

Credit Card Fraud Detection

The Problem

Our goal is to predict if a credit card transaction is fraudulent or not. Because this data is so skewed, we'll use different techniques to try to make our predictive model work even when there is a majority class.

Collecting Data

This data set came from Kaggle (<https://www.kaggle.com/datasets/mig-ulb/creditcardfraud> (<https://www.kaggle.com/datasets/mig-ulb/creditcardfraud>)). We'll need to use certain libraries to accomplish our goal moving forward. I'll import Pandas, Numpy, Matplotlib and Seaborn initially.

Importing Libraries

```
In [9]: 1 # scientific computing libraries
2 import pandas as pd
3 import numpy as np
4
5 # data mining libraries
6 from sklearn.tree import DecisionTreeClassifier
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.preprocessing import LabelEncoder
9 from sklearn.decomposition import PCA#, FastICA
10 from sklearn.model_selection import train_test_split, KFold, StratifiedKFold, GridSearchCV, learning_curve
11 from sklearn import svm
12 from sklearn.linear_model import LogisticRegression
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.metrics import roc_curve, auc, confusion_matrix, accuracy_score, f1_score, precision_score, recall_score, r
15
16 from imblearn.pipeline import make_pipeline, Pipeline
17 from imblearn.over_sampling import SMOTE
18
19 #plot libraries
20 import plotly
21 import plotly.graph_objs as go
22 import plotly.figure_factory as ff
23 from plotly.offline import init_notebook_mode
24 init_notebook_mode(connected=True) # to show plots in notebook
25 from matplotlib import pyplot as plt
26 %matplotlib inline
27 import seaborn as sns
28
29 # offline plotly
30 from plotly.offline import plot, iplot
31
32 # do not show any warnings
33 import warnings
34 warnings.filterwarnings('ignore')
35
36 SEED = 17 # specify seed for reproducible results
37 pd.set_option('display.max_columns', None) # prevents abbreviation (with '...') of columns in prints
38
39 import os
40 for dirname, _, filenames in os.walk('/kaggle/input'):
41     for filename in filenames:
42         print(os.path.join(dirname, filename))
```

Loading the Data frame

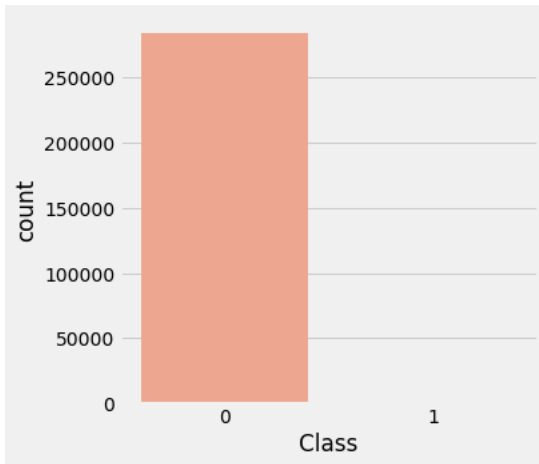
```
In [10]: 1 df = pd.read_csv("creditcard.csv")
2 df.head()
```

Out[10]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.3111
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.1437
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.1659
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.2879
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.1196

Checking the target feature, class as Fraud or not fraud

```
In [11]: 1 # apply the Fivethirtyeight style to plots.
2 plt.style.use("fivethirtyeight")
3
4 # display a frequency distribution for churn.
5 plt.figure(figsize=(5, 5))
6 ax = sns.countplot(x=df['Class'], palette='Reds', linewidth=1)
7 plt.show()
```



The plot shows a class imbalance of the data between Class (fraud) and class (not-fraud). To address this, resampling would be a suitable approach. To keep this case simple, the imbalance is kept forward and specific metrics are chosen for model evaluations.

```
In [12]: 1 df['Class'].value_counts(normalize=True)
```

```
Out[12]: 0    0.998273
1    0.001727
Name: Class, dtype: float64
```

It looks like they're losing 99.8% of their customers transactions are not fraudulent. Only 0.17% are.

