

Models

In [38]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib
4 import matplotlib.pyplot as plt
5 import itertools
6 from imblearn.over_sampling import SMOTE
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.preprocessing import StandardScaler
10
11 import statsmodels.api as sm
12 from statsmodels.api import OLS
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.linear_model import LogisticRegressionCV
15 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
16 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
17 from sklearn.neighbors import KNeighborsClassifier
18 from sklearn.tree import DecisionTreeClassifier
19 from sklearn.ensemble import RandomForestClassifier
20 from sklearn.ensemble import AdaBoostClassifier
21
22 from sklearn.model_selection import train_test_split
23 from sklearn.model_selection import cross_val_score
24
25 from sklearn.metrics import precision_recall_fscore_support
26 from sklearn.metrics import confusion_matrix
27 from sklearn.metrics import accuracy_score
28 from sklearn.metrics import auc, roc_curve
29
30 import warnings
31 warnings.filterwarnings('ignore')
32
33 import keras
34 from keras.models import Sequential
35 from keras.layers import Dense
36 from keras.utils import to_categorical
37
38 import seaborn as sns
39 pd.set_option('display.width', 500)
40 pd.set_option('display.max_columns', 500)
41
42 % matplotlib inline
```

Data Preparation

There are too many observations in the dataset, which will take hours or even days to fit for some sophisticated algorithms. Thus, we decided to take a stratified sample on year of raw dataset in the model fitting stage.

```
In [2]: 1 df_fe = pd.read_csv("data/output_fe.csv")
        2 print("Shape of the raw input file: {}".format(df_fe.shape))
```

Shape of the raw input file: (95765, 196)

```
In [3]: 1 df_clean = df_fe.copy(deep=True)
        2 df_clean = df_fe.groupby('year', group_keys=False).apply(lambda x: x.sample(f
        3 df_clean.drop(['year'], axis=1, inplace=True)
        4
        5 print("Shape of the clean input file: {}".format(df_clean.shape))
```

Shape of the clean input file: (4789, 195)

Stratified Sampling

We split training and test dataset by stratifying on the response variable.

```
In [4]: 1 df_train, df_test = train_test_split(df_clean,
        2                                     test_size = .5,
        3                                     stratify = df_clean['response'],
        4                                     random_state=90)
        5
        6 print("Shape of the training set: {}".format(df_train.shape))
        7 print("Shape of the test set: {}".format(df_test.shape))
```

Shape of the training set: (2394, 195)

Shape of the test set: (2395, 195)

```
In [5]: 1 def split_columns(df, target_col, drop_columns):
        2     # Get the response variable
        3     y_train = df[[target_col]]
        4
        5     # Drop the required columns
        6     X_train = df.drop(drop_columns, axis=1)
        7
        8     return X_train, y_train
```

```
In [6]: 1 X_train, y_train = split_columns(df_train, target_col='response', drop_column
        2 X_test, y_test = split_columns(df_test, target_col='response', drop_columns=
```

Standardization

We standardize all the predictors that are not dummy variables.

```
In [7]: 1 def scale_datasets(train_data, test_data, cols_to_scale):
2         train = train_data.copy()
3         test = test_data.copy()
4
5         # Fit the scaler on the training data
6         scaler = StandardScaler().fit(train[cols_to_scale])
7
8         # Scale both the test and training data.
9         train[cols_to_scale] = scaler.transform(train[cols_to_scale])
10        test[cols_to_scale] = scaler.transform(test[cols_to_scale])
11
12        return train, test
```

```
In [8]: 1 X_train, X_test = scale_datasets(X_train, X_test, list(X_train.columns))
```

Custom Functions

```
In [9]: 1 def plot_confusion_matrix(cm, classes,
2                             normalize=False,
3                             title='Confusion matrix',
4                             cmap=plt.cm.Blues,
5                             fontsize=16):
6
7     """
8     This function plots the confusion matrix.
9     Normalization can be applied by setting `normalize=True`.
10    """
11    if normalize:
12        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
13
14    plt.imshow(cm, interpolation='nearest', cmap=cmap)
15    plt.title(title, fontsize=fontsize)
16    plt.colorbar()
17    tick_marks = np.arange(len(classes))
18    plt.xticks(tick_marks, classes, fontsize=fontsize)
19    plt.yticks(tick_marks, classes, fontsize=fontsize)
20
21    fmt = '.2f' if normalize else 'd'
22    thresh = cm.max() / 2.
23    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
24        plt.text(j, i, format(cm[i, j], fmt),
25                horizontalalignment="center",
26                color="white" if cm[i, j] > thresh else "black",
27                fontsize=fontsize)
28
29    plt.ylabel('True label', fontsize=fontsize)
30    plt.xlabel('Predicted label', fontsize=fontsize)
31    plt.tight_layout()
```

Classification of Good and Bad Loans

Baseline Model

For classification, a simple baseline is always predicting the most common class, which is good loans in our dataset.

```
In [10]: 1 base_train_acc = y_train.response.value_counts()[0] / len(y_train)
2 base_test_acc = y_test.response.value_counts()[0] / len(y_test)
3 print("Baseline model accuracy in training set: {:.2%}".format(base_train_acc))
4 print("Baseline model accuracy in test set: {:.2%}".format(base_test_acc))
```

Baseline model accuracy in training set: 84.17%

Baseline model accuracy in test set: 84.18%

Oversampling

As we can see from the baseline model accuracy, the lending club dataset is an imbalanced one with 84% of majority class and only 16% minority class. A common problem with imbalanced dataset is that the model will simply predict the majority class with a high accuracy score and ignore the minority class. However, accuracy or precision may not be the only concern we have as data scientists.

For example, if we were building a model for cancer detection, we would want to capture all the patients that do have cancer, even at the cost of some misclassification of healthy people, as these people can be examined further by doctors. As investors on lending club, what we really care about is the default risk. The key question we need to ask ourselves is that among all the bad loans, how many of them we can predict correctly? This is the so called recall rate.

Oversampling is a technique that can deal with this imbalanced dataset. The basic idea is to oversample the minority class, so that the model can achieve higher recall at the cost of precision, and that's exactly what we want. In the modeling stage, We used Synthetic Minority Oversampling Technique (SMOTE) to oversample the training set.

```
In [11]: 1 sm = SMOTE(random_state=12, ratio = 1.0)
2 X_train_sm, y_train_sm = sm.fit_sample(X_train, y_train)
3 print("Oversampled Training Set:\nNumber of bad loans: {}\nNumber of good loans: {}".format(
4     np.count_nonzero(y_train_sm == 1), np.count_nonzero(y_train_sm == 0))
5 )
```

Oversampled Training Set:

Number of bad loans: 2015

Number of good loans: 2015

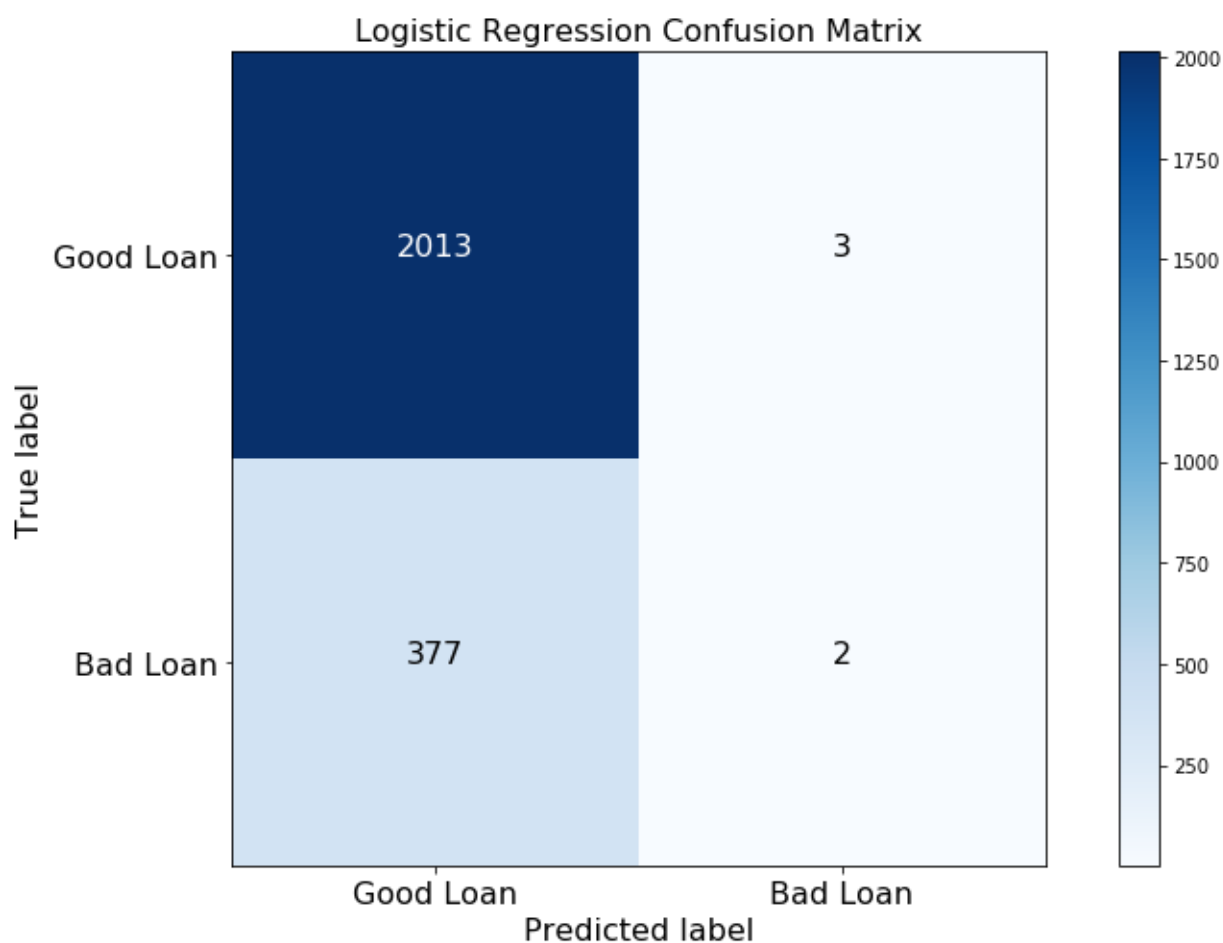
Logistic Regression

Raw training set

```
In [12]: 1 # Raw training set
2 logit = LogisticRegressionCV(cv=5, random_state=0, penalty='l2').fit(X_train,
3
4 train_acc = logit.score(X_train, y_train)
5 test_acc = logit.score(X_test, y_test)
6 report_lr = precision_recall_fscore_support(y_test, logit.predict(X_test), av
7
8 logit_results = {
9     'model': 'Logistic',
10    'train_acc': train_acc,
11    'test_acc': test_acc,
12    'precision': report_lr[0],
13    'recall': report_lr[1],
14    'F1': report_lr[2]
15 }
16
17 print('Logistic Regression: accuracy on train={:.2%}, test={:.2%}, precision=
18       format(logit_results['train_acc'],
19               logit_results['test_acc'],
20               logit_results['precision'],
21               logit_results['recall'],
22               logit_results['F1']))
```

Logistic Regression: accuracy on train=84.34%, test=84.13%, precision=0.40, recall=0.01, F1=0.01

```
In [13]: 1 # Plot confusion matrix
2 cnf_matrix = confusion_matrix(y_test, logit.predict(X_test))
3 np.set_printoptions(precision=2)
4
5 plt.figure(figsize=(11,7))
6 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normaliz
7                        title='Logistic Regression Confusion Matrix')
8
9 plt.show()
```

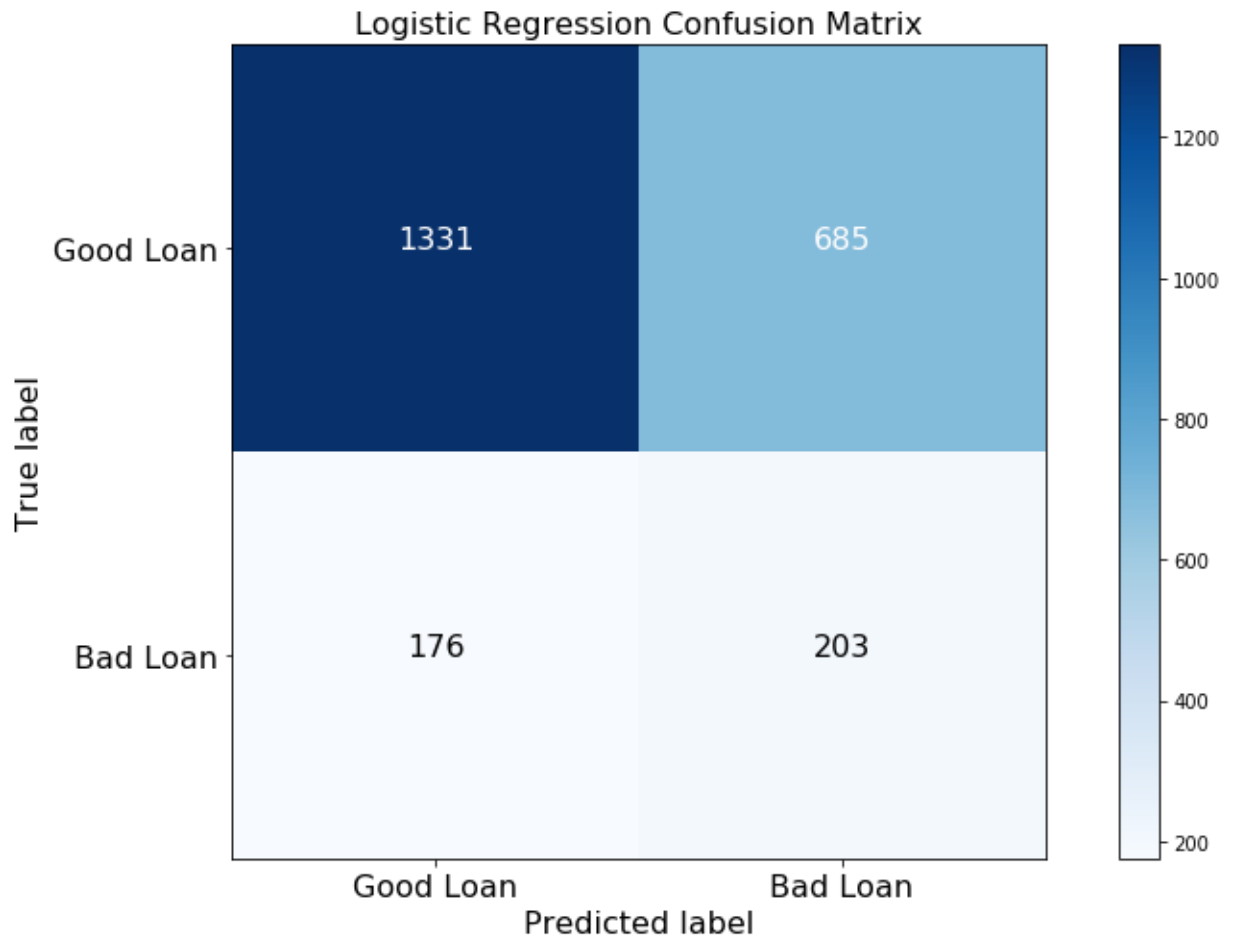


Oversampled Training set

```
In [33]: 1 # SMOTE training set
2 logit_sm = LogisticRegressionCV(cv=5, random_state=0, penalty='l2').fit(X_train_sm, y_train_sm)
3
4 train_acc, test_acc = logit_sm.score(X_train_sm, y_train_sm), logit_sm.score(X_test_sm, y_test_sm)
5 report_lr = precision_recall_fscore_support(y_test, logit_sm.predict(X_test), average='weighted', zero_division=0)
6
7 logit_sm_results = {
8     'model': 'Logistic',
9     'train_acc': train_acc,
10    'test_acc': test_acc,
11    'precision': report_lr[0],
12    'recall': report_lr[1],
13    'F1': report_lr[2]
14 }
15
16 print('Logistic Regression: accuracy on train={:.2%}, test={:.2%}, precision={:.2%}, recall={:.2%}, F1={:.2%}'.format(
17     train_acc, test_acc, report_lr[0], report_lr[1], report_lr[2]))
18
19 # Print the results
20
21
```

Logistic Regression: accuracy on train=73.05%, test=64.05%, precision=0.23, recall=0.54, F1=0.32

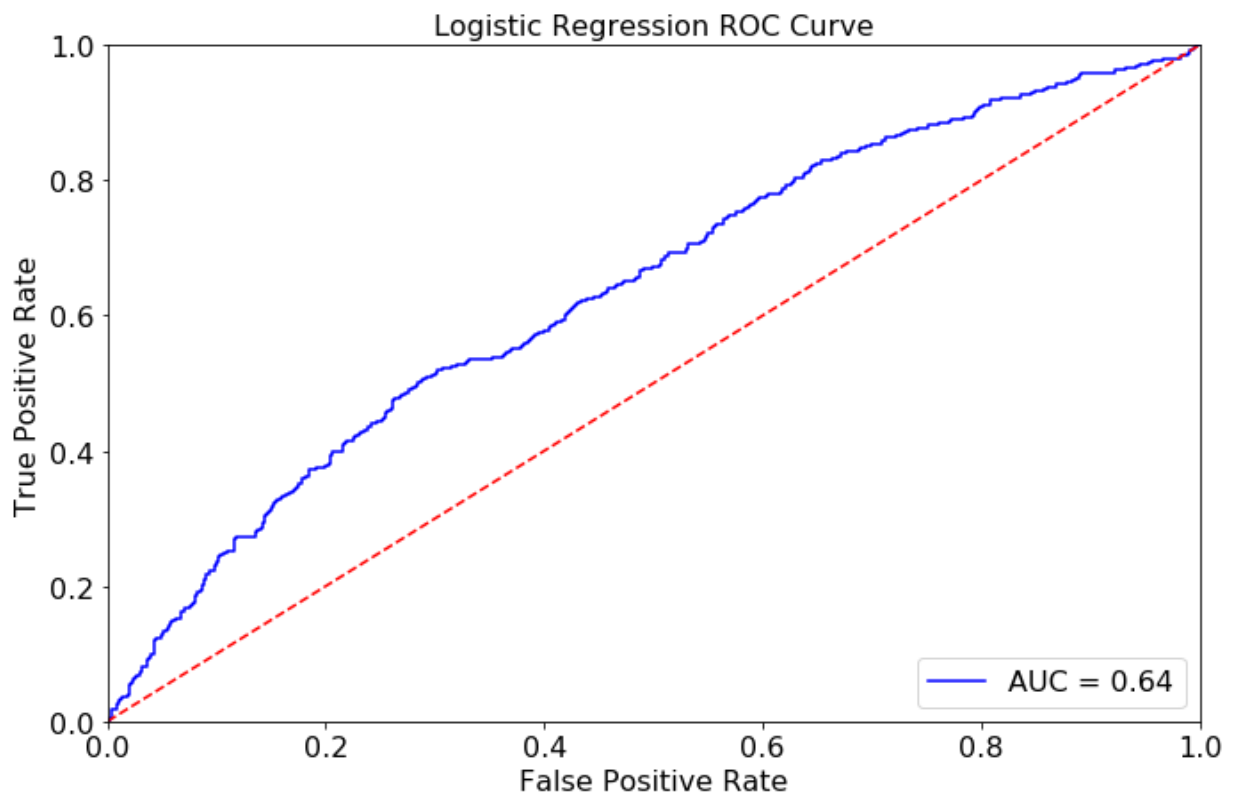
```
In [34]: 1 # Plot confusion matrix
2 cnf_matrix = confusion_matrix(y_test, logit_sm.predict(X_test))
3 np.set_printoptions(precision=2)
4
5 plt.figure(figsize=(11,7))
6 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normaliz
7                        title='Logistic Regression Confusion Matrix')
8
9 plt.show()
```




