline_report.iloc[2,1],  
3     smote_report.iloc[2,1],  
4     totek_report.iloc[2,1],  
5     nm_report.iloc[2,1],  
6     enn_report.iloc[2,1],  
7     smotek_report.iloc[2,1],  
8     SMOTENN_report.iloc[2,1],  
9     SMOTENN_CV_report.iloc[2,1],  
10    nm_test_report.iloc[2,1]  
11 ]  
12 np.mean(f1_scores)
```

Out[52]: 0.7945008205358247

```
In [59]: 1 plt.figure(figsize=(12,6))  
2 x = ['baseline', 'SMOTE', 'Tomek Links', 'Near Miss', 'ENN', 'SMOTETomek', 'SMO  
3 y = f1_scores  
4 plt.plot(x,y)  
5 plt.ylabel('F1 Score',fontsize=14)  
6 plt.xlabel('Resample Method',fontsize=14)  
7 plt.xticks(fontsize=12)  
8 plt.title("F1 Score by Resample Strategy",fontsize=18)  
9 plt.show()
```



```
In [60]: ▶ 1 roc_auc_scores = [  
2     baseline_auc,  
3     smote_auc,  
4     tl_auc,  
5     nm_auc,  
6     enn_auc,  
7     smotek_auc,  
8     smotenn_auc,  
9     smotenn_cv_auc,  
10    nm_cv_auc  
11 ]  
12 np.mean(roc_auc_scores)
```

Out[60]: 0.9153062130149175

**plot ROC AUC curves**

```

In [69]: ▶ 1 plt.figure(figsize=(10,8))
2 lw = 2
3
4 print('Baseline Model AUC: {}'.format(baseline_auc))
5 print('SMOTE resample AUC: {}'.format(smote_auc))
6
7 plt.plot(train_fpr, train_tpr, color='darkorange',
8          lw=lw, label='Baseline model ROC curve')
9 plt.plot(smote_fpr, smote_tpr, color='blue',
10          lw=lw, label='SMOTE ROC curve')
11
12
13 print('Tomek Links AUC: {}'.format(tl_auc))
14 print('Near Miss AUC: {}'.format(nm_auc))
15
16 plt.plot(tl_fpr, tl_tpr, color='yellow',
17          lw=lw, label='Tomek Links ROC curve')
18 plt.plot(nm_fpr, nm_tpr, color='gold',
19          lw=lw, label='Near Miss ROC curve')
20
21
22 print('Edited Nearest Neighbors AUC: {}'.format(enn_auc))
23 print('SMOTETomek AUC: {}'.format(smotek_auc))
24
25 plt.plot(enn_fpr, enn_tpr, color='red',
26          lw=lw, label='Edited Nearest Neighbors ROC curve')
27 plt.plot(smotek_fpr, smotek_tpr, color='purple',
28          lw=lw, label='SMOTETomek ROC curve')
29
30
31 print('SMOTEENN AUC: {}'.format(smotenn_auc))
32 plt.plot(smotenn_fpr, smotenn_tpr, color='darkblue',
33          lw=lw, label='SMOTEENN ROC curve')
34
35
36 # Formatting
37 plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
38 plt.xlim([0.0, 1.0])
39 plt.ylim([0.0, 1.05])
40 plt.yticks([i/20.0 for i in range(21)])
41 plt.xticks([i/20.0 for i in range(21)])
42 plt.xlabel('False Positive Rate', fontsize=14)
43 plt.ylabel('True Positive Rate', fontsize=14)
44 plt.title('Training ROC Curve by Resample Strategy', fontsize=18)
45 plt.legend(loc="lower right")
46 plt.show()