So our original dataset is very imbalanced. Since nearly all transactions are non-fraudulent, our predictive models might make assumptions. We don't want that. We want it to actually catch fraudulent transactions.

To handle class imbalance in this classification problem, several strategies can be employed:

- Collect more data (not applicable in this case)
- Change performance metric:
- Precision, Recall, and F1 Score using confusion matrix
- Kappa for normalized accuracy
- ROC curves to calculate sensitivity/specificity ratio
- Resample the dataset:
- Over-sampling by adding copies of under-represented class (for little data)
- Under-sampling by removing instances from over-represented class (for lots of data)

```
In [13]: 1 from sklearn.preprocessing import StandardScaler
2
3 amount_col = df['Amount'].values.reshape(-1, 1)
4 scaler = StandardScaler().fit(amount_col)
5 df['normAmount'] = scaler.transform(amount_col)
6 df.drop(['Time', 'Amount'], axis=1, inplace=True)
7 df.head()
```

```
Out[13]:
```

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14
0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.311169
1	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.143772
2	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.165946
3	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.287924
4	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.119670

Assign X and Y without resampling

```
In [14]: 1 X = df.loc[:, df.columns != 'Class']
2 y = df.loc[:, df.columns == 'Class']
```

Undersample

```
In [15]: 1 # Find the indices of the minority class
2 fraud_indices = df[df.Class == 1].index
3
4 # Find the indices of the majority class
5 normal_indices = df[df.Class == 0].index
6
7 # Sample an equal number of instances from the majority class as the minority class
8 normal_sample = df.loc[np.random.choice(normal_indices, len(fraud_indices), replace=False), :]
9
10 # Concatenate the minority class and the sampled majority class
11 under_sample_data = pd.concat([df.loc[fraud_indices], normal_sample])
12
13 # Split the data into X (features) and y (Labels)
14 X_undersample = under_sample_data.drop("Class", axis=1)
15 y_undersample = under_sample_data["Class"]
16
17 # Calculate class proportions
18 normal_prop = len(under_sample_data[under_sample_data.Class == 0]) / len(under_sample_data)
19 fraud_prop = len(under_sample_data[under_sample_data.Class == 1]) / len(under_sample_data)
20
21 # Print class proportions
22 print("Percentage of normal transactions: {:.2f}%".format(normal_prop * 100))
23 print("Percentage of fraud transactions: {:.2f}%".format(fraud_prop * 100))
24 print("Total number of transactions in resampled data: {}".format(len(under_sample_data)))
```