```

In [48]: 1 # Plot ROC curve
2 probs = logit_sm.predict_proba(X_test)
3 preds = probs[:,1]
4 fpr, tpr, threshold = roc_curve(y_test, preds)
5 roc_auc = auc(fpr, tpr)
6
7 f, ax = plt.subplots(figsize=(11,7))
8 plt.title('Logistic Regression ROC Curve', fontsize=16)
9 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
10 plt.legend(loc = 'lower right', fontsize=16)
11 plt.plot([0, 1], [0, 1], 'r--')
12 plt.xlim([0, 1])
13 plt.ylim([0, 1])
14 plt.ylabel('True Positive Rate', fontsize=16)
15 plt.xlabel('False Positive Rate', fontsize=16)
16 plt.tick_params(labelsize=16)
17 plt.show()
18

```



kNN

Raw training set

```

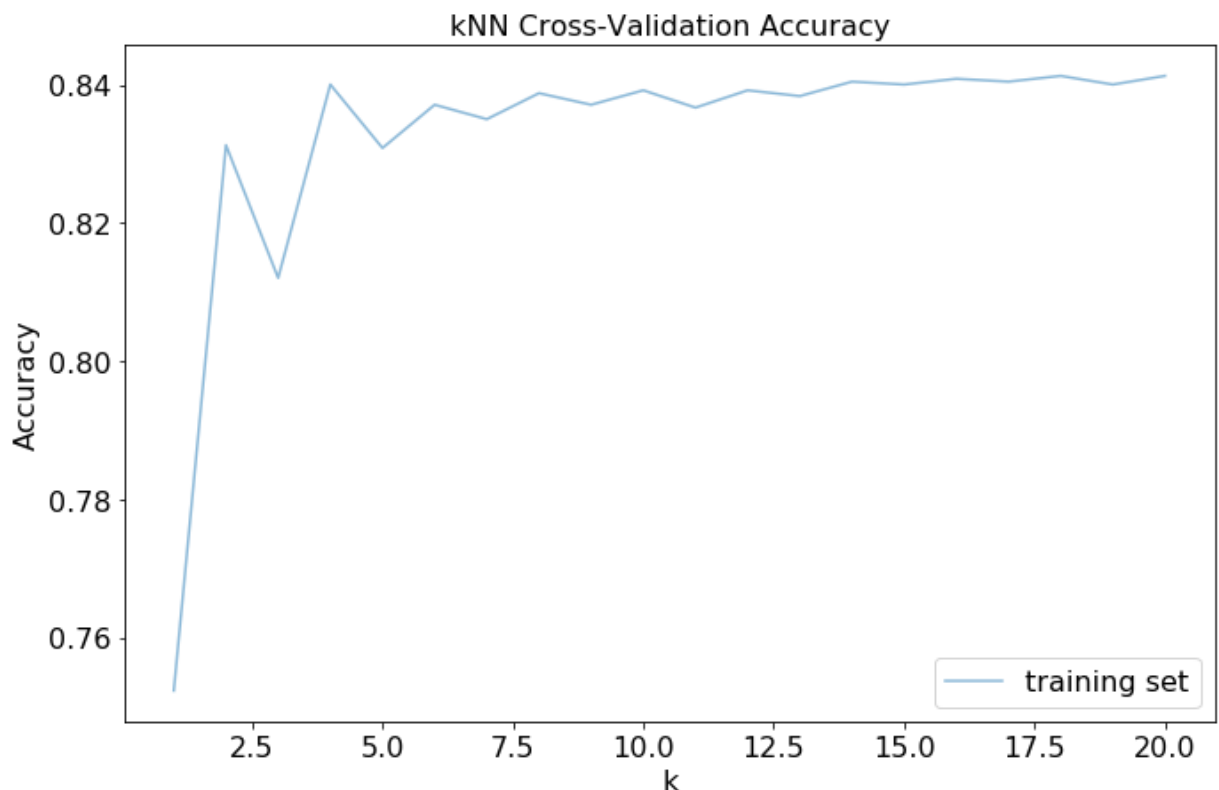
In [16]: 1 max_score = 0
          2 max_k = 0
          3 scores = []
          4
          5 for k in range(1, 21):
          6     knn = KNeighborsClassifier(n_neighbors = k)
          7     score = cross_val_score(knn, X_train, y_train, cv=5).mean()
          8
          9     scores.append(score)
         10     if score > max_score:
         11         max_k = k
         12         max_score = score
         13
         14 scores = pd.DataFrame({'k': range(1, 21), 'accuracy': scores})

```

```

In [17]: 1 fig, ax = plt.subplots(1, 1, figsize=(11, 7))
          2 ax.plot(scores['k'], scores['accuracy'], alpha=0.5, label='training set')
          3 ax.set_title('kNN Cross-Validation Accuracy', fontsize=16)
          4 ax.set_xlabel('k', fontsize=16)
          5 ax.set_ylabel('Accuracy', fontsize=16)
          6 ax.legend(fontsize=16)
          7 ax.tick_params(labelsize=16)
          8 plt.show()

```

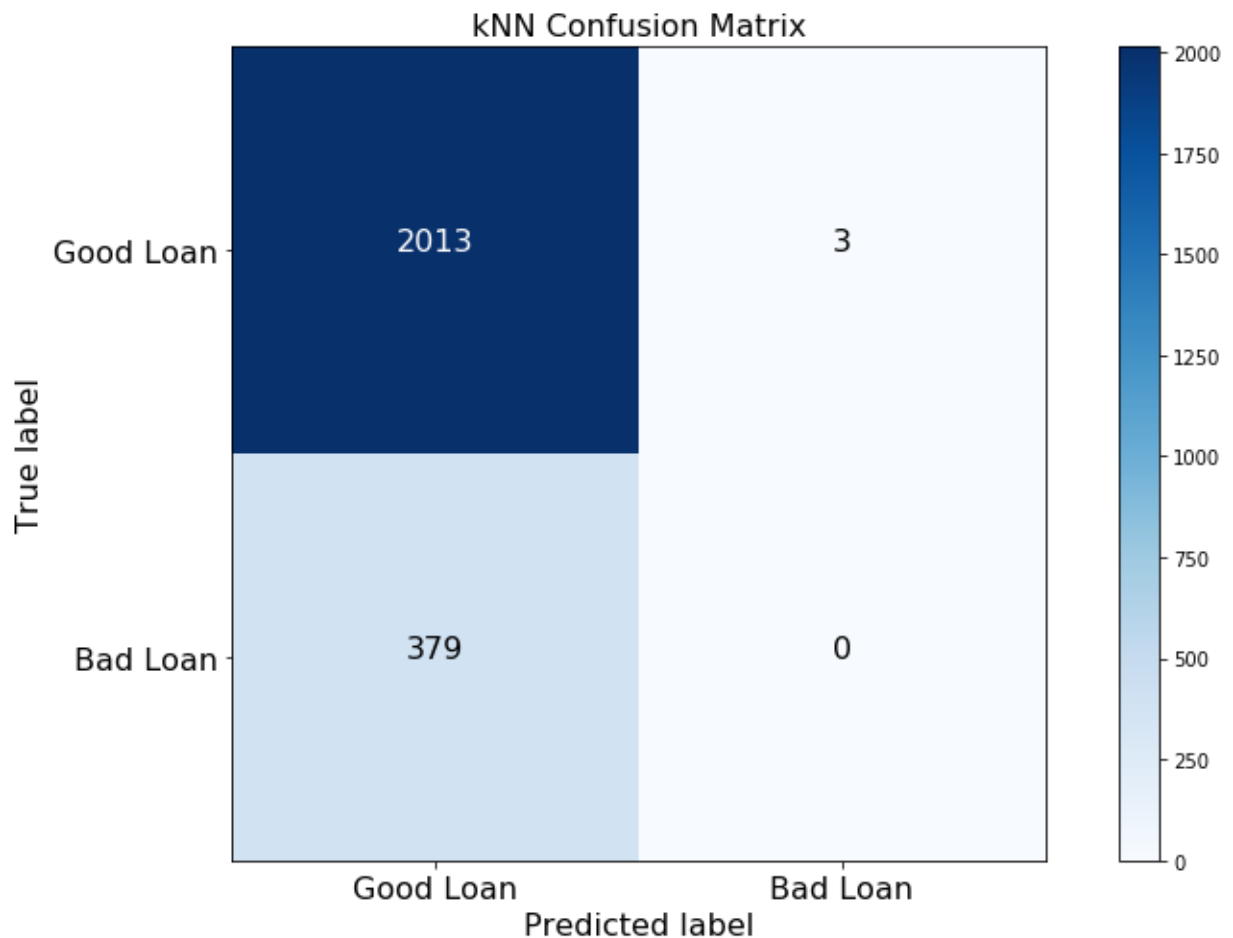


```
In [18]: 1 knn = KNeighborsClassifier(n_neighbors = max_k)
2 knn.fit(X_train, y_train)
3
4 train_acc, test_acc = knn.score(X_train, y_train), knn.score(X_test, y_test)
5
6 print('kNN: Optimal k={}'.format(max_k))
7
8 report_lr = precision_recall_fscore_support(y_test, knn.predict(X_test), aver
9
10 knn_results = {
11     'model': 'kNN',
12     'train_acc': train_acc,
13     'test_acc': test_acc,
14     'precision': report_lr[0],
15     'recall': report_lr[1],
16     'F1': report_lr[2]
17 }
18
19 print('kNN: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={
20     format(knn_results['train_acc'],
21           knn_results['test_acc'],
22           knn_results['precision'],
23           knn_results['recall'],
24           knn_results['F1']))
```

kNN: Optimal k=18

kNN: accuracy on train=84.21%, test=84.05%, precision=0.00, recall=0.00, F1=0.00

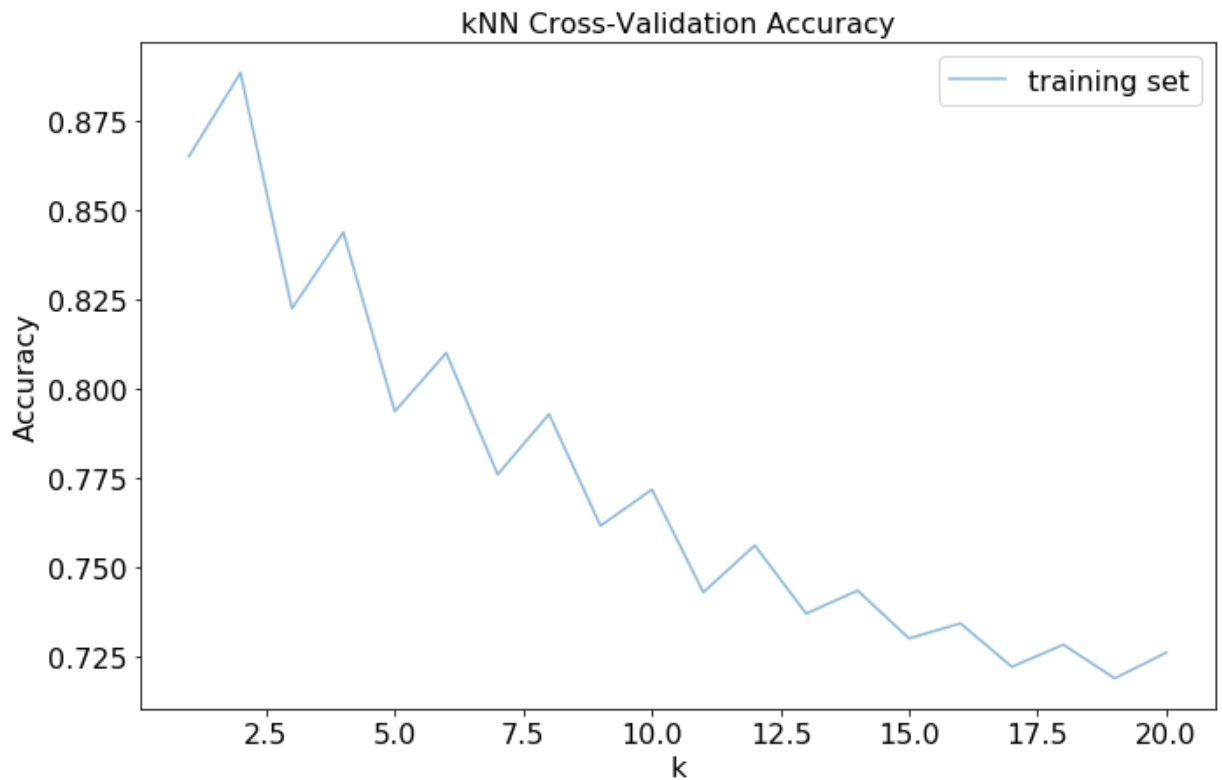
```
In [19]: 1 cnf_matrix = confusion_matrix(y_test, knn.predict(X_test))
2 np.set_printoptions(precision=2)
3
4 plt.figure(figsize=(11,7))
5 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
6                        title='kNN Confusion Matrix')
7
8 plt.show()
```



Oversampled Training set

```
In [20]: 1 max_score = 0
2 max_k = 0
3 scores = []
4
5 for k in range(1, 21):
6     knn = KNeighborsClassifier(n_neighbors = k)
7     score = cross_val_score(knn, X_train_sm, y_train_sm, cv=5).mean()
8
9     scores.append(score)
10    if score > max_score:
11        max_k = k
12        max_score = score
13
14 scores = pd.DataFrame({'k': range(1, 21), 'accuracy': scores})
```

```
In [21]: 1 fig, ax = plt.subplots(1, 1, figsize=(11, 7))
2         ax.plot(scores['k'], scores['accuracy'], alpha=0.5, label='training set')
3         ax.set_title('kNN Cross-Validation Accuracy', fontsize=16)
4         ax.set_xlabel('k', fontsize=16)
5         ax.set_ylabel('Accuracy', fontsize=16)
6         ax.legend(fontsize=16)
7         ax.tick_params(labelsize=16)
8         plt.show()
```



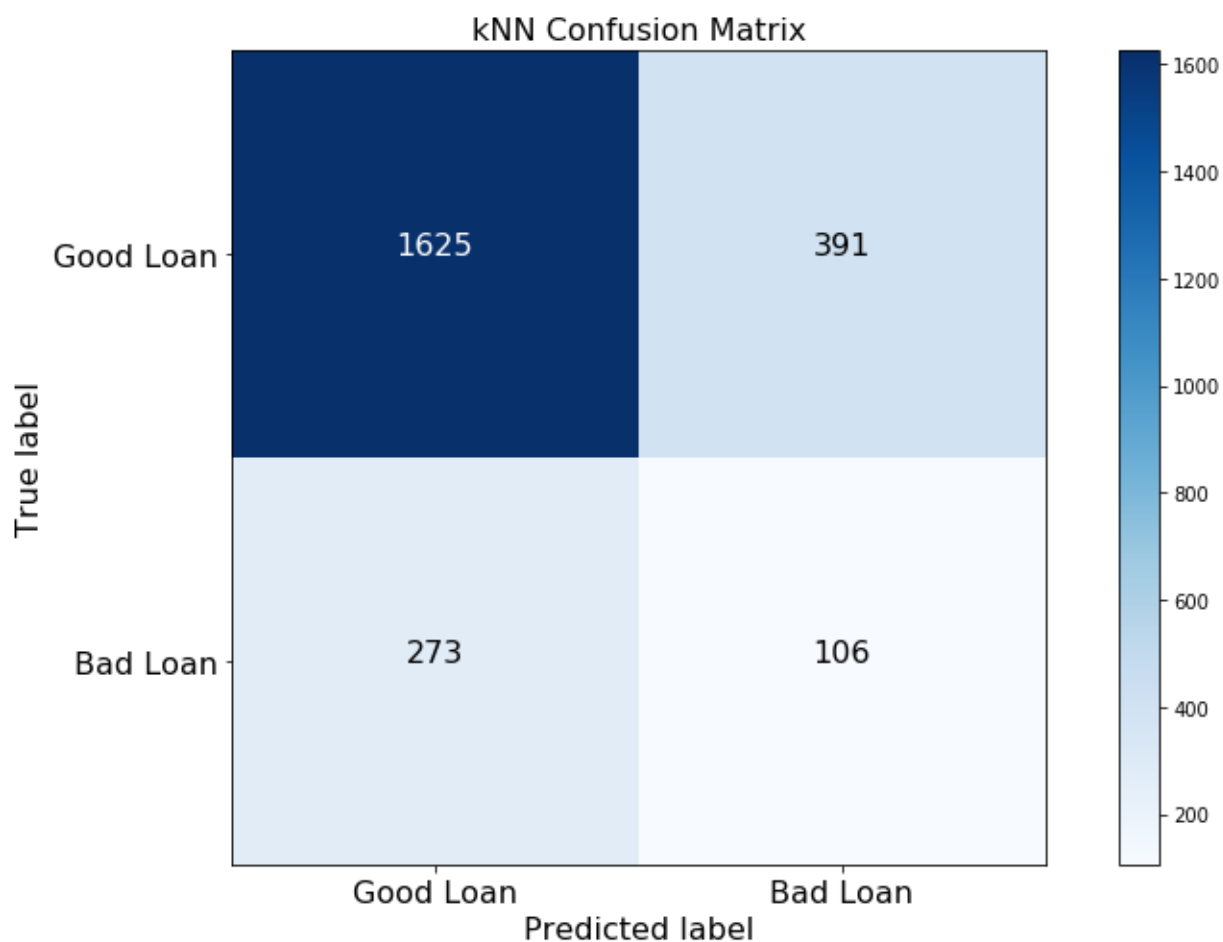
```
In [22]: 1 knn_sm = KNeighborsClassifier(n_neighbors = max_k)
2 knn_sm.fit(X_train_sm, y_train_sm)
3
4 train_acc, test_acc = knn_sm.score(X_train_sm, y_train_sm), knn_sm.score(X_te
5
6 print('kNN: Optimal k={}'.format(max_k))
7
8 report_lr = precision_recall_fscore_support(y_test, knn_sm.predict(X_test), a
9
10 knn_sm_results = {
11     'model': 'kNN',
12     'train_acc': train_acc,
13     'test_acc': test_acc,
14     'precision': report_lr[0],
15     'recall': report_lr[1],
16     'F1': report_lr[2]
17 }
18
19 print('kNN: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={
20     format(knn_sm_results['train_acc'],
21           knn_sm_results['test_acc'],
22           knn_sm_results['precision'],
23           knn_sm_results['recall'],
24           knn_sm_results['F1']))
```

kNN: Optimal k=2

kNN: accuracy on train=99.80%, test=72.28%, precision=0.21, recall=0.28, F1=0.2

4