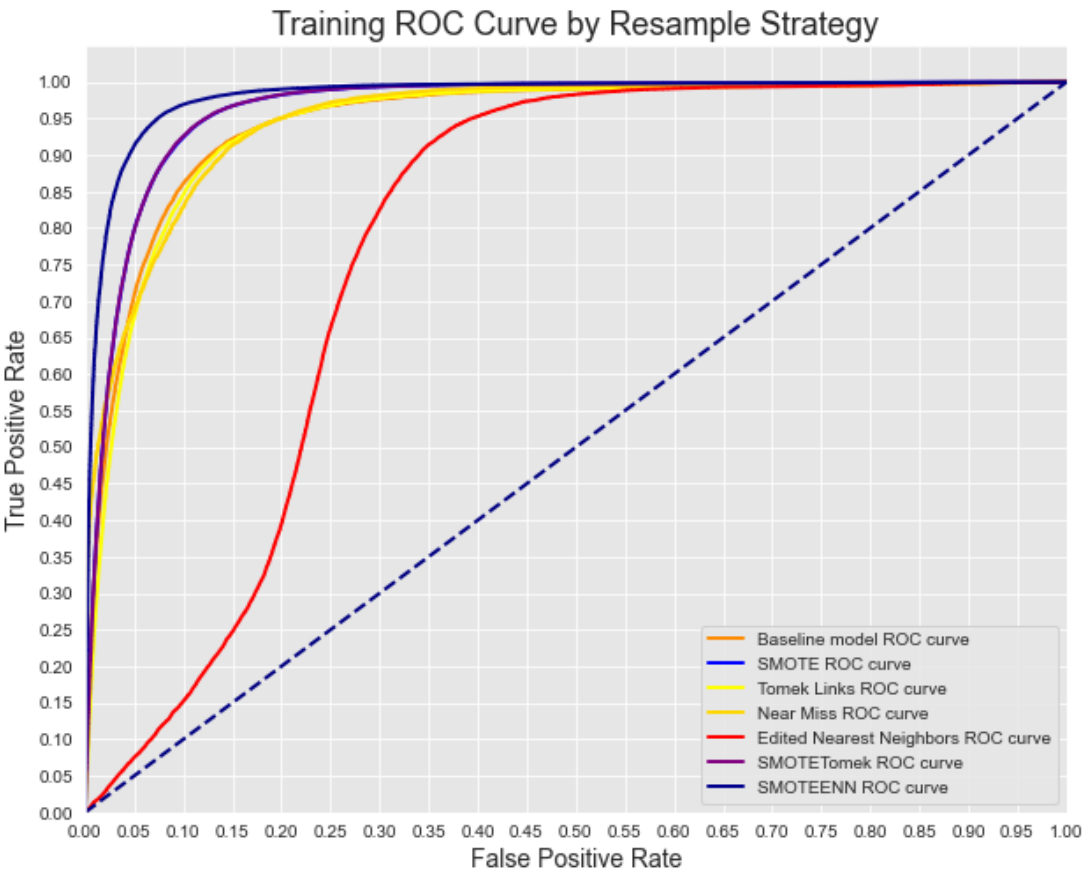
```

```

Baseline Model AUC: 0.9471988366981169
SMOTE resample AUC: 0.9645822106754068
Tomek Links AUC: 0.9434908347763292
Near Miss AUC: 0.9499421309957607
Edited Nearest Neighbors AUC: 0.7817683396992051
SMOTETomek AUC: 0.9647099965956508
SMOTEENN AUC: 0.9810652094683898

```

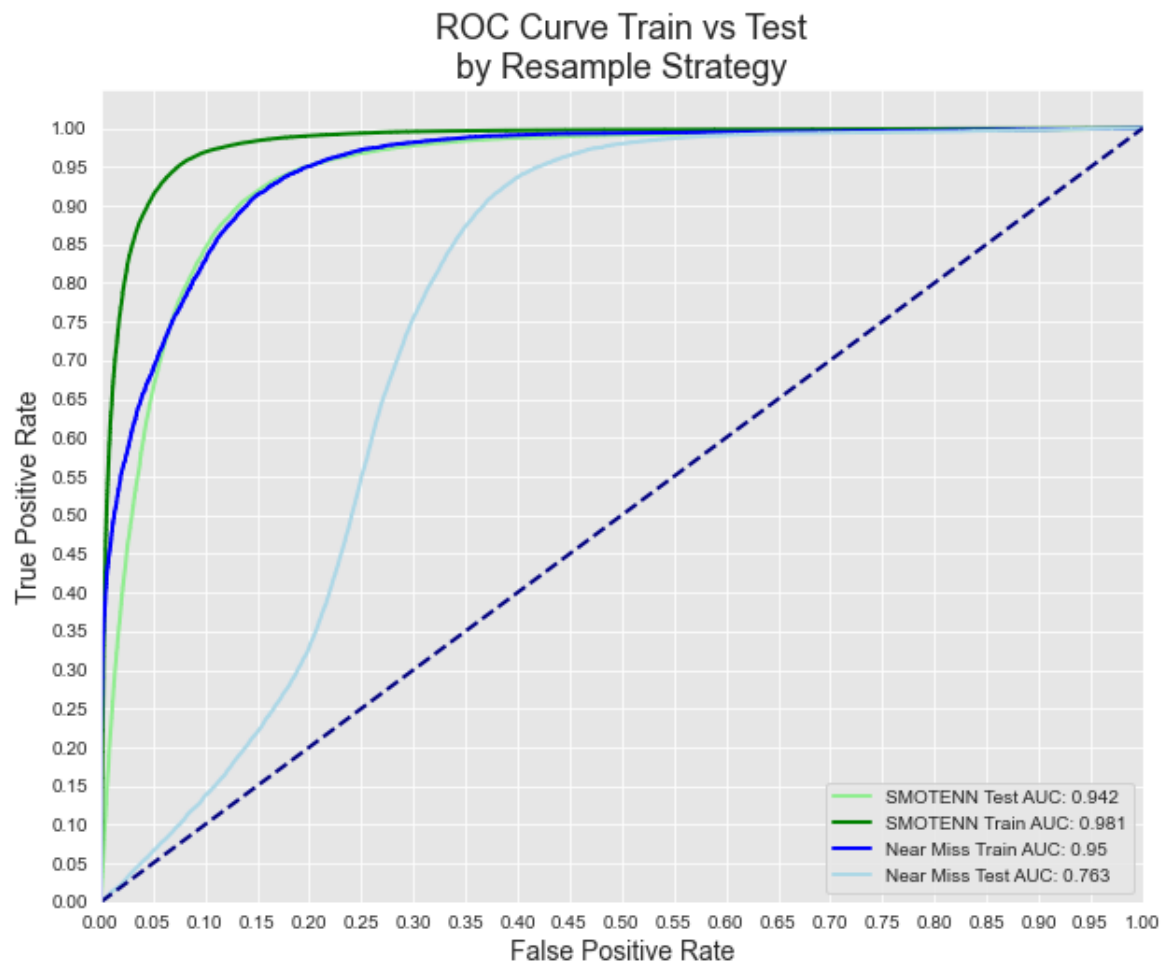




```

In [73]: ▶ 1 plt.figure(figsize=(10,8))
2 lw = 2
3
4 plt.plot(smotenn_cv_fpr, smotenn_cv_tpr, color='lightgreen',
5          lw=lw, label=f'SMOTENN Test AUC: {round(smotenn_cv_auc,3)}')
6 plt.plot(smotenn_fpr,smotenn_tpr,color='green',
7          lw=lw, label=f'SMOTENN Train AUC: {round(smotenn_auc,3)}')
8
9 plt.plot(nm_fpr, nm_tpr, color='blue',
10          lw=lw, label=f'Near Miss Train AUC: {round(nm_auc,3)}')
11 plt.plot(nm_cv_fpr, nm_cv_tpr, color='lightblue',
12          lw=lw, label=f'Near Miss Test AUC: {round(nm_cv_auc,3)}')
13
14 # Formatting
15 plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
16 plt.xlim([0.0, 1.0])
17 plt.ylim([0.0, 1.05])
18 plt.yticks([i/20.0 for i in range(21)])
19 plt.xticks([i/20.0 for i in range(21)])
20 plt.xlabel('False Positive Rate',fontsize=14)
21 plt.ylabel('True Positive Rate',fontsize=14)
22 plt.title('ROC Curve Train vs Test\nby Resample Strategy',fontsize=18)
23 plt.legend(loc="lower right")
24 plt.show()

```



## Final Observations

The ensemble of oversampling with SMOTE and undersampling with Edited Nearest Neighbors is clearly the strongest performer among the resampling methods explored in this notebook, mainly based on F1 score and ROC AUC. I bothered to also test Near Miss, because it is the best performing non-ensemble method and it only undersamples (there is no data used that didn't exist in the first place); I won't be moving forward with Near Miss because it under-performed against SMOTEENN consistently with the training models.

The next step is to write a function in `preprocessor.py` to make this exact resample strategy portable between notebooks.

In the next development notebook I will use the preprocessing from EDA, and the SMOTEENN resampling method found here to train and optimize via gridsearching a decision tree and/or random forest. To use as my final model.