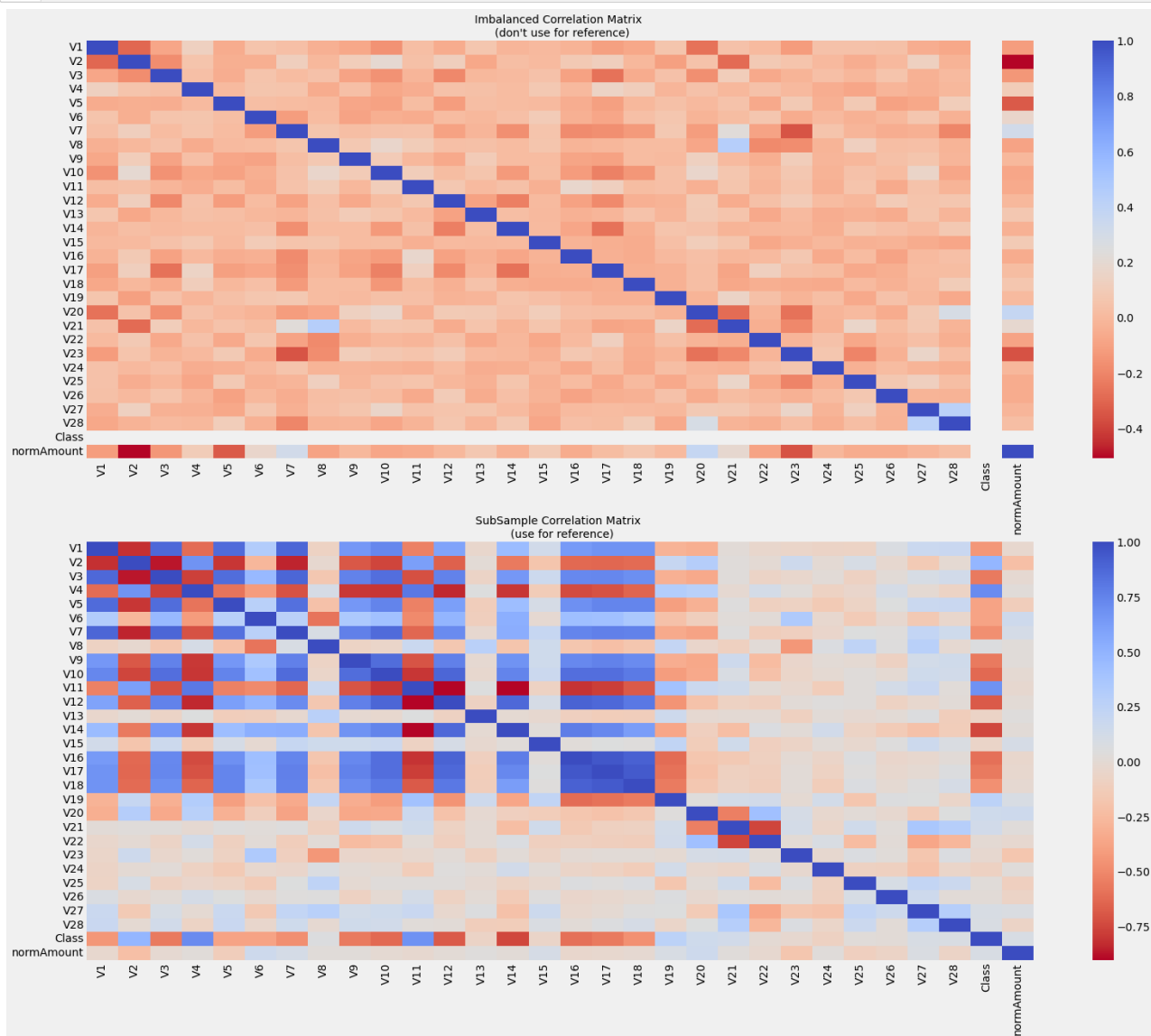
```
Percentage of normal transactions: 50.00%
Percentage of fraud transactions: 50.00%
Total number of transactions in resampled data: 984
```

Much better!

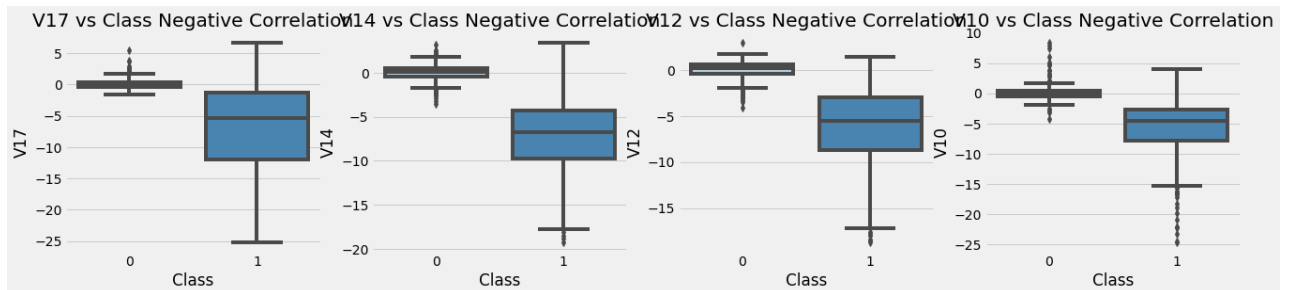
```

In [16]: 1 # Make sure we use the subsample in our correlation
2
3 f, (ax1, ax2) = plt.subplots(2, 1, figsize=(24,20))
4
5 # Entire DataFrame
6 corr = normal_sample.corr()
7 sns.heatmap(corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax1)
8 ax1.set_title("Imbalanced Correlation Matrix \n (don't use for reference)", fontsize=14)
9
10 sub_sample_corr = under_sample_data.corr()
11 sns.heatmap(sub_sample_corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax2)
12 ax2.set_title('SubSample Correlation Matrix \n (use for reference)', fontsize=14)
13 plt.show()
14