```
In [23]: 1 cnf_matrix = confusion_matrix(y_test, knn_sm.predict(X_test))
2 np.set_printoptions(precision=2)
3
4 plt.figure(figsize=(11,7))
5 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
6                        title='kNN Confusion Matrix')
7
8 plt.show()
```



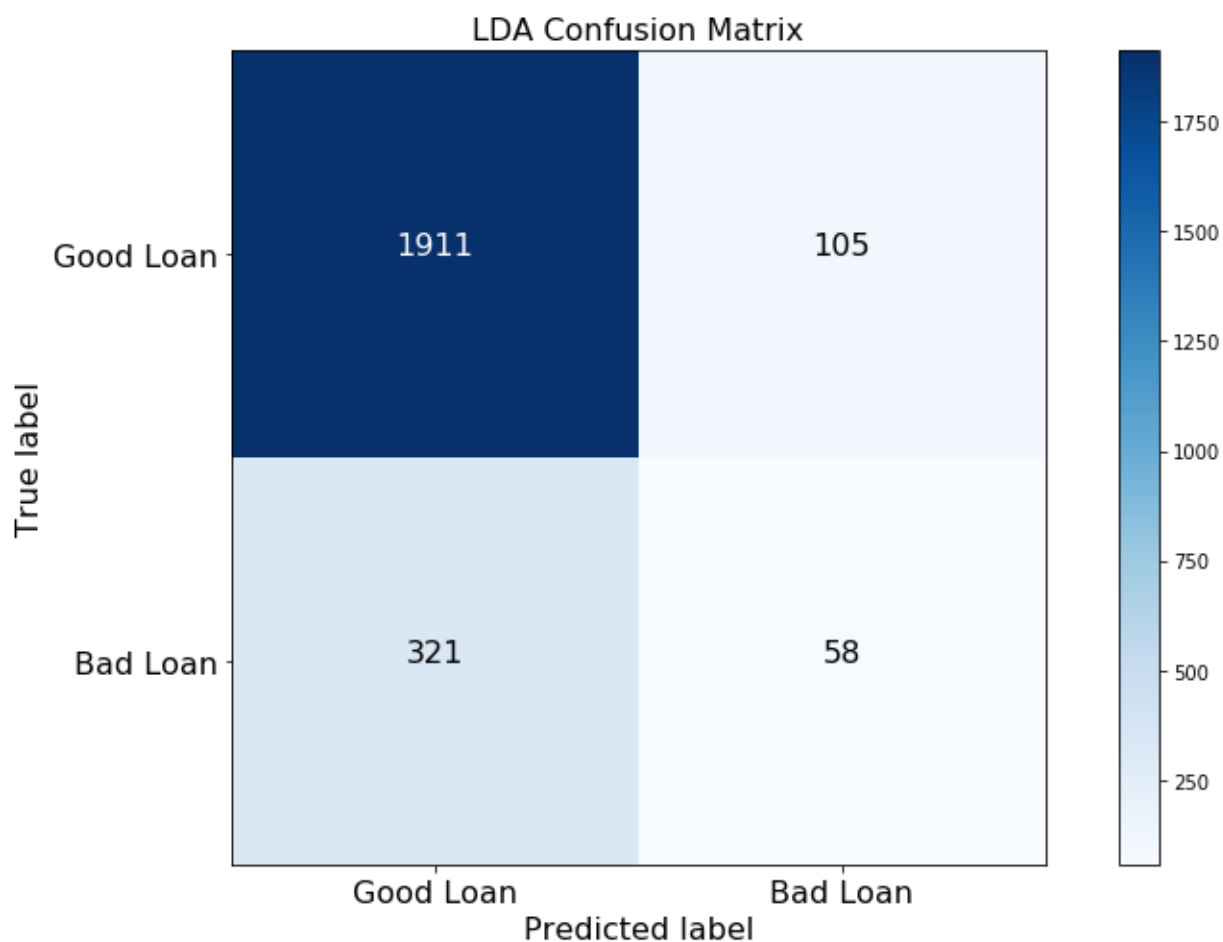
LDA

Raw training set

```
In [24]: 1 lda = LinearDiscriminantAnalysis()
2         lda.fit(X_train, y_train)
3
4         train_acc, test_acc = lda.score(X_train, y_train), lda.score(X_test, y_test)
5         report_lr = precision_recall_fscore_support(y_test, lda.predict(X_test), aver
6
7         lda_results = {
8             'model': 'LDA',
9             'train_acc': train_acc,
10            'test_acc': test_acc,
11            'precision': report_lr[0],
12            'recall': report_lr[1],
13            'F1': report_lr[2]
14        }
15
16         print('LDA: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={
17             format(lda_results['train_acc'],
18                   lda_results['test_acc'],
19                   lda_results['precision'],
20                   lda_results['recall'],
21                   lda_results['F1']))
```

LDA: accuracy on train=85.13%, test=82.21%, precision=0.36, recall=0.15, F1=0.2
1


```
In [25]: 1 cnf_matrix = confusion_matrix(y_test, lda.predict(X_test))
2         np.set_printoptions(precision=2)
3
4         plt.figure(figsize=(11,7))
5         plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize=True,
6                               title='LDA Confusion Matrix')
7
8         plt.show()
```

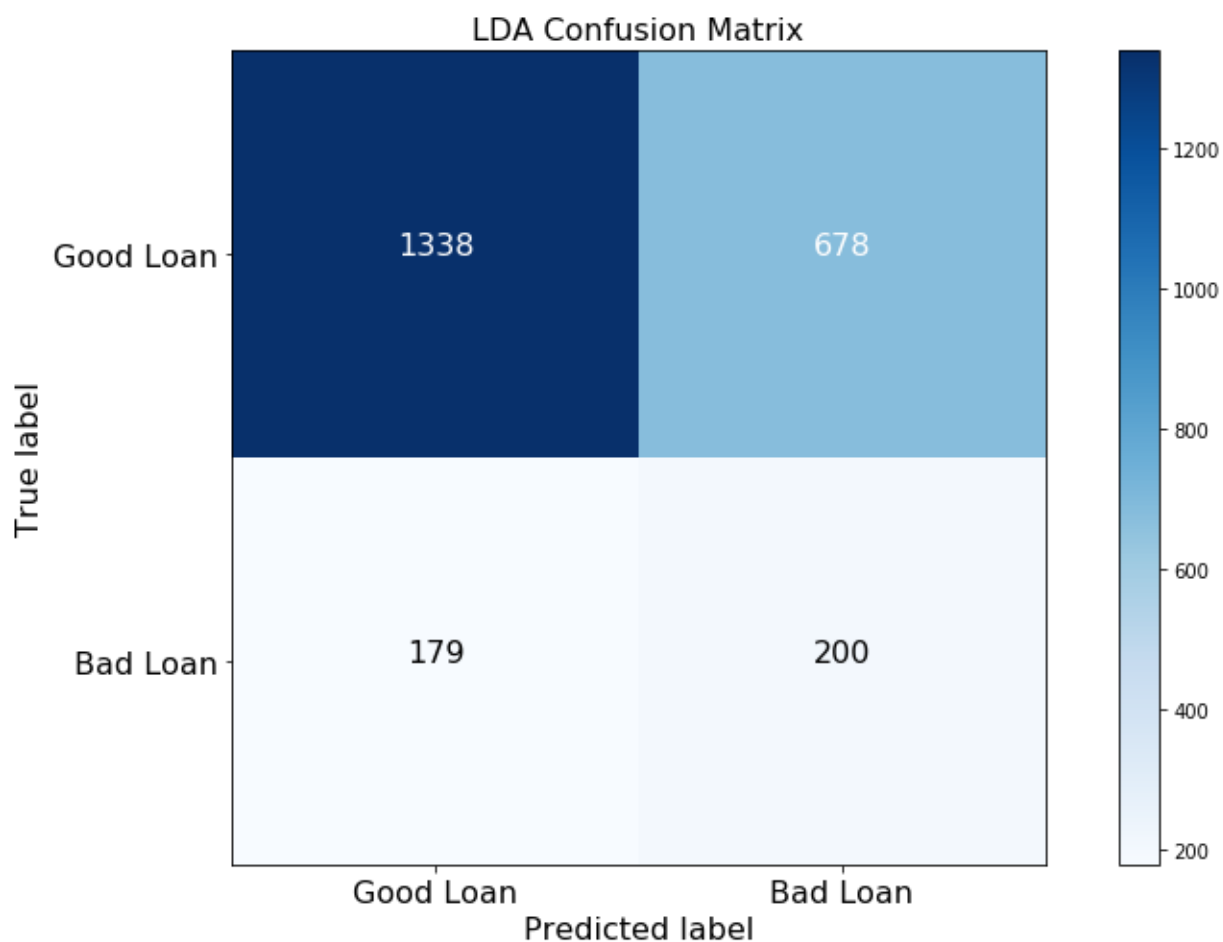


Oversampled Training set

```
In [26]: 1 lda_sm = LinearDiscriminantAnalysis()
2         lda_sm.fit(X_train_sm, y_train_sm)
3
4         train_acc, test_acc = lda_sm.score(X_train_sm, y_train_sm), lda_sm.score(X_test_sm, y_test_sm)
5         report_lr = precision_recall_fscore_support(y_test, lda_sm.predict(X_test), average='weighted')
6
7         lda_sm_results = {
8             'model': 'LDA',
9             'train_acc': train_acc,
10            'test_acc': test_acc,
11            'precision': report_lr[0],
12            'recall': report_lr[1],
13            'F1': report_lr[2]
14        }
15
16         print('LDA: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={:.2f}, F1={:.2f}'.format(
17             train_acc, test_acc, report_lr[0], report_lr[1], report_lr[2]))
18
19         # Print the confusion matrix
20         print('Confusion Matrix:')
21         print(confusion_matrix(y_test, lda_sm.predict(X_test)))
```

LDA: accuracy on train=72.61%, test=64.22%, precision=0.23, recall=0.53, F1=0.32

```
In [27]: 1 cnf_matrix = confusion_matrix(y_test, lda_sm.predict(X_test))
2         np.set_printoptions(precision=2)
3
4         plt.figure(figsize=(11,7))
5         plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize=True,
6                               title='LDA Confusion Matrix')
7
8         plt.show()
```



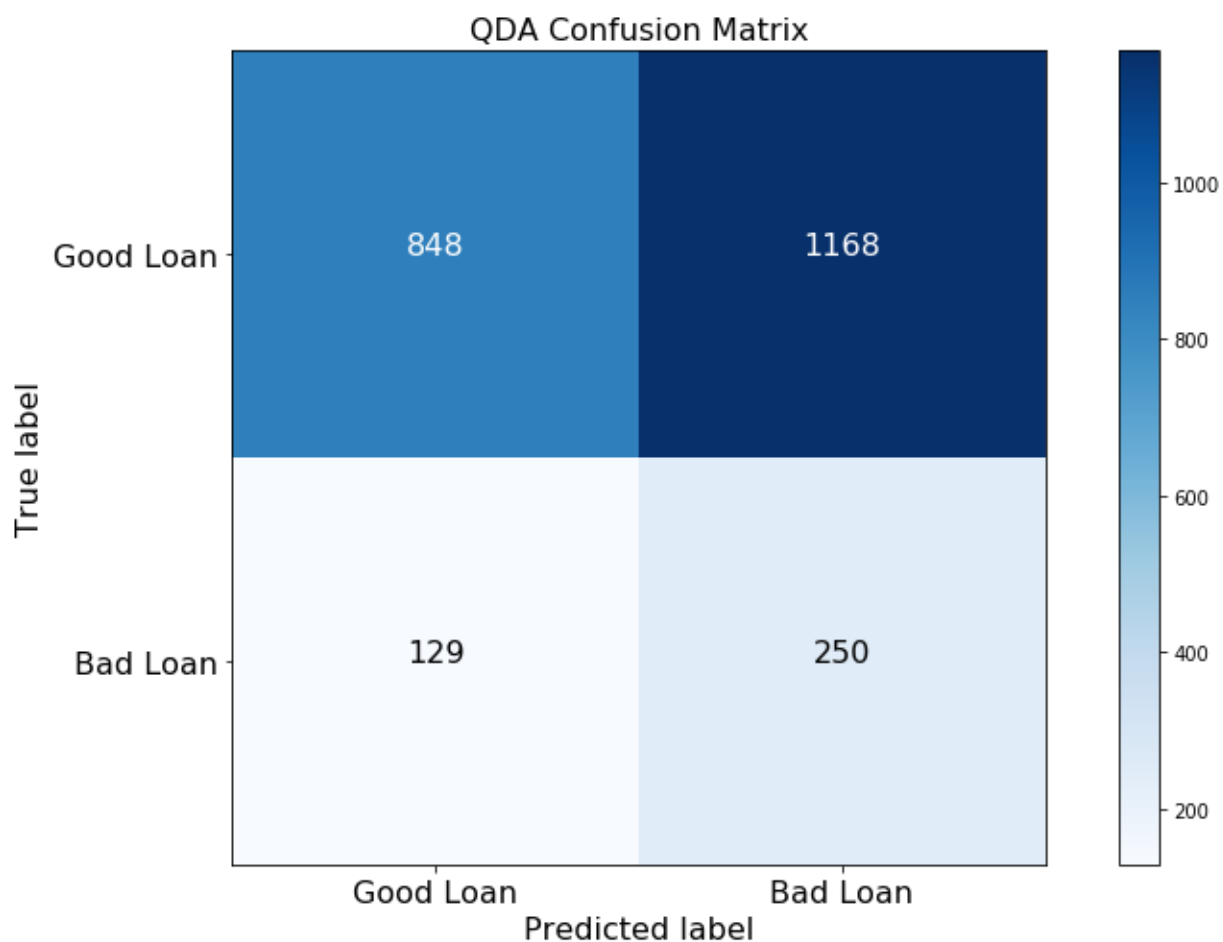
QDA

Raw training set

```
In [28]: 1 qda = QuadraticDiscriminantAnalysis()
2 qda.fit(X_train, y_train)
3
4 train_acc, test_acc = qda.score(X_train, y_train), qda.score(X_test, y_test)
5 report_lr = precision_recall_fscore_support(y_test, qda.predict(X_test), aver
6
7 qda_results = {
8     'model': 'QDA',
9     'train_acc': train_acc,
10    'test_acc': test_acc,
11    'precision': report_lr[0],
12    'recall': report_lr[1],
13    'F1': report_lr[2]
14 }
15
16 print('LDA: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={
17       format(qda_results['train_acc'],
18             qda_results['test_acc'],
19             qda_results['precision'],
20             qda_results['recall'],
21             qda_results['F1']))
```

LDA: accuracy on train=53.63%, test=45.85%, precision=0.18, recall=0.66, F1=0.2
8

```
In [29]: 1 cnf_matrix = confusion_matrix(y_test, qda.predict(X_test))
2         np.set_printoptions(precision=2)
3
4         plt.figure(figsize=(11,7))
5         plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize=True,
6                               title='QDA Confusion Matrix')
7
8         plt.show()
```



Oversampled Training set

```

In [30]: 1 qda_sm = QuadraticDiscriminantAnalysis()
2 qda_sm.fit(X_train_sm, y_train_sm)
3
4 train_acc, test_acc = qda_sm.score(X_train_sm, y_train_sm), qda_sm.score(X_te
5 report_lr = precision_recall_fscore_support(y_test, qda_sm.predict(X_test), a
6
7 qda_sm_results = {
8     'model': 'QDA',
9     'train_acc': train_acc,
10    'test_acc': test_acc,
11    'precision': report_lr[0],
12    'recall': report_lr[1],
13    'F1': report_lr[2]
14 }
15
16 print('LDA: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={
17       format(qda_sm_results['train_acc'],
18             qda_sm_results['test_acc'],
19             qda_sm_results['precision'],
20             qda_sm_results['recall'],
21             qda_sm_results['F1']))

```

LDA: accuracy on train=76.72%, test=53.32%, precision=0.17, recall=0.51, F1=0.26

Single Decision Tree

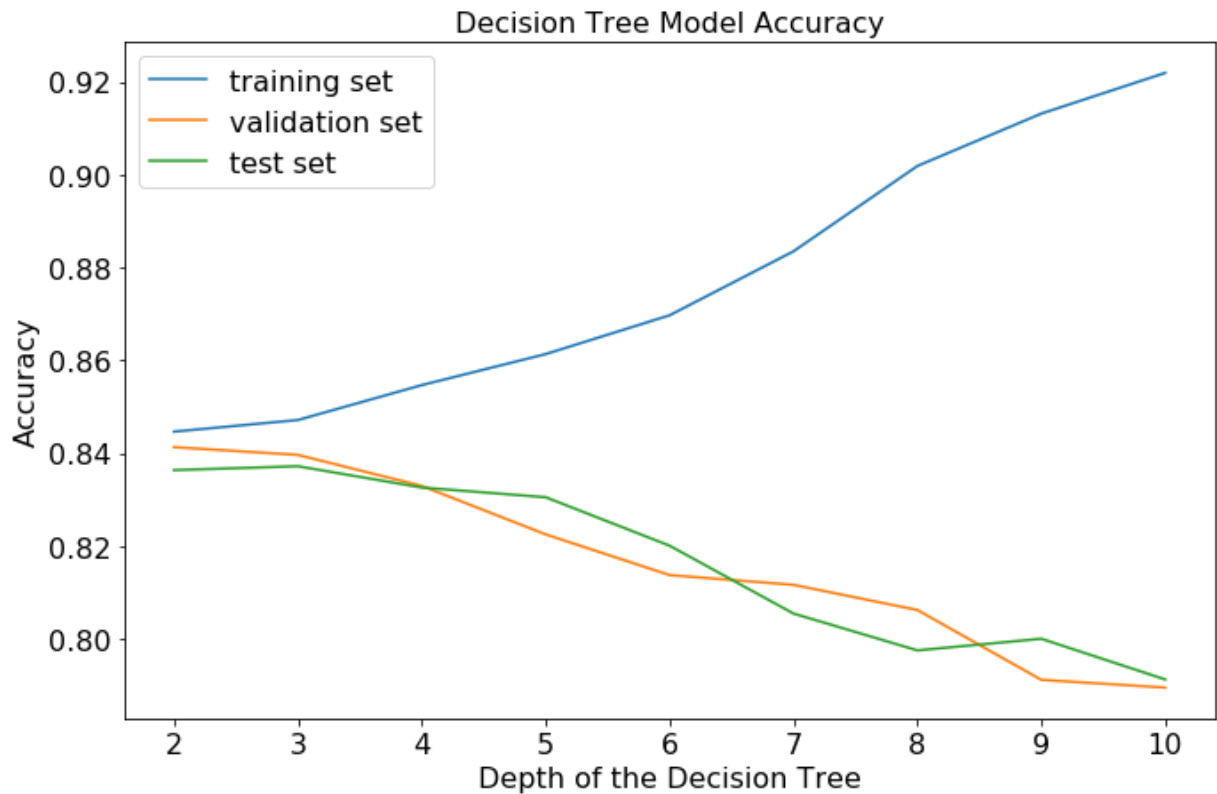
Raw training set

```

In [31]: 1 train_scores = []
2 validation_scores = []
3 test_scores = []
4
5 best_score = 0
6 best_depth = 0
7
8 depths = [i for i in range(2, 11)]
9
10 for depth in depths:
11     tree = DecisionTreeClassifier(max_depth = depth)
12     tree.fit(X_train, y_train)
13
14     train_scores.append(tree.score(X_train, y_train))
15     test_scores.append(tree.score(X_test, y_test))
16
17     val_score = cross_val_score(estimator=tree, X=X_train, y=y_train, cv=5).m
18     validation_scores.append(val_score)
19
20     if val_score > best_score:
21         best_depth = depth
22         best_score = score

```