```



```
In [17]: 1 f, axes = plt.subplots(ncols=4, figsize=(20,4))
2
3 # Negative Correlations with our Class (The lower our feature value the more likely it will be a fraud transaction)
4 sns.boxplot(x="Class", y="V17", data=under_sample_data, palette='Blues', ax=axes[0])
5 axes[0].set_title('V17 vs Class Negative Correlation')
6
7 sns.boxplot(x="Class", y="V14", data=under_sample_data, palette='Blues', ax=axes[1])
8 axes[1].set_title('V14 vs Class Negative Correlation')
9
10
11 sns.boxplot(x="Class", y="V12", data=under_sample_data, palette='Blues', ax=axes[2])
12 axes[2].set_title('V12 vs Class Negative Correlation')
13
14
15 sns.boxplot(x="Class", y="V10", data=under_sample_data, palette='Blues', ax=axes[3])
16 axes[3].set_title('V10 vs Class Negative Correlation')
17
18 plt.show()
```



Train Test Split

```
In [18]: 1 from sklearn.model_selection import train_test_split
2
3 # Split the whole dataset into training and testing
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
5
6 # Print the size of each dataset
7 print("Number of transactions in the training dataset: ", len(X_train))
8 print("Number of transactions in the testing dataset: ", len(X_test))
9 print("Total number of transactions: ", len(X_train) + len(X_test))
10
11 # Split the undersampled dataset into training and testing
12 X_train_undersample, X_test_undersample, y_train_undersample, y_test_undersample = train_test_split(X_undersample, y_undersample, test_size=0.3, random_state=0)
13
14 # Print the size of each dataset
15 print("\nNumber of transactions in the training dataset: ", len(X_train_undersample))
16 print("Number of transactions in the testing dataset: ", len(X_test_undersample))
17 print("Total number of transactions: ", len(X_train_undersample) + len(X_test_undersample))
```

Number of transactions in the training dataset: 199364
 Number of transactions in the testing dataset: 85443
 Total number of transactions: 284807

Number of transactions in the training dataset: 688
 Number of transactions in the testing dataset: 296
 Total number of transactions: 984

Logistic Regression

We are highly focused on recall score as it helps us identify the maximum number of fraudulent transactions. Recall, along with accuracy and precision, is a key metric for evaluating a confusion matrix. Recall calculates the number of true positive instances divided by the sum of true positive instances and false negative instances.

Since our dataset is imbalanced, many normal transactions may be predicted as fraudulent, leading to false negatives. Recall helps us address this issue. Increasing recall may lead to a decrease in precision, but that is acceptable in our scenario, as predicting a fraudulent transaction as normal is not as critical as the opposite.

We can also apply a cost function to assign different weights to false negatives and false positives, but we'll leave that for another time.

```
In [19]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.model_selection import KFold, cross_val_score
3 from sklearn.metrics import confusion_matrix, precision_recall_curve, auc, roc_auc_score, roc_curve, recall_score, classification_report
```

```
In [20]: 1 from sklearn.model_selection import KFold
2 from sklearn.linear_model import LogisticRegression
3 import numpy as np
4 import pandas as pd
5 from sklearn.metrics import recall_score
6
7 def print_kfold_scores(x_train_data, y_train_data):
8     k_fold = KFold(n_splits=5, shuffle=False)
9     c_param_range = [0.01, 0.1, 1, 10, 100]
10    results_table = pd.DataFrame(index=range(len(c_param_range)), columns=['C_parameter', 'Mean recall score'])
11    j = 0
12    for c_param in c_param_range:
13        print('C parameter: ', c_param)
14        recall_accs = []
15        for train, test in k_fold.split(x_train_data):
16            lr = LogisticRegression(C=c_param, penalty='l2')
17            lr.fit(x_train_data.iloc[train], y_train_data.iloc[train])
18            y_pred = lr.predict(x_train_data.iloc[test])
19            recall = recall_score(y_train_data.iloc[test], y_pred)
20            recall_accs.append(recall)
21            print('Recall: ', recall)
22            results_table.loc[j, 'C_parameter'] = c_param
23            results_table.loc[j, 'Mean recall score'] = np.mean(recall_accs)
24            j += 1
25    best_c = results_table.loc[results_table['Mean recall score'].astype(float).idxmax()]['C_parameter']
26    print('Best model to choose from cross validation is with C parameter =', best_c)
27    return best_c
28
29 best_c = print_kfold_scores(X_train_undersample, y_train_undersample)
30 best_c
```

```
C parameter: 0.01
Recall: 0.821917808219178
Recall: 0.8493150684931506
Recall: 0.9152542372881356
Recall: 0.9324324324324325
Recall: 0.8787878787878788
C parameter: 0.1
Recall: 0.8356164383561644
Recall: 0.863013698630137
Recall: 0.9491525423728814
Recall: 0.9324324324324325
Recall: 0.8939393939393939
C parameter: 1
Recall: 0.8356164383561644
Recall: 0.863013698630137
Recall: 0.9661016949152542
Recall: 0.9324324324324325
Recall: 0.8939393939393939
C parameter: 10
Recall: 0.8493150684931506
Recall: 0.863013698630137
Recall: 0.9491525423728814
Recall: 0.9324324324324325
Recall: 0.8939393939393939
C parameter: 100
Recall: 0.8493150684931506
Recall: 0.863013698630137
Recall: 0.9661016949152542
Recall: 0.9459459459459459
Recall: 0.8787878787878788
Best model to choose from cross validation is with C parameter = 100
```