```
In [32]: 1 fig, ax = plt.subplots(1, 1, figsize=(11, 7))
2 ax.plot(depths, train_scores, label='training set')
3 ax.plot(depths, validation_scores, label='validation set')
4 ax.plot(depths, test_scores, label='test set')
5 ax.tick_params(labelsize=16)
6 ax.set_title('Decision Tree Model Accuracy', fontsize=16)
7 ax.set_xlabel('Depth of the Decision Tree', fontsize=16)
8 ax.set_ylabel('Accuracy', fontsize=16)
9 ax.legend(fontsize=16)
10 plt.show()
```

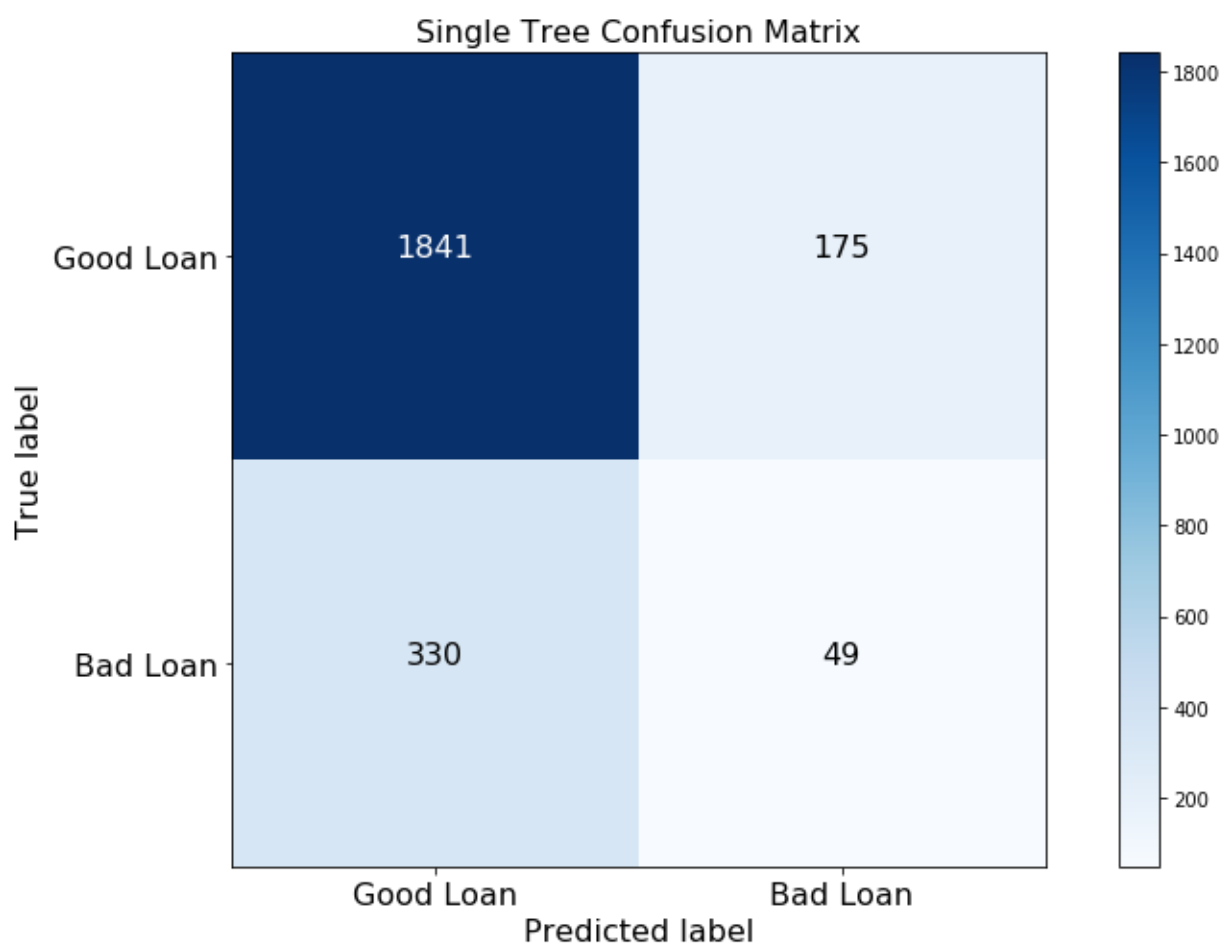


```
In [33]: 1 tree = DecisionTreeClassifier(max_depth = best_depth)
2 tree.fit(X_train, y_train)
3
4 train_acc, test_acc = tree.score(X_train, y_train), tree.score(X_test, y_test)
5
6 print('Single Tree: Optimal depth={}'.format(best_depth))
7
8 report_lr = precision_recall_fscore_support(y_test, tree.predict(X_test), average='micro')
9
10 tree_results = {
11     'model': 'Single Tree',
12     'train_acc': train_acc,
13     'test_acc': test_acc,
14     'precision': report_lr[0],
15     'recall': report_lr[1],
16     'F1': report_lr[2]
17 }
18
19 print('Single Tree: accuracy on train={:.2%}, test={:.2%}, precision={:.2f},
20       format(tree_results['train_acc'],
21               tree_results['test_acc'],
22               tree_results['precision'],
23               tree_results['recall'],
24               tree_results['F1']))
```

Single Tree: Optimal depth=10

Single Tree: accuracy on train=92.19%, test=78.91%, precision=0.22, recall=0.13, F1=0.16

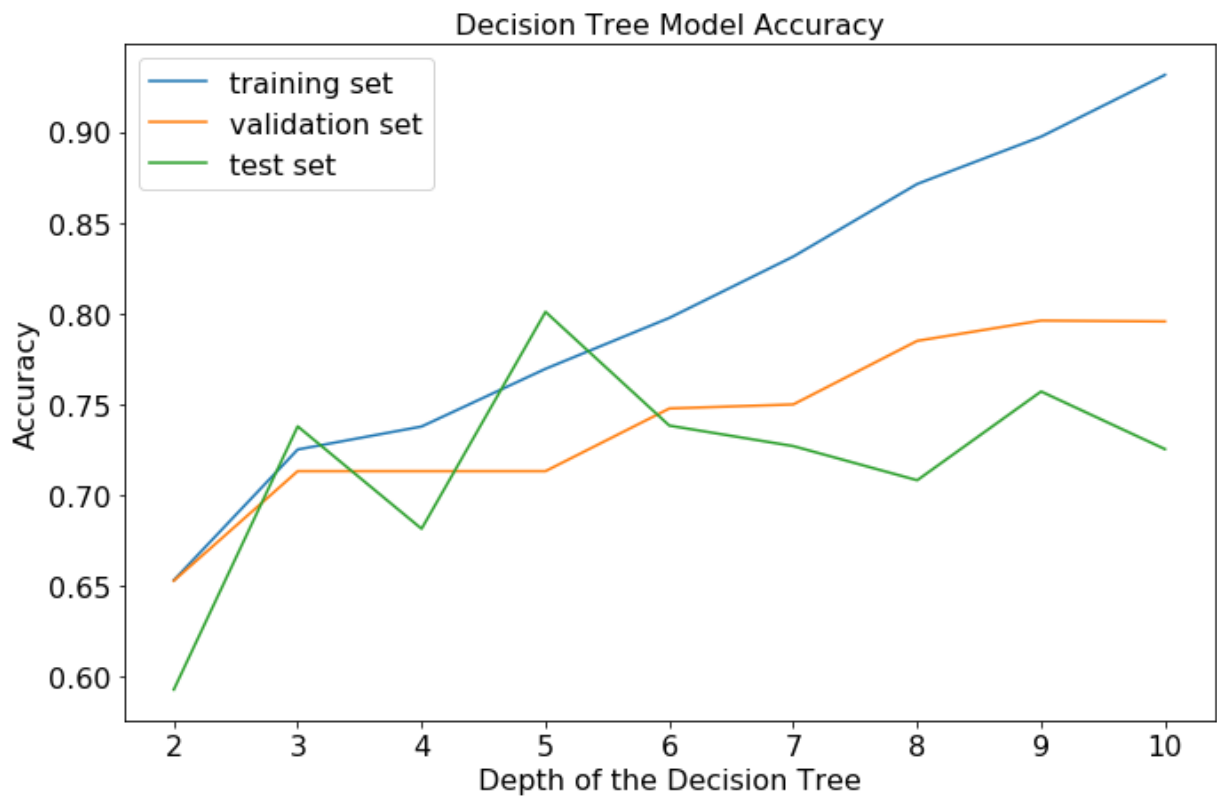

```
In [34]: 1 # Plot confusion matrix
2 y_pred = tree.predict(X_test)
3
4 cnf_matrix = confusion_matrix(y_test, y_pred)
5 np.set_printoptions(precision=2)
6
7 plt.figure(figsize=(11,7))
8 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
9                       title='Single Tree Confusion Matrix')
10
11 plt.show()
```



Oversampled Training set

```
In [35]: 1 train_scores = []
2         validation_scores = []
3         test_scores = []
4
5         best_score = 0
6         best_depth = 0
7
8         depths = [i for i in range(2, 11)]
9
10        for depth in depths:
11            tree = DecisionTreeClassifier(max_depth = depth)
12            tree.fit(X_train_sm, y_train_sm)
13
14            train_scores.append(tree.score(X_train_sm, y_train_sm))
15            test_scores.append(tree.score(X_test, y_test))
16
17            val_score = cross_val_score(estimator=tree, X=X_train_sm, y=y_train_sm, c
18            validation_scores.append(val_score)
19
20            if val_score > best_score:
21                best_depth = depth
22                best_score = score
```

```
In [36]: 1 fig, ax = plt.subplots(1, 1, figsize=(11, 7))
2 ax.plot(depths, train_scores, label='training set')
3 ax.plot(depths, validation_scores, label='validation set')
4 ax.plot(depths, test_scores, label='test set')
5 ax.tick_params(labelsize=16)
6 ax.set_title('Decision Tree Model Accuracy', fontsize=16)
7 ax.set_xlabel('Depth of the Decision Tree', fontsize=16)
8 ax.set_ylabel('Accuracy', fontsize=16)
9 ax.legend(fontsize=16)
10 plt.show()
```

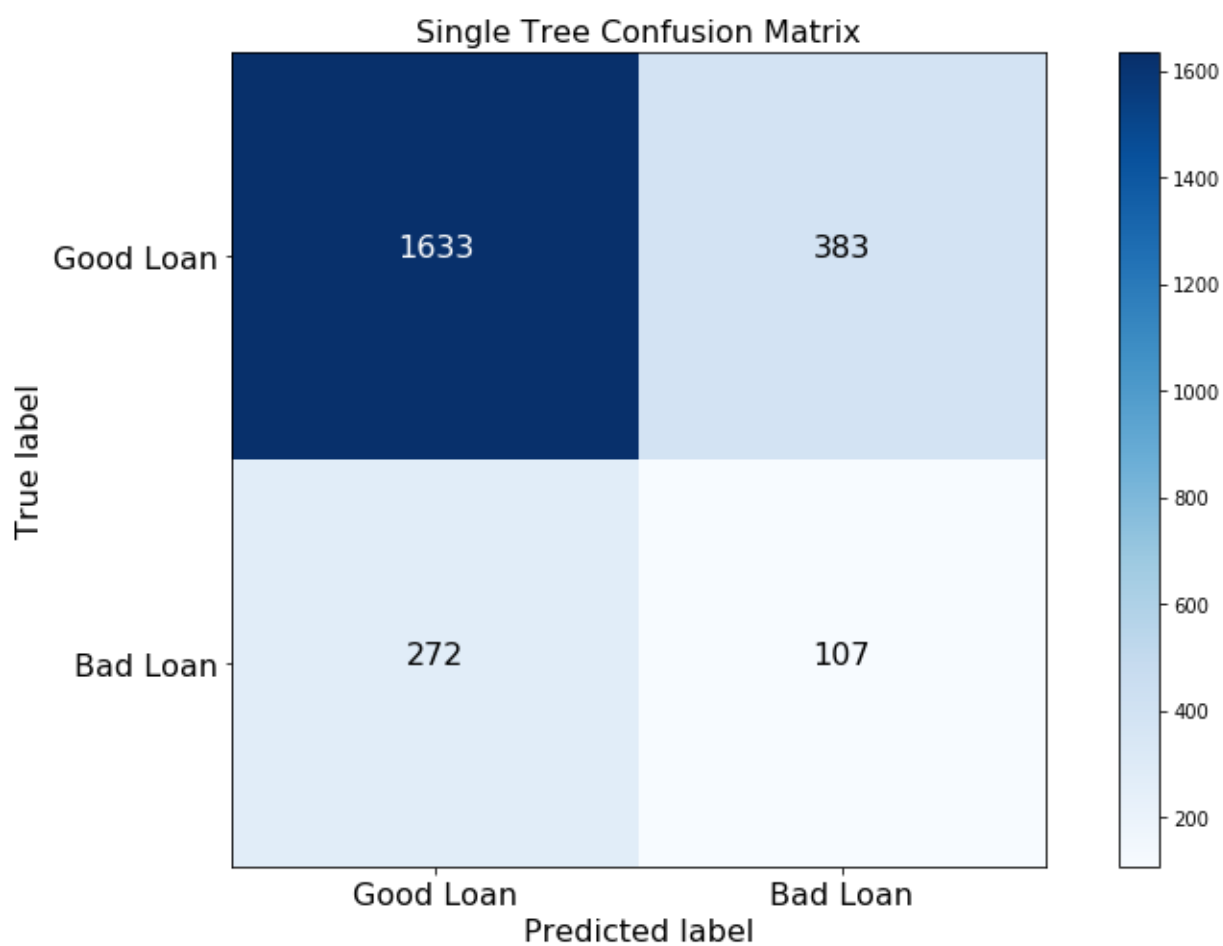


```
In [37]: 1 tree_sm = DecisionTreeClassifier(max_depth = best_depth)
2 tree_sm.fit(X_train_sm, y_train_sm)
3
4 train_acc, test_acc = tree_sm.score(X_train_sm, y_train_sm), tree_sm.score(X_
5
6 print('Single Tree: Optimal depth={}'.format(best_depth))
7
8 report_lr = precision_recall_fscore_support(y_test, tree_sm.predict(X_test),
9
10 tree_sm_results = {
11     'model': 'Single Tree',
12     'train_acc': train_acc,
13     'test_acc': test_acc,
14     'precision': report_lr[0],
15     'recall': report_lr[1],
16     'F1': report_lr[2]
17 }
18
19 print('Single Tree: accuracy on train={:.2%}, test={:.2%}, precision={:.2f},
20       format(tree_sm_results['train_acc'],
21             tree_sm_results['test_acc'],
22             tree_sm_results['precision'],
23             tree_sm_results['recall'],
24             tree_sm_results['F1']))
```

Single Tree: Optimal depth=10

Single Tree: accuracy on train=93.08%, test=72.65%, precision=0.22, recall=0.28, F1=0.25

```
In [38]: 1 # Plot confusion matrix
2 y_pred = tree_sm.predict(X_test)
3
4 cnf_matrix = confusion_matrix(y_test, y_pred)
5 np.set_printoptions(precision=2)
6
7 plt.figure(figsize=(11,7))
8 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
9                       title='Single Tree Confusion Matrix')
10
11 plt.show()
```

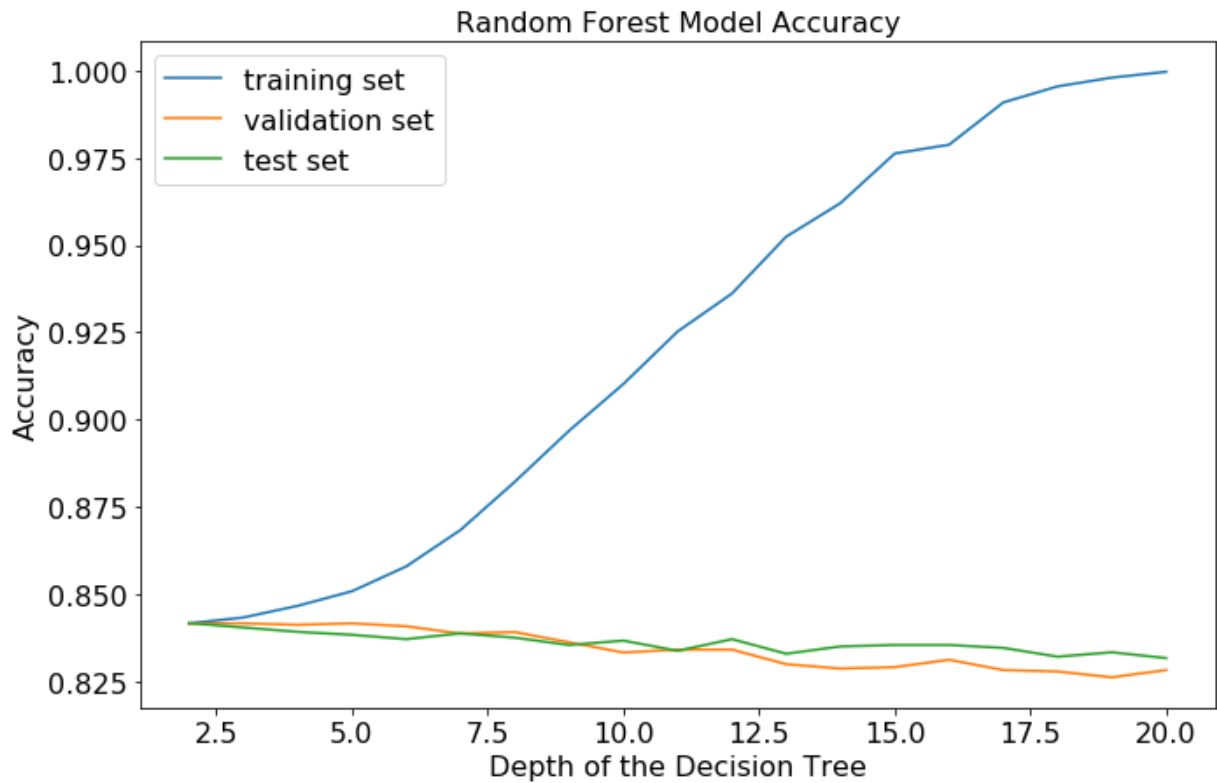


Random Forest

Raw training set

```
In [39]: 1 train_scores = []
2         validation_scores = []
3         test_scores = []
4
5         best_score = 0
6         best_depth = 0
7
8         n_trees = 100
9         depths = [i for i in range(2, 21)]
10
11        for depth in depths:
12            rf = RandomForestClassifier(n_estimators=n_trees, max_depth=depth, n_jobs=
13            rf.fit(X_train, y_train)
14
15            train_scores.append(rf.score(X_train, y_train))
16            test_scores.append(rf.score(X_test, y_test))
17
18            val_score = cross_val_score(estimator=rf, X=X_train, y=y_train, cv=5).mean()
19            validation_scores.append(val_score)
20
21            if val_score > best_score:
22                best_depth = depth
23                best_score = val_score
```

```
In [40]: 1 fig, ax = plt.subplots(1, 1, figsize=(11, 7))
2         ax.plot(depths, train_scores, label='training set')
3         ax.plot(depths, validation_scores, label='validation set')
4         ax.plot(depths, test_scores, label='test set')
5         ax.tick_params(labelsize=16)
6         ax.set_title('Random Forest Model Accuracy', fontsize=16)
7         ax.set_xlabel('Depth of the Decision Tree', fontsize=16)
8         ax.set_ylabel('Accuracy', fontsize=16)
9         ax.legend(fontsize=16)
10        plt.show()
```

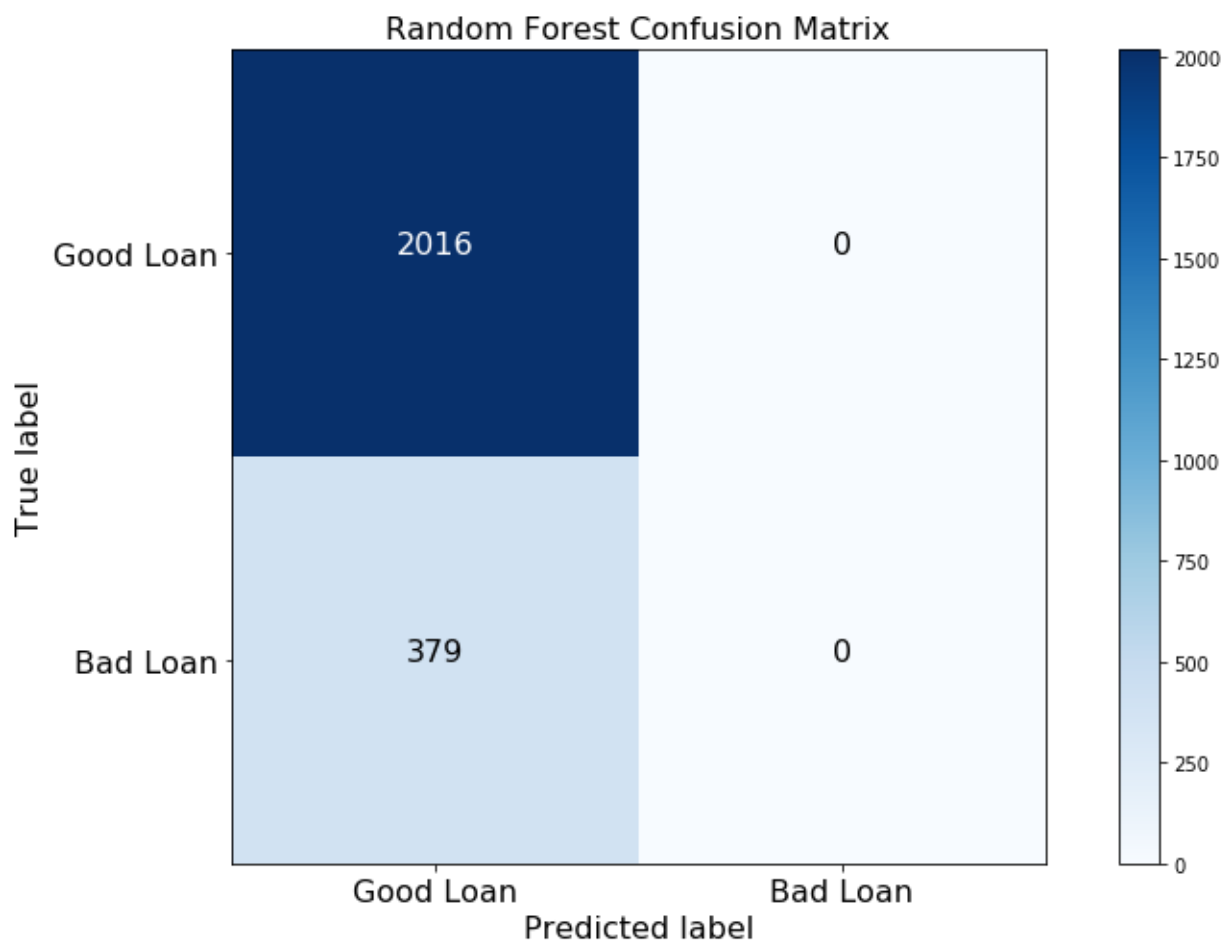


```
In [41]: 1 rf = RandomForestClassifier(n_estimators=n_trees, max_depth=best_depth, n_job
2 rf.fit(X_train, y_train)
3
4 train_acc, test_acc = rf.score(X_train, y_train), rf.score(X_test, y_test)
5
6 print('Random Forest: Optimal depth={}'.format(best_depth))
7
8 report_lr = precision_recall_fscore_support(y_test, rf.predict(X_test), avera
9
10 rf_results = {
11     'model': 'Random Forest',
12     'train_acc': train_acc,
13     'test_acc': test_acc,
14     'precision': report_lr[0],
15     'recall': report_lr[1],
16     'F1': report_lr[2]
17 }
18
19 print('Random Forest: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}
20       format(rf_results['train_acc'],
21             rf_results['test_acc'],
22             rf_results['precision'],
23             rf_results['recall'],
24             rf_results['F1']))
```

Random Forest: Optimal depth=2

Random Forest: accuracy on train=84.17%, test=84.18%, precision=0.00, recall=0.00, F1=0.00


```
In [42]: 1 # Plot confusion matrix
2 y_pred = rf.predict(X_test)
3
4 cnf_matrix = confusion_matrix(y_test, y_pred)
5 np.set_printoptions(precision=2)
6
7 plt.figure(figsize=(11,7))
8 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
9                       title='Random Forest Confusion Matrix')
10
11 plt.show()
```

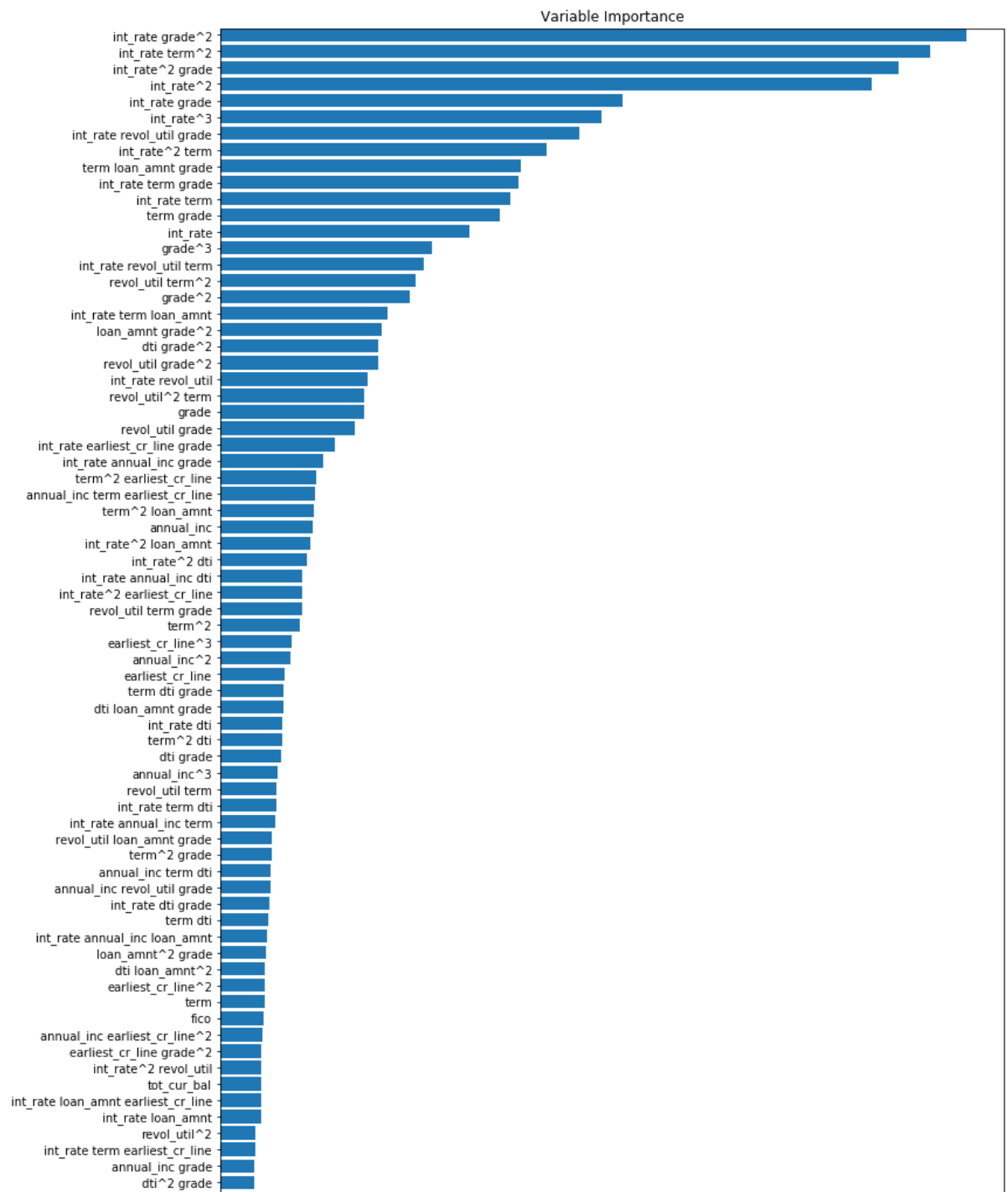


In [43]:

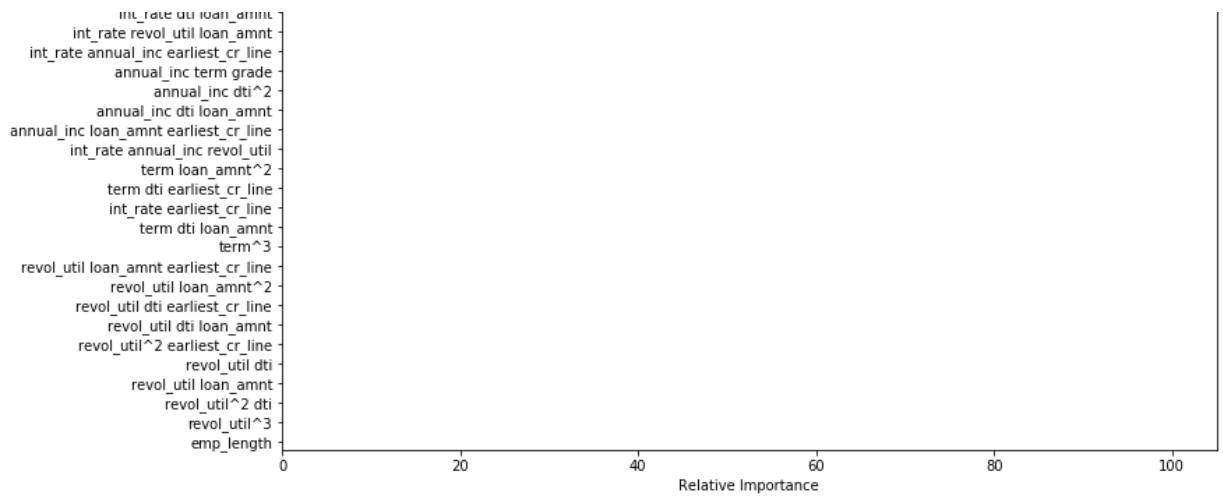
```

1 # Random Forest Feature Importance
2 feature_importance = rf.feature_importances_
3 feature_importance = 100.0 * (feature_importance / feature_importance.max())
4 sorted_idx = np.argsort(feature_importance)
5 pos = np.arange(sorted_idx.shape[0])
6
7 # Plot
8 plt.figure(figsize=(12,50))
9 plt.barh(pos, feature_importance[sorted_idx], align='center')
10 plt.yticks(pos, X_train.columns[sorted_idx])
11 plt.xlabel('Relative Importance')
12 plt.title('Variable Importance')
13 plt.margins(y=0)
14 plt.show()

```



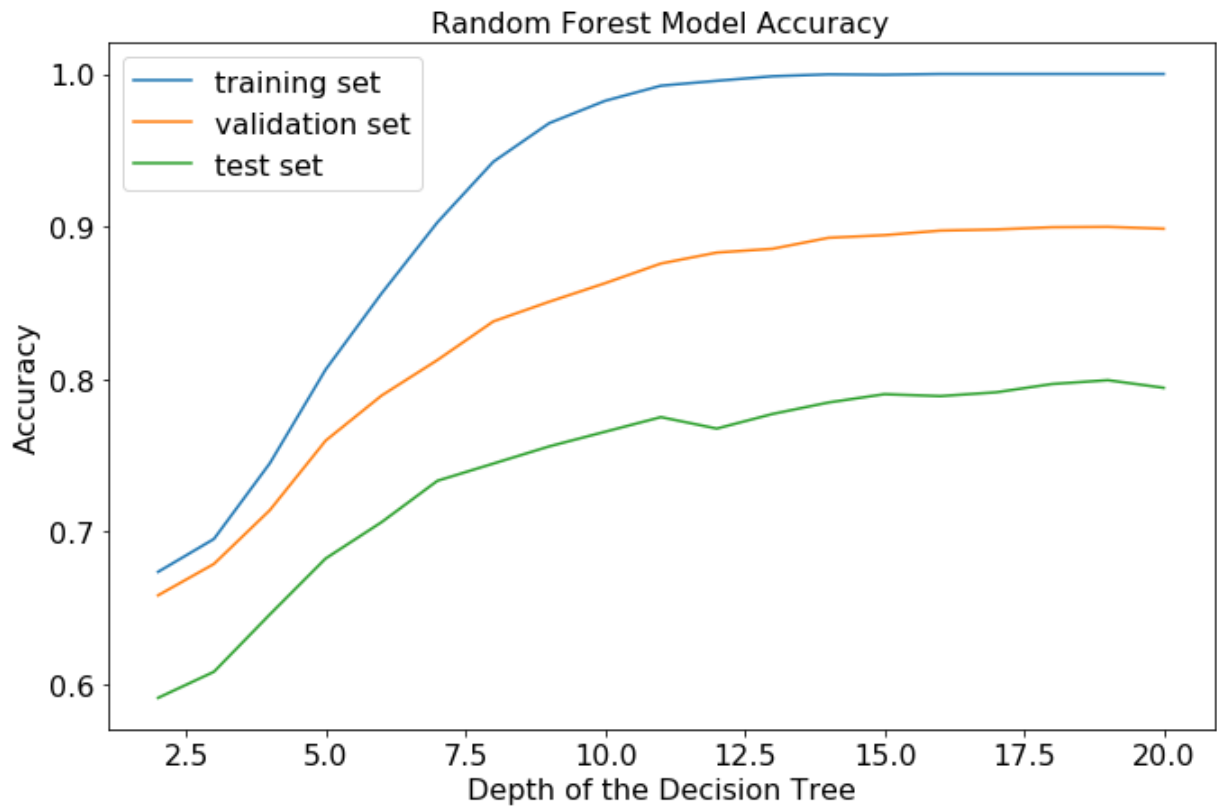
annual_inc earliest_cr_line
 term dti^2
 int_rate revol_util^2
 revol_util term dti
 annual_inc^2 dti
 annual_inc^2 term
 revol_util earliest_cr_line grade
 revol_util term loan_amnt
 revol_util dti^2
 annual_inc loan_amnt^2
 annual_inc^2 earliest_cr_line
 annual_inc earliest_cr_line grade
 int_rate revol_util earliest_cr_line
 term earliest_cr_line^2
 term earliest_cr_line grade
 annual_inc grade^2
 revol_util earliest_cr_line
 annual_inc loan_amnt grade
 annual_inc revol_util loan_amnt
 annual_inc^2 loan_amnt
 revol_util^2 grade
 annual_inc^2 revol_util
 annual_inc dti grade
 annual_inc revol_util
 term loan_amnt
 dti earliest_cr_line^2
 int_rate revol_util dti
 dti earliest_cr_line
 annual_inc^2 grade
 revol_util earliest_cr_line^2
 revol_util term earliest_cr_line
 annual_inc revol_util^2
 earliest_cr_line^2 grade
 int_rate dti^2
 term earliest_cr_line
 annual_inc dti
 int_rate earliest_cr_line^2
 annual_inc dti earliest_cr_line
 annual_inc term^2
 dti
 int_rate loan_amnt^2
 fico_rmg
 annual_inc term
 purpose_house
 loan_amnt earliest_cr_line^2
 revol_util dti grade
 revol_util^2 loan_amnt
 annual_inc loan_amnt
 dti loan_amnt earliest_cr_line
 purpose_moving
 purpose_other
 purpose_renewable_energy
 purpose_small_business
 purpose_vacation
 loan_amnt
 purpose_wedding
 purpose_major_purchase
 dti^2 earliest_cr_line
 loan_amnt^2 earliest_cr_line
 loan_amnt^3
 revol_util
 dti earliest_cr_line grade
 purpose_medical
 int_rate annual_inc^2
 dti^3
 purpose_home_improvement
 delinq_2yrs
 inq_last_6mths
 open_acc
 pub_rec
 application_type
 acc_now_delinq
 tot_coll_amt
 loan_amnt earliest_cr_line grade
 home_ownership_NONE
 home_ownership_OTHER
 home_ownership_OWN
 home_ownership_RENT
 ver_Source_Verified
 ver_Verified
 purpose_credit_card
 purpose_debt_consolidation
 purpose_educational
 dti^2 loan_amnt
 term loan_amnt earliest_cr_line
 int_rate annual_inc
 annual_inc revol_util dti
 dti^2
 dti loan_amnt
 annual_inc term loan_amnt
 annual_inc revol_util earliest_cr_line
 loan_amnt^2
 loan_amnt earliest_cr_line
 loan_amnt grade
 earliest_cr_line grade
 term grade^2
 annual_inc revol_util term
 int_rate loan_amnt grade
 int_rate^2 annual_inc
 int_rate dti earliest_cr_line
 int_rate dti loan_amnt



Oversampled Training set

```
In [44]: 1 train_scores = []
2 validation_scores = []
3 test_scores = []
4
5 best_score = 0
6 best_depth = 0
7
8 n_trees = 100
9 depths = [i for i in range(2, 21)]
10
11 for depth in depths:
12     rf = RandomForestClassifier(n_estimators=n_trees, max_depth=depth, n_jobs=
13     rf.fit(X_train_sm, y_train_sm)
14
15     train_scores.append(rf.score(X_train_sm, y_train_sm))
16     test_scores.append(rf.score(X_test, y_test))
17
18     val_score = cross_val_score(estimator=rf, X=X_train_sm, y=y_train_sm, cv=
19     validation_scores.append(val_score)
20
21     if val_score > best_score:
22         best_depth = depth
23         best_score = val_score
```