Out[20]: 100

Make a function to create confusion Matrix's

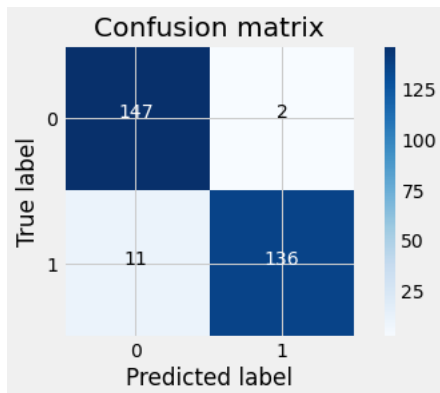
```
In [21]: 1 import itertools
2
3 def plot_confusion_matrix(cm, classes,
4                           normalize=False,
5                           title='Confusion matrix',
6                           cmap=plt.cm.Blues):
7     """
8     This function prints and plots the confusion matrix.
9     Normalization can be applied by setting `normalize=True`.
10    """
11    plt.imshow(cm, interpolation='nearest', cmap=cmap)
12    plt.title(title)
13    plt.colorbar()
14    tick_marks = np.arange(len(classes))
15    plt.xticks(tick_marks, classes, rotation=0)
16    plt.yticks(tick_marks, classes)
17
18    if normalize:
19        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
20        #print("Normalized confusion matrix")
21    else:
22        #print('Confusion matrix, without normalization')
23
24    #print(cm)
25
26    thresh = cm.max() / 2.
27    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
28        plt.text(j, i, cm[i, j],
29                horizontalalignment="center",
30                color="white" if cm[i, j] > thresh else "black")
31
32    plt.tight_layout()
33    plt.ylabel('True label')
34    plt.xlabel('Predicted label')
```

```

In [22]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import confusion_matrix
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 def plot_confusion_matrix(cm, classes,
7                           title='Confusion matrix',
8                           cmap=plt.cm.Blues):
9     plt.imshow(cm, interpolation='nearest', cmap=cmap)
10    plt.title(title)
11    plt.colorbar()
12    tick_marks = np.arange(len(classes))
13    plt.xticks(tick_marks, classes, rotation=0)
14    plt.yticks(tick_marks, classes)
15
16    thresh = cm.max() / 2.
17    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
18        plt.text(j, i, cm[i, j],
19                horizontalalignment="center",
20                color="white" if cm[i, j] > thresh else "black")
21
22    plt.tight_layout()
23    plt.ylabel('True label')
24    plt.xlabel('Predicted label')
25
26 # Use this C_parameter to build the final model with the whole training dataset and predict the classes in the test
27 # dataset
28 lr = LogisticRegression(C=best_c, penalty='l2')
29 lr.fit(X_train_undersample, y_train_undersample.values.ravel())
30 y_pred_undersample = lr.predict(X_test_undersample.values)
31
32 # Compute confusion matrix
33 cnf_matrix = confusion_matrix(y_test_undersample, y_pred_undersample)
34 np.set_printoptions(precision=2)
35
36 recall = cnf_matrix[1, 1] / (cnf_matrix[1, 0] + cnf_matrix[1, 1])
37 print("Recall metric in the testing dataset: ", recall)
38
39 # Plot non-normalized confusion matrix
40 class_names = [0, 1]
41 plt.figure()
42 plot_confusion_matrix(cnf_matrix, classes=class_names, title='Confusion matrix')
43 plt.show()

```

Recall metric in the testing dataset: 0.9251700680272109



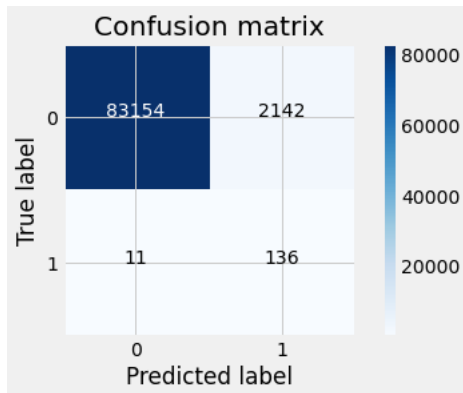
Our model predicts a 93.2% recall on the undersampled test set. We'll try it now with our whole data to see if it still works.


```

In [23]: 1 def plot_confusion_matrix(cm, classes,
2         title='Confusion matrix',
3         cmap=plt.cm.Blues):
4     plt.imshow(cm, interpolation='nearest', cmap=cmap)
5     plt.title(title)
6     plt.colorbar()
7     tick_marks = np.arange(len(classes))
8     plt.xticks(tick_marks, classes, rotation=0)
9     plt.yticks(tick_marks, classes)
10
11     thresh = cm.max() / 2.
12     for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
13         plt.text(j, i, cm[i, j],
14                 horizontalalignment="center",
15                 color="white" if cm[i, j] > thresh else "black")
16
17     plt.tight_layout()
18     plt.ylabel('True label')
19     plt.xlabel('Predicted label')
20
21 # Use this C_parameter to build the final model with the whole training dataset and predict the classes in the test
22 # dataset
23 lr = LogisticRegression(C=best_c, penalty='l2')
24 lr.fit(X_train_undersample, y_train_undersample.values.ravel())
25 y_pred = lr.predict(X_test.values)
26
27 # Compute confusion matrix
28 cnf_matrix = confusion_matrix(y_test, y_pred)
29 np.set_printoptions(precision=2)
30
31 recall = cnf_matrix[1, 1] / (cnf_matrix[1, 0] + cnf_matrix[1, 1])
32 print("Recall metric in the testing dataset: ", recall)
33
34 # Plot non-normalized confusion matrix
35 class_names = [0, 1]
36 plt.figure()
37 plot_confusion_matrix(cnf_matrix, classes=class_names, title='Confusion matrix')
38 plt.show()
39