```
In [45]: 1 fig, ax = plt.subplots(1, 1, figsize=(11, 7))
2         ax.plot(depths, train_scores, label='training set')
3         ax.plot(depths, validation_scores, label='validation set')
4         ax.plot(depths, test_scores, label='test set')
5         ax.tick_params(labelsize=16)
6         ax.set_title('Random Forest Model Accuracy', fontsize=16)
7         ax.set_xlabel('Depth of the Decision Tree', fontsize=16)
8         ax.set_ylabel('Accuracy', fontsize=16)
9         ax.legend(fontsize=16)
10        plt.show()
```

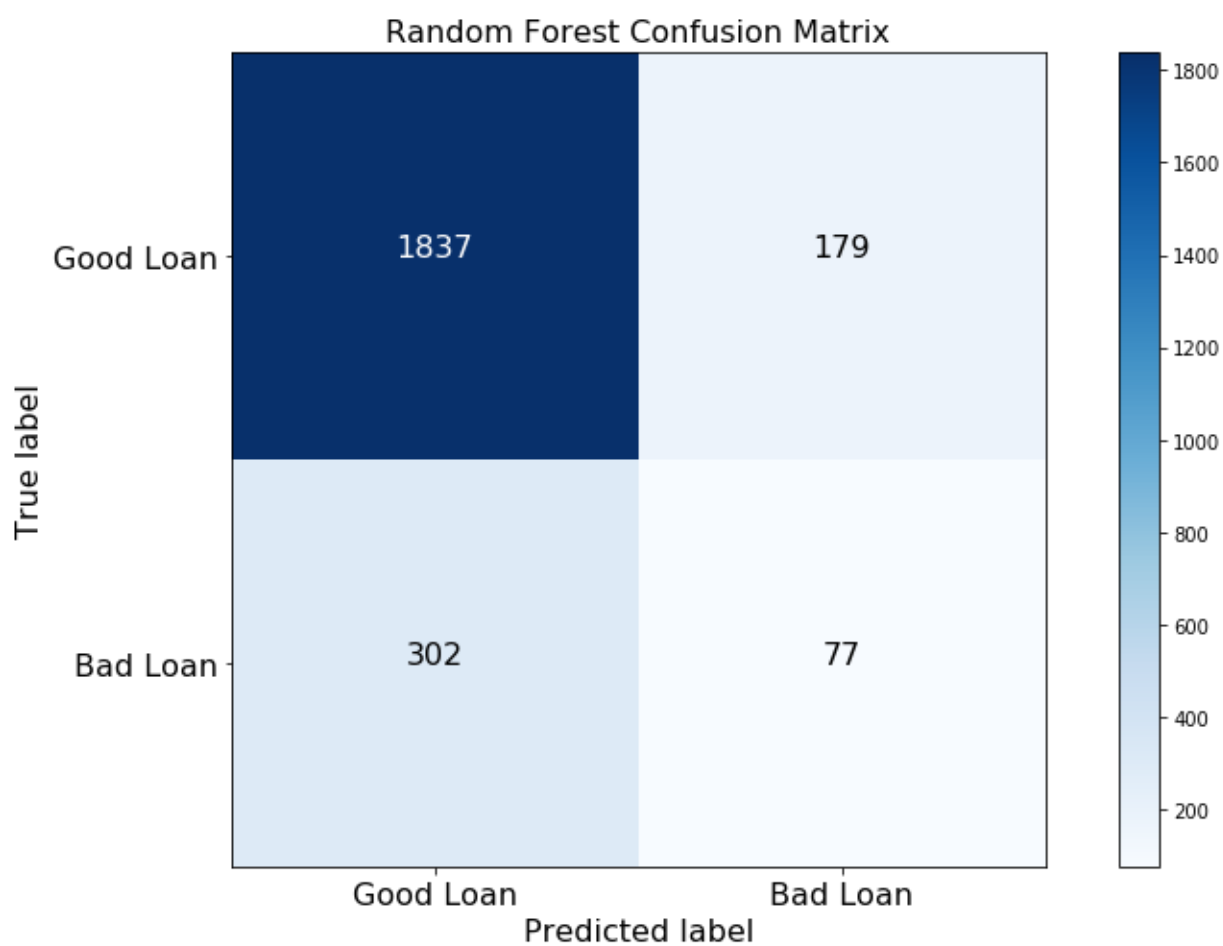


```
In [46]: 1 rf_sm = RandomForestClassifier(n_estimators=n_trees, max_depth=best_depth, n_
2 rf_sm.fit(X_train_sm, y_train_sm)
3
4 train_acc, test_acc = rf_sm.score(X_train_sm, y_train_sm), rf_sm.score(X_test
5
6 print('Random Forest: Optimal depth={}'.format(best_depth))
7
8 report_lr = precision_recall_fscore_support(y_test, rf_sm.predict(X_test), av
9
10 rf_sm_results = {
11     'model': 'Random Forest',
12     'train_acc': train_acc,
13     'test_acc': test_acc,
14     'precision': report_lr[0],
15     'recall': report_lr[1],
16     'F1': report_lr[2]
17 }
18
19 print('Random Forest: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}
20       format(rf_sm_results['train_acc'],
21             rf_sm_results['test_acc'],
22             rf_sm_results['precision'],
23             rf_sm_results['recall'],
24             rf_sm_results['F1']))
```

Random Forest: Optimal depth=19

Random Forest: accuracy on train=100.00%, test=79.92%, precision=0.30, recall=0.20, F1=0.24

```
In [47]: 1 # Plot confusion matrix
2 y_pred = rf_sm.predict(X_test)
3
4 cnf_matrix = confusion_matrix(y_test, y_pred)
5 np.set_printoptions(precision=2)
6
7 plt.figure(figsize=(11,7))
8 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
9                       title='Random Forest Confusion Matrix')
10
11 plt.show()
```

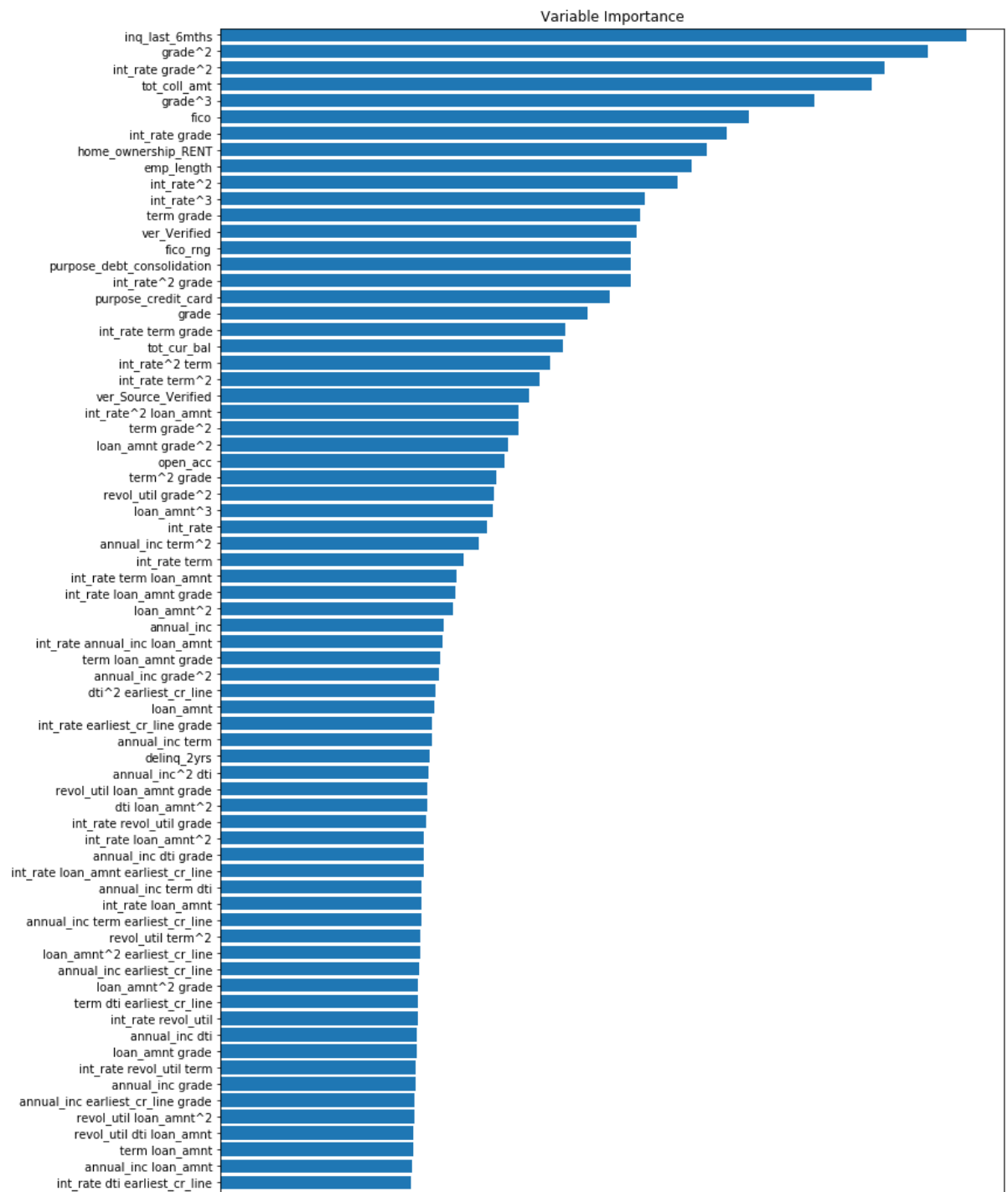


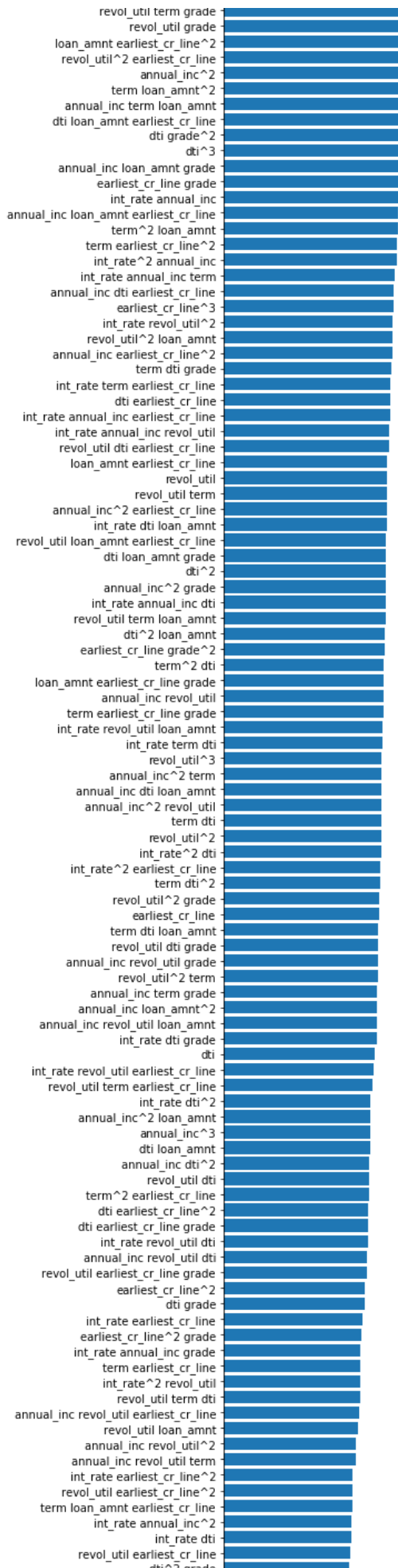
In [48]:

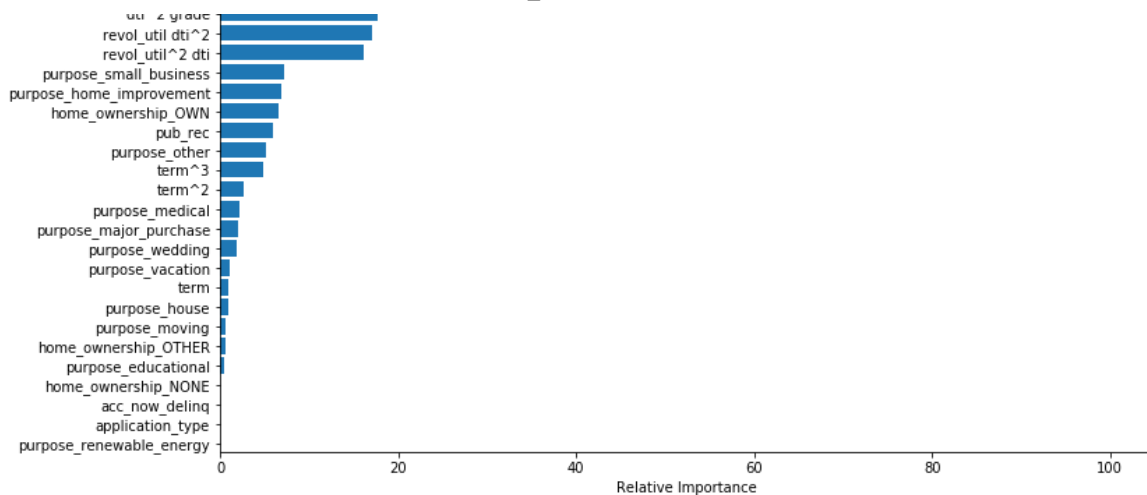
```

1 # Random Forest Feature Importance
2 feature_importance = rf_sm.feature_importances_
3 feature_importance = 100.0 * (feature_importance / feature_importance.max())
4 sorted_idx = np.argsort(feature_importance)
5 pos = np.arange(sorted_idx.shape[0])
6
7 # Plot
8 plt.figure(figsize=(12,50))
9 plt.barh(pos, feature_importance[sorted_idx], align='center')
10 plt.yticks(pos, X_train.columns[sorted_idx])
11 plt.xlabel('Relative Importance')
12 plt.title('Variable Importance')
13 plt.margins(y=0)
14 plt.show()

```







AddBoost

Raw training set

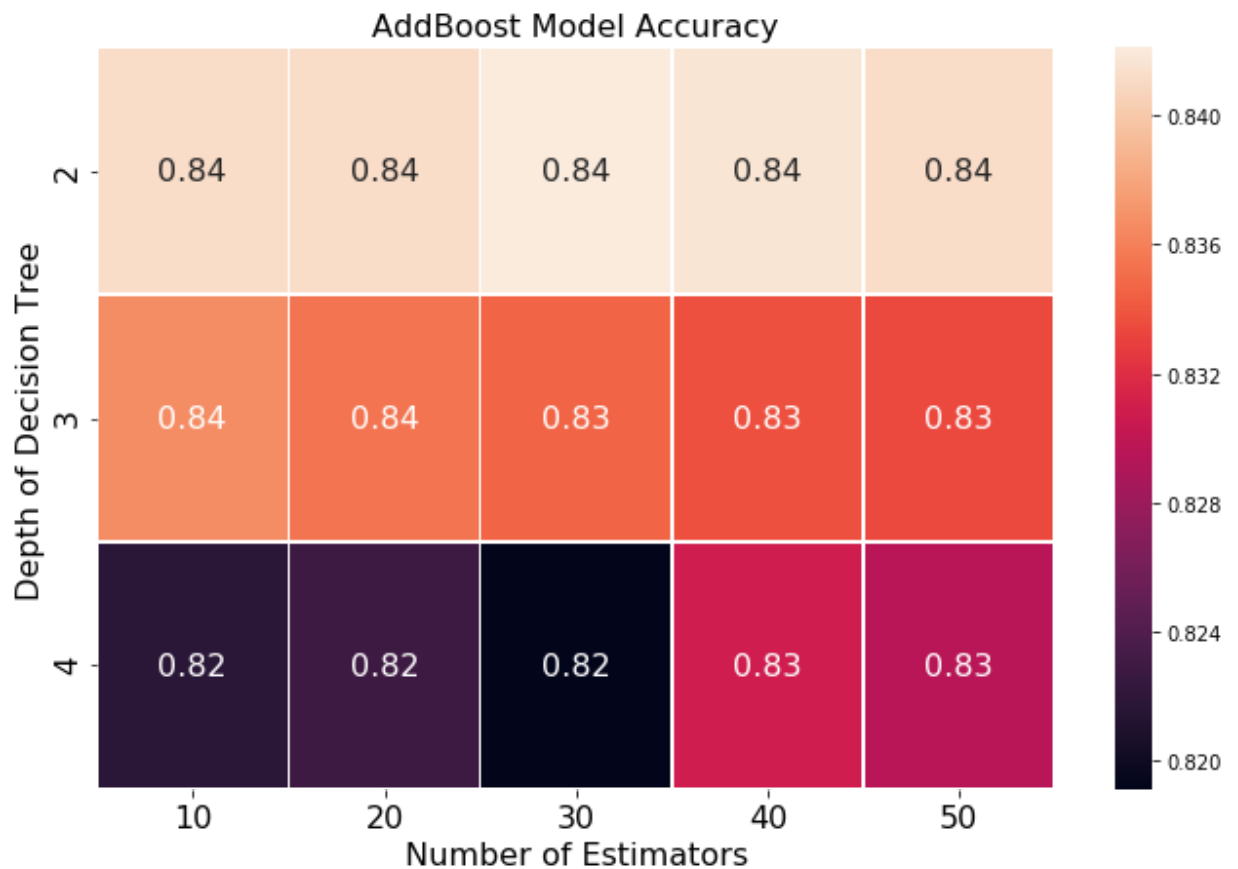
In [49]:

```

1  estimators = [10, 20, 30, 40, 50]
2  depths = list(range(2, 5))
3
4  train_scores = pd.DataFrame(index=depths, columns=estimators)
5  validation_scores = pd.DataFrame(index=depths, columns=estimators)
6  test_scores = pd.DataFrame(index=depths, columns=estimators)
7
8  best_score = 0
9  best_depth = 0
10 best_estimator = 0
11
12 for e in estimators:
13     for d in depths:
14         ab = AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_dep
15                                 n_estimators=e, learning_rate=0.05)
16         ab.fit(X_train, y_train)
17
18         train_scores.loc[d, e] = ab.score(X_train, y_train)
19         test_scores.loc[d, e] = ab.score(X_test, y_test)
20
21         val_score = cross_val_score(estimator=ab, X=X_train, y=y_train, cv=5)
22         validation_scores.loc[d, e] = val_score
23
24         if val_score > best_score:
25             best_depth = d
26             best_estimator = e
27             best_score = val_score

```

```
In [50]: 1 validation_scores = validation_scores.astype(float)
2
3 f, ax = plt.subplots(figsize=(11, 7))
4 sns.heatmap(validation_scores, annot=True, linewidths=.5, ax=ax, annot_kws={"
5 ax.set_ylabel('Depth of Decision Tree', fontsize=16)
6 ax.set_xlabel('Number of Estimators', fontsize=16)
7 ax.set_title('AddBoost Model Accuracy', fontsize=16)
8 ax.tick_params(labelsize=16)
```



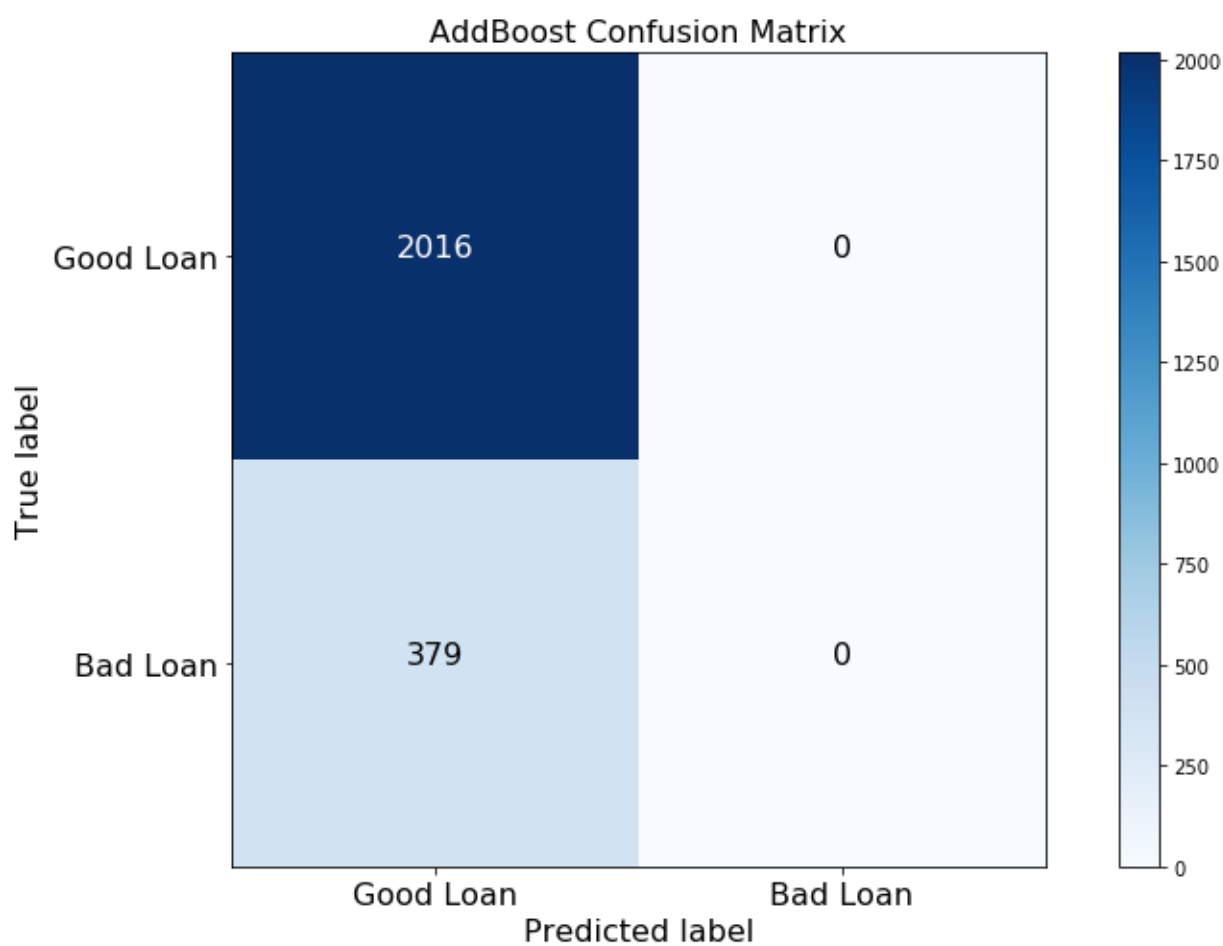
```
In [51]: 1 print('AddBoost: Optimal depth={}'.format(best_depth))
2 print('AddBoost: Optimal number of estimators={}'.format(best_estimator))
3
4 ab = AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=best_
5                         n_estimators=best_estimator, learning_rate=0.05)
6 ab.fit(X_train, y_train)
7
8 train_acc, test_acc = ab.score(X_train, y_train), ab.score(X_test, y_test)
9
10 report_lr = precision_recall_fscore_support(y_test, ab.predict(X_test), avera
11
12 ab_results = {
13     'model': 'AddBoost',
14     'train_acc': train_acc,
15     'test_acc': test_acc,
16     'precision': report_lr[0],
17     'recall': report_lr[1],
18     'F1': report_lr[2]
19 }
20
21 print('Random Forest: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}
22       format(ab_results['train_acc'],
23             ab_results['test_acc'],
24             ab_results['precision'],
25             ab_results['recall'],
26             ab_results['F1']))
```

AddBoost: Optimal depth=2

AddBoost: Optimal number of estimators=30

Random Forest: accuracy on train=84.17%, test=84.18%, precision=0.00, recall=0.00, F1=0.00

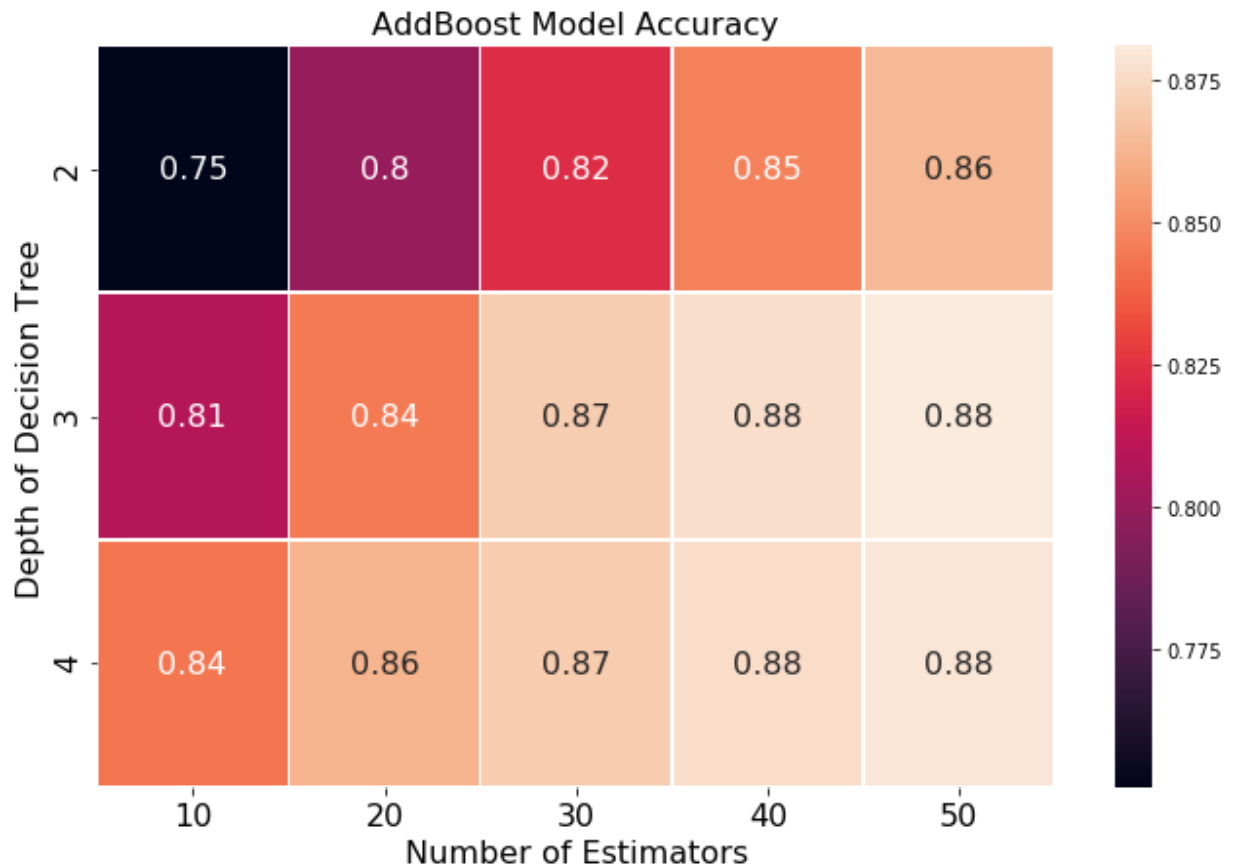
```
In [52]: 1 # Plot confusion matrix
2 y_pred = ab.predict(X_test)
3
4 cnf_matrix = confusion_matrix(y_test, y_pred)
5 np.set_printoptions(precision=2)
6
7 plt.figure(figsize=(11,7))
8 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
9                       title='AddBoost Confusion Matrix')
10
11 plt.show()
```



Oversampled Training set

```
In [53]: 1 estimators = [10, 20, 30, 40, 50]
2 depths = list(range(2, 5))
3
4 train_scores = pd.DataFrame(index=depths, columns=estimators)
5 validation_scores = pd.DataFrame(index=depths, columns=estimators)
6 test_scores = pd.DataFrame(index=depths, columns=estimators)
7
8 best_score = 0
9 best_depth = 0
10 best_estimator = 0
11
12 for e in estimators:
13     for d in depths:
14         ab = AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_dep
15                                 n_estimators=e, learning_rate=0.05)
16         ab.fit(X_train_sm, y_train_sm)
17
18         train_scores.loc[d, e] = ab.score(X_train_sm, y_train_sm)
19         test_scores.loc[d, e] = ab.score(X_test, y_test)
20
21         val_score = cross_val_score(estimator=ab, X=X_train_sm, y=y_train_sm,
22                                     validation_scores.loc[d, e] = val_score
23
24         if val_score > best_score:
25             best_depth = d
26             best_estimator = e
27             best_score = val_score
```

```
In [54]: 1 validation_scores = validation_scores.astype(float)
2
3 f, ax = plt.subplots(figsize=(11, 7))
4 sns.heatmap(validation_scores, annot=True, linewidths=.5, ax=ax, annot_kws={"
5 ax.set_ylabel('Depth of Decision Tree', fontsize=16)
6 ax.set_xlabel('Number of Estimators', fontsize=16)
7 ax.set_title('AddBoost Model Accuracy', fontsize=16)
8 ax.tick_params(labelsize=16)
```



```

In [55]: 1 print('AddBoost: Optimal depth={}'.format(best_depth))
2 print('AddBoost: Optimal number of estimators={}'.format(best_estimator))
3
4 ab_sm = AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=be
5                             n_estimators=best_estimator, learning_rate=0.05)
6 ab_sm.fit(X_train_sm, y_train_sm)
7
8 train_acc, test_acc = ab_sm.score(X_train_sm, y_train_sm), ab.score(X_test, y
9
10 report_lr = precision_recall_fscore_support(y_test, ab_sm.predict(X_test), av
11
12 ab_sm_results = {
13     'model': 'AddBoost',
14     'train_acc': train_acc,
15     'test_acc': test_acc,
16     'precision': report_lr[0],
17     'recall': report_lr[1],
18     'F1': report_lr[2]
19 }
20
21 print('Random Forest: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}
22       format(ab_sm_results['train_acc'],
23             ab_sm_results['test_acc'],
24             ab_sm_results['precision'],
25             ab_sm_results['recall'],
26             ab_sm_results['F1']))

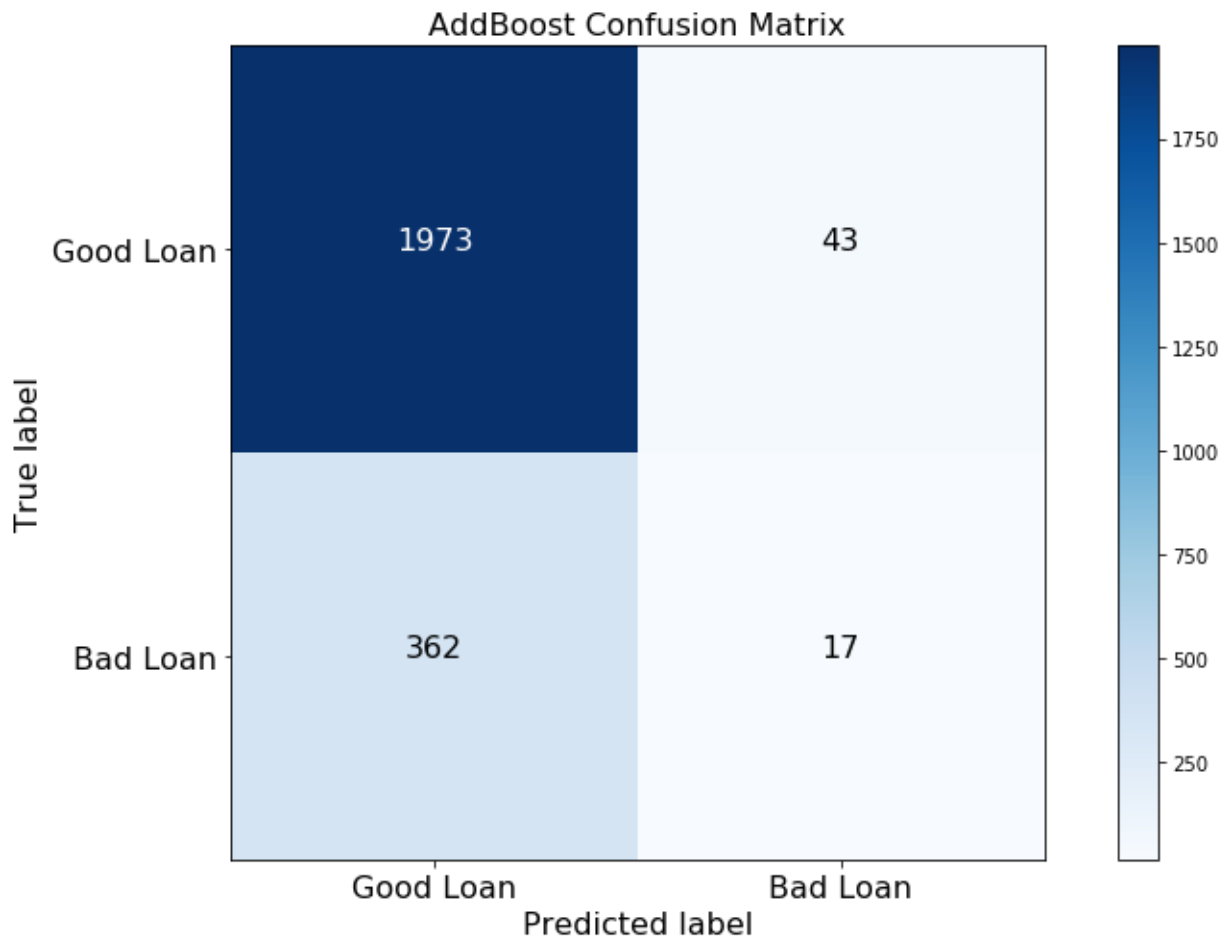
```

AddBoost: Optimal depth=3

AddBoost: Optimal number of estimators=50

Random Forest: accuracy on train=90.79%, test=82.13%, precision=0.28, recall=0.04, F1=0.08


```
In [56]: 1 # Plot confusion matrix
2 y_pred = ab_sm.predict(X_test)
3
4 cnf_matrix = confusion_matrix(y_test, y_pred)
5 np.set_printoptions(precision=2)
6
7 plt.figure(figsize=(11,7))
8 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
9                       title='AddBoost Confusion Matrix')
10
11 plt.show()
```



Neural Network

Raw training set

```
In [12]: 1 X_train_arr = X_train.values
2 y_train_arr = to_categorical(y_train.values)
3 X_test_arr = X_test.values
4 y_test_arr = to_categorical(y_test.values)
```

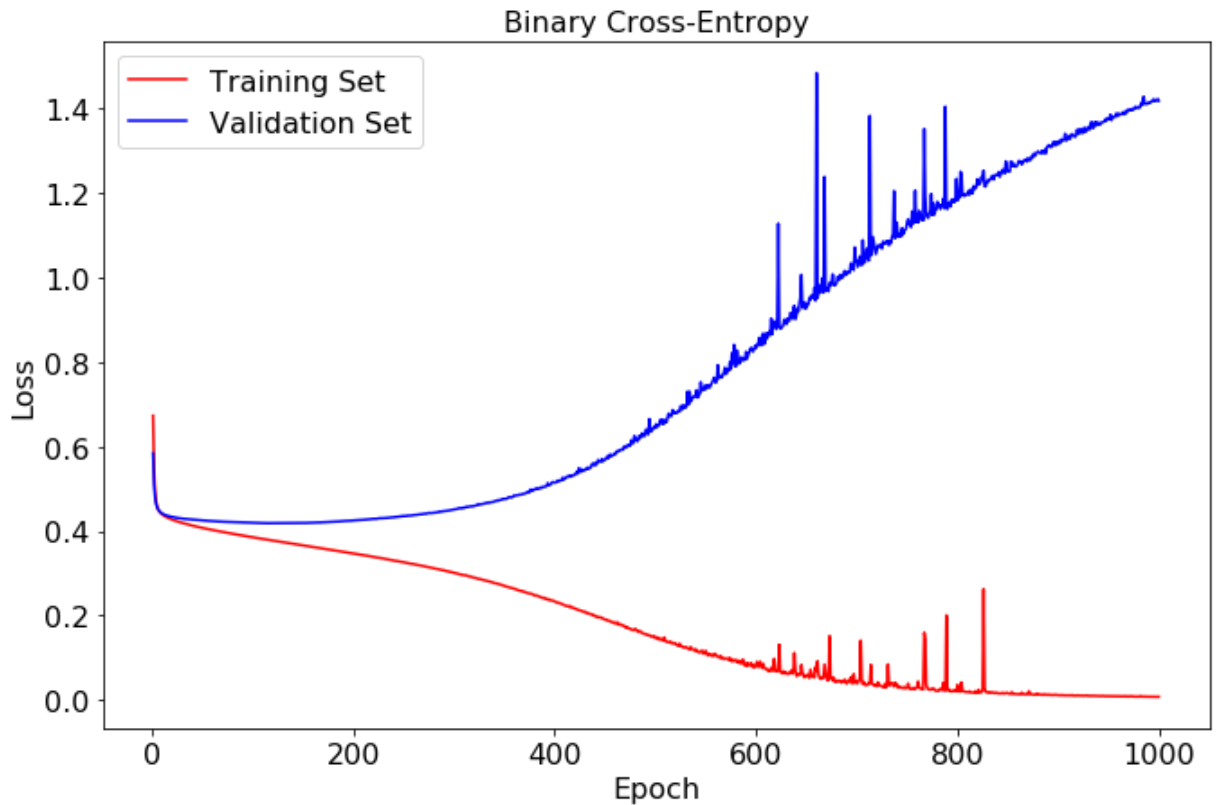
```

In [13]: 1 H = 100 # number of nodes in the layer
2 input_dim = X_train.shape[1] # input dimension
3 output_dim = 2 # output dimension
4
5 nn = Sequential() # create sequential multi-layer perceptron
6
7 # Layer 0, our hidden layer
8 nn.add(Dense(H, input_dim=input_dim, activation='relu'))
9
10 # Layer 1, our hidden layer
11 nn.add(Dense(H, activation='relu'))
12
13 # Layer 2, our hidden layer
14 nn.add(Dense(H, activation='relu'))
15
16 # Layer 3
17 nn.add(Dense(output_dim, activation='sigmoid'))
18
19 # compile the model
20 nn.compile(loss='binary_crossentropy', optimizer='sgd')
21 nn.summary()
22
23 epochs = 1000
24 batch_size = 128
25 validation_split = 0.5
26
27 nn_history = nn.fit(X_train_arr, y_train_arr,
28                    batch_size=batch_size,
29                    epochs=epochs, verbose=False,
30                    shuffle = True, validation_data = (X_test_arr, y_test_arr))

```

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 100)	19500
dense_2 (Dense)	(None, 100)	10100
dense_3 (Dense)	(None, 100)	10100
dense_4 (Dense)	(None, 2)	202
Total params: 39,902		
Trainable params: 39,902		
Non-trainable params: 0		

```
In [14]: 1 fig, ax = plt.subplots(1, 1, figsize=(11,7))
2 ax.plot(range(1, epochs + 1), nn_history.history['loss'], 'r', label='Trainin
3 ax.plot(range(1, epochs + 1), nn_history.history['val_loss'], 'b', label='Val
4 ax.set_title('Binary Cross-Entropy', fontsize=16)
5 ax.set_xlabel('Epoch', fontsize=16)
6 ax.set_ylabel('Loss', fontsize=16)
7 ax.legend(fontsize=16)
8 ax.tick_params(labelsize=16)
```



```

In [15]: 1 H = 100 # number of nodes in the layer
          2 input_dim = X_train.shape[1] # input dimension
          3 output_dim = 2 # output dimension
          4
          5 nn = Sequential() # create sequential multi-layer perceptron
          6
          7 # Layer 0, our hidden layer
          8 nn.add(Dense(H, input_dim=input_dim, activation='relu'))
          9
          10 # Layer 1, our hidden layer
          11 nn.add(Dense(H, activation='relu'))
          12
          13 # Layer 2, our hidden layer
          14 nn.add(Dense(H, activation='relu'))
          15
          16 # Layer 3
          17 nn.add(Dense(output_dim, activation='sigmoid'))
          18
          19 # compile the model
          20 nn.compile(loss='binary_crossentropy', optimizer='sgd')
          21 nn.summary()
          22
          23 epochs = 160
          24 batch_size = 128
          25
          26 nn_history = nn.fit(X_train_arr, y_train_arr,
          27                   batch_size=batch_size,
          28                   epochs=epochs, verbose=False,
          29                   shuffle = True, validation_data = (X_test_arr, y_test_arr))

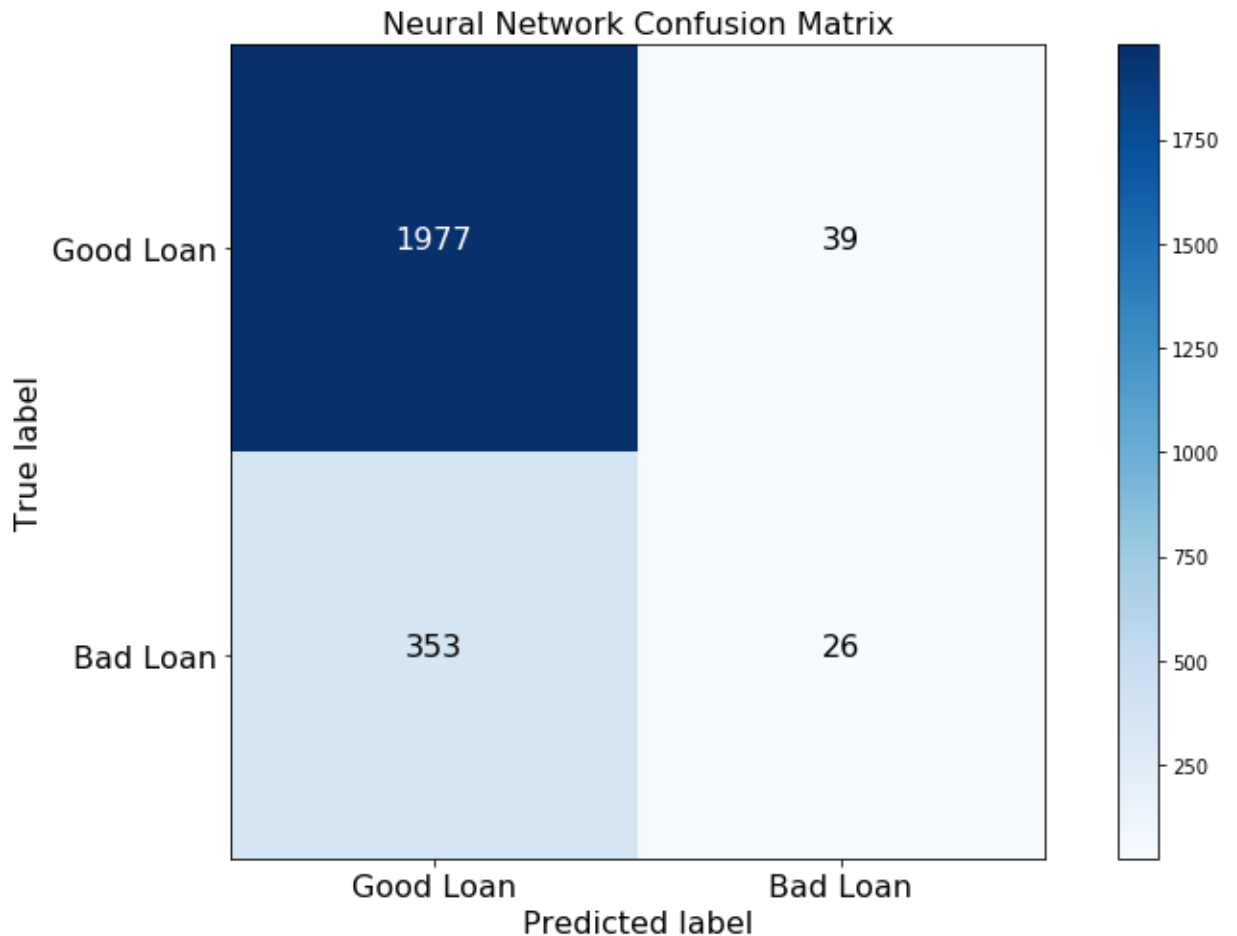
```

Layer (type)	Output Shape	Param #
=====		
dense_5 (Dense)	(None, 100)	19500
dense_6 (Dense)	(None, 100)	10100
dense_7 (Dense)	(None, 100)	10100
dense_8 (Dense)	(None, 2)	202
=====		
Total params: 39,902		
Trainable params: 39,902		
Non-trainable params: 0		

```
In [16]: 1 nn_train_df = pd.DataFrame(nn.predict(X_train_arr))
2 nn_train_df['pred'] = nn_train_df.apply(lambda x: 1 if x[0] < x[1] else 0, axis=1)
3 nn_train_acc = accuracy_score(y_train, nn_train_df['pred'].values)
4
5 nn_test_df = pd.DataFrame(nn.predict(X_test_arr))
6 nn_test_df['pred'] = nn_test_df.apply(lambda x: 1 if x[0] < x[1] else 0, axis=1)
7 nn_test_acc = accuracy_score(y_test, nn_test_df['pred'].values)
8
9 report_lr = precision_recall_fscore_support(y_test, nn_test_df['pred'], average='micro')
10
11 nn_results = {
12     'model': 'Neural Network',
13     'train_acc': nn_train_acc,
14     'test_acc': nn_test_acc,
15     'precision': report_lr[0],
16     'recall': report_lr[1],
17     'F1': report_lr[2]
18 }
19
20 print('Neural Network: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={:.2f}, F1={:.2f}'.format(
21     nn_results['train_acc'],
22     nn_results['test_acc'],
23     nn_results['precision'],
24     nn_results['recall'],
25     nn_results['F1']))
```

Neural Network: accuracy on train=85.96%, test=83.63%, precision=0.40, recall=0.07, F1=0.12

```
In [17]: 1 # Plot confusion matrix
2 cnf_matrix = confusion_matrix(y_test, nn_test_df['pred'].values)
3 np.set_printoptions(precision=2)
4
5 plt.figure(figsize=(11,7))
6 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
7                        title='Neural Network Confusion Matrix')
8
9 plt.show()
```



Oversampled Training set

```
In [18]: 1 X_train_arr = X_train_sm
2 y_train_arr = to_categorical(y_train_sm)
3 X_test_arr = X_test.values
4 y_test_arr = to_categorical(y_test.values)
```

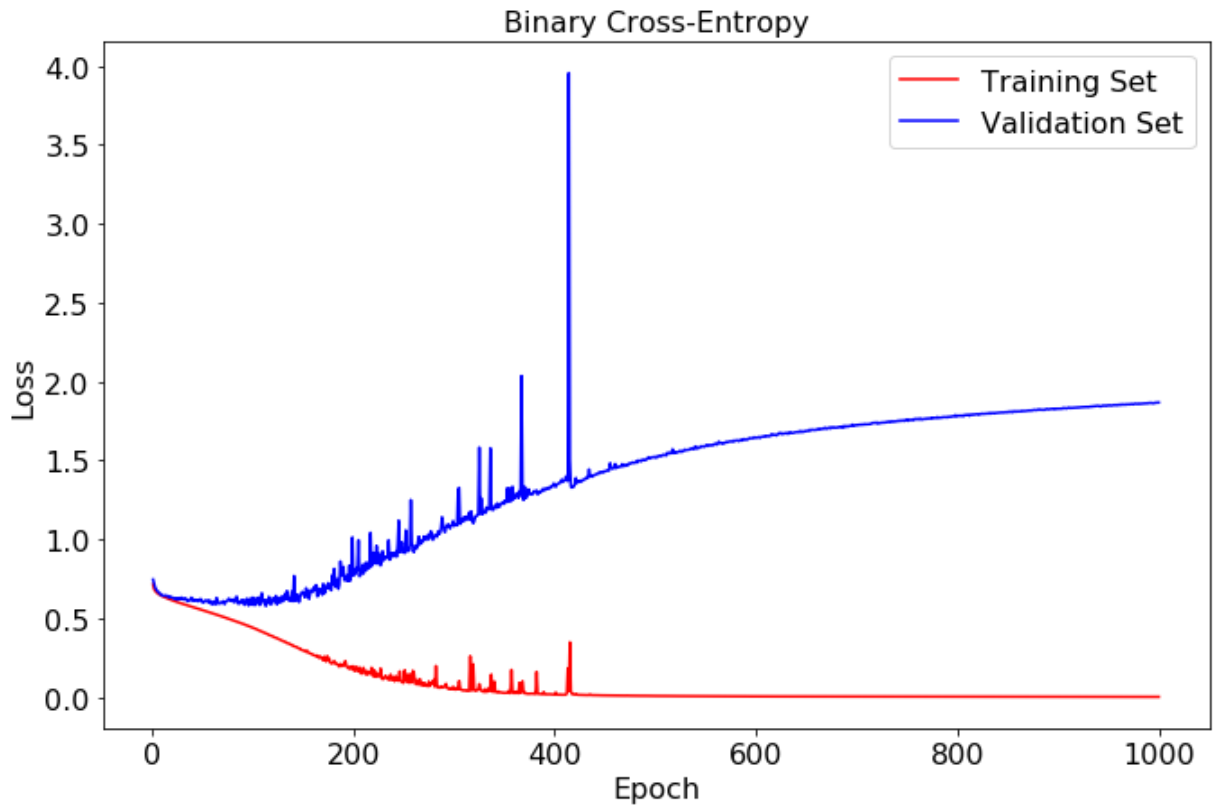
```

In [19]: 1 H = 100 # number of nodes in the layer
2 input_dim = X_train_sm.shape[1] # input dimension
3 output_dim = 2 # output dimension
4
5 nn = Sequential() # create sequential multi-layer perceptron
6
7 # Layer 0, our hidden layer
8 nn.add(Dense(H, input_dim=input_dim, activation='relu'))
9
10 # Layer 1, our hidden layer
11 nn.add(Dense(H, activation='relu'))
12
13 # Layer 2, our hidden layer
14 nn.add(Dense(H, activation='relu'))
15
16 # Layer 3
17 nn.add(Dense(output_dim, activation='sigmoid'))
18
19 # compile the model
20 nn.compile(loss='binary_crossentropy', optimizer='sgd')
21 nn.summary()
22
23 epochs = 1000
24 batch_size = 128
25 validation_split = 0.5
26
27 nn_history = nn.fit(X_train_arr, y_train_arr,
28                     batch_size=batch_size,
29                     epochs=epochs, verbose=False,
30                     shuffle = True, validation_data = (X_test_arr, y_test_arr))

```

Layer (type)	Output Shape	Param #
=====	=====	=====
dense_9 (Dense)	(None, 100)	19500
dense_10 (Dense)	(None, 100)	10100
dense_11 (Dense)	(None, 100)	10100
dense_12 (Dense)	(None, 2)	202
=====	=====	=====
Total params: 39,902		
Trainable params: 39,902		
Non-trainable params: 0		

```
In [20]: 1 fig, ax = plt.subplots(1, 1, figsize=(11,7))
2 ax.plot(range(1, epochs + 1), nn_history.history['loss'], 'r', label='Trainin
3 ax.plot(range(1, epochs + 1), nn_history.history['val_loss'], 'b', label='Val
4 ax.set_title('Binary Cross-Entropy', fontsize=16)
5 ax.set_xlabel('Epoch', fontsize=16)
6 ax.set_ylabel('Loss', fontsize=16)
7 ax.legend(fontsize=16)
8 ax.tick_params(labelsize=16)
```




```

In [21]: 1 H = 100 # number of nodes in the layer
          2 input_dim = X_train_sm.shape[1] # input dimension
          3 output_dim = 2 # output dimension
          4
          5 nn_sm = Sequential() # create sequential multi-layer perceptron
          6
          7 # Layer 0, our hidden layer
          8 nn_sm.add(Dense(H, input_dim=input_dim, activation='relu'))
          9
          10 # Layer 1, our hidden layer
          11 nn_sm.add(Dense(H, activation='relu'))
          12
          13 # Layer 2, our hidden layer
          14 nn_sm.add(Dense(H, activation='relu'))
          15
          16 # Layer 3
          17 nn_sm.add(Dense(output_dim, activation='sigmoid'))
          18
          19 # compile the model
          20 nn_sm.compile(loss='binary_crossentropy', optimizer='sgd')
          21 nn_sm.summary()
          22
          23 epochs = 100
          24 batch_size = 128
          25
          26 nn_sm_history = nn_sm.fit(X_train_arr, y_train_arr,
          27                          batch_size=batch_size,
          28                          epochs=epochs, verbose=False,
          29                          shuffle = True, validation_data = (X_test_arr, y_te

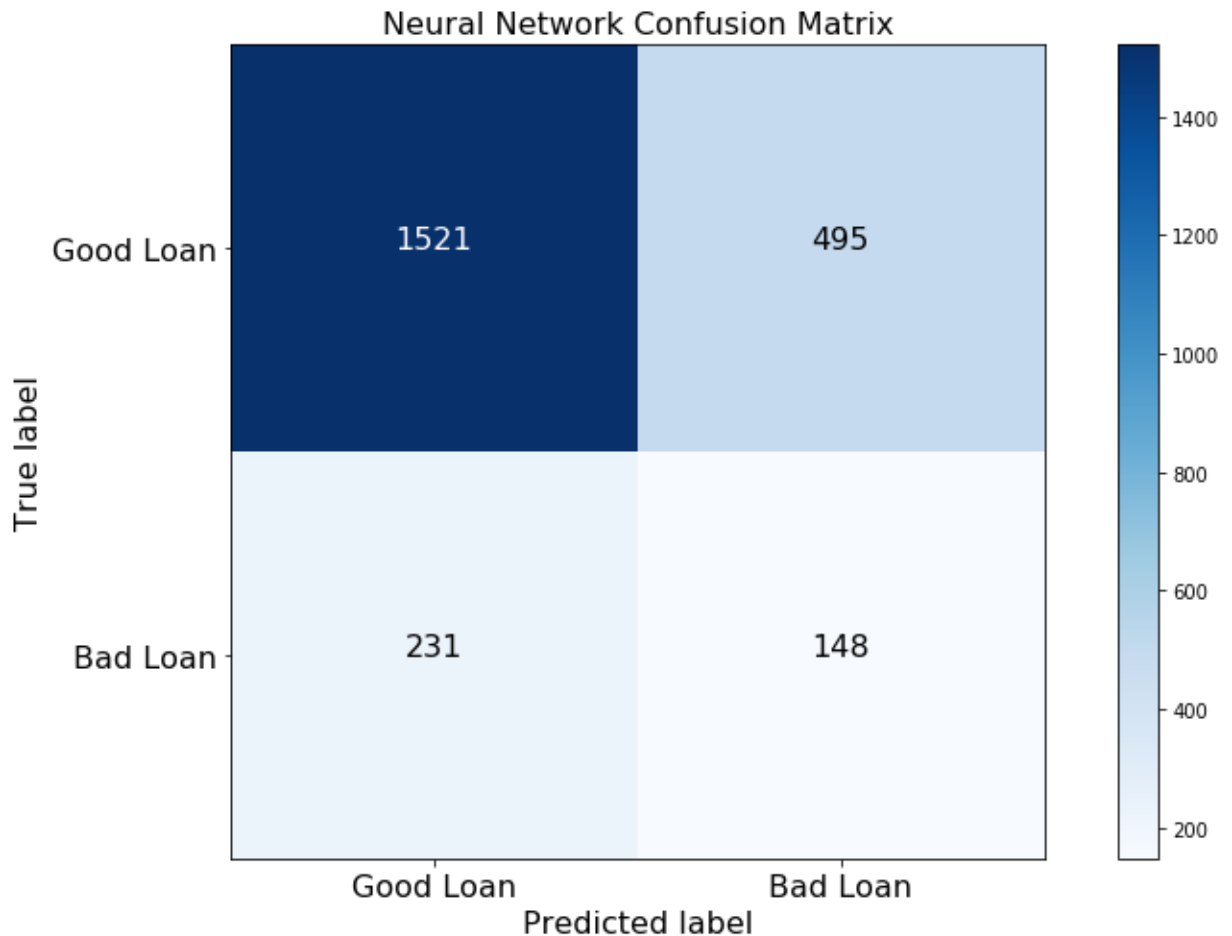
```

Layer (type)	Output Shape	Param #
=====		
dense_13 (Dense)	(None, 100)	19500
dense_14 (Dense)	(None, 100)	10100
dense_15 (Dense)	(None, 100)	10100
dense_16 (Dense)	(None, 2)	202
=====		
Total params: 39,902		
Trainable params: 39,902		
Non-trainable params: 0		

```
In [22]: 1 nn_train_df = pd.DataFrame(nn_sm.predict(X_train_arr))
2 nn_train_df['pred'] = nn_train_df.apply(lambda x: 1 if x[0] < x[1] else 0, axis=1)
3 nn_train_acc = accuracy_score(y_train_sm, nn_train_df['pred'].values)
4
5 nn_test_df = pd.DataFrame(nn_sm.predict(X_test_arr))
6 nn_test_df['pred'] = nn_test_df.apply(lambda x: 1 if x[0] < x[1] else 0, axis=1)
7 nn_test_acc = accuracy_score(y_test, nn_test_df['pred'].values)
8
9 report_lr = precision_recall_fscore_support(y_test, nn_test_df['pred'], average='micro')
10
11 nn_sm_results = {
12     'model': 'Neural Network',
13     'train_acc': nn_train_acc,
14     'test_acc': nn_test_acc,
15     'precision': report_lr[0],
16     'recall': report_lr[1],
17     'F1': report_lr[2]
18 }
19
20 print('Neural Network: accuracy on train={:.2%}, test={:.2%}, precision={:.2f}, recall={:.2f}, F1={:.2f}'.format(
21     nn_sm_results['train_acc'],
22     nn_sm_results['test_acc'],
23     nn_sm_results['precision'],
24     nn_sm_results['recall'],
25     nn_sm_results['F1']))
```

Neural Network: accuracy on train=84.74%, test=69.69%, precision=0.23, recall=0.39, F1=0.29

```
In [23]: 1 # Plot confusion matrix
2 cnf_matrix = confusion_matrix(y_test, nn_test_df['pred'].values)
3 np.set_printoptions(precision=2)
4
5 plt.figure(figsize=(11,7))
6 plot_confusion_matrix(cnf_matrix, classes=["Good Loan", "Bad Loan"], normalize
7                        title='Neural Network Confusion Matrix')
8
9 plt.show()
```



Output Model Results

```
In [25]: 1 # Output results
2 results = pd.DataFrame([logit_results,
3                          knn_results,
4                          lda_results,
5                          qda_results,
6                          tree_results,
7                          rf_results,
8                          ab_results,
9                          nn_results])
10 results.set_index('model', inplace=True)
11 results.to_csv("data/results.csv", index=True)
```

```
In [70]: 1 # Output results
2 results_sm = pd.DataFrame([logit_sm_results,
3                             knn_sm_results,
4                             lda_sm_results,
5                             qda_sm_results,
6                             tree_sm_results,
7                             rf_sm_results,
8                             ab_sm_results,
9                             nn_sm_results])
10 results_sm.set_index('model', inplace=True)
11 results_sm.to_csv("data/results_sm.csv", index=True)
```

```
In [ ]: 1
```