```

Recall metric in the testing dataset: 0.9251700680272109



Still almost identical! That's great.

We'll check it again with our skewed data using Logistic Regression

```
In [24]: 1 best_c = print_kfold_scores(X_train,y_train)
```

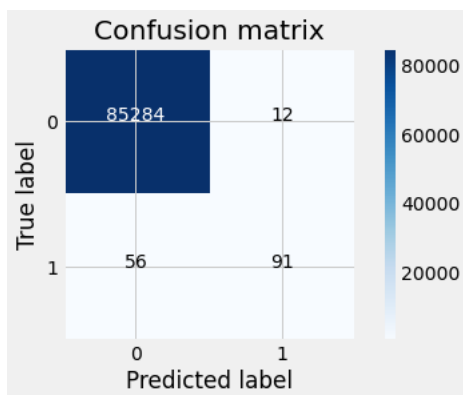
```
C parameter: 0.01
Recall: 0.5373134328358209
Recall: 0.6164383561643836
Recall: 0.6666666666666666
Recall: 0.6
Recall: 0.5
C parameter: 0.1
Recall: 0.5522388059701493
Recall: 0.6164383561643836
Recall: 0.7166666666666667
Recall: 0.6153846153846154
Recall: 0.5625
C parameter: 1
Recall: 0.5522388059701493
Recall: 0.6164383561643836
Recall: 0.7333333333333333
Recall: 0.6153846153846154
Recall: 0.575
C parameter: 10
Recall: 0.5522388059701493
Recall: 0.6164383561643836
Recall: 0.7333333333333333
Recall: 0.6153846153846154
Recall: 0.575
C parameter: 100
Recall: 0.5522388059701493
Recall: 0.6164383561643836
Recall: 0.7333333333333333
Recall: 0.6153846153846154
Recall: 0.575
Best model to choose from cross validation is with C parameter = 1
```

```

In [25]: 1 def plot_confusion_matrix(cm, classes, title):
2         """
3         Plot confusion matrix using matplotlib.
4         """
5         plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
6         plt.title(title)
7         plt.colorbar()
8         tick_marks = np.arange(len(classes))
9         plt.xticks(tick_marks, classes, rotation=0)
10        plt.yticks(tick_marks, classes)
11
12        thresh = cm.max() / 2.
13        for i, j in np.ndindex(cm.shape):
14            plt.text(j, i, format(cm[i, j], 'd'),
15                    horizontalalignment="center",
16                    color="white" if cm[i, j] > thresh else "black")
17
18        plt.tight_layout()
19        plt.ylabel('True label')
20        plt.xlabel('Predicted label')
21
22        # Fit logistic regression model using "best_c" as the regularization strength
23        lr = LogisticRegression(C=best_c, penalty='l2')
24        lr.fit(X_train, y_train.values.ravel())
25
26        # Make predictions on the test data
27        y_pred = lr.predict(X_test.values)
28
29        # Compute confusion matrix
30        cnf_matrix = confusion_matrix(y_test, y_pred)
31        np.set_printoptions(precision=2)
32
33        # Calculate recall metric
34        recall = cnf_matrix[1, 1] / (cnf_matrix[1, 0] + cnf_matrix[1, 1])
35        print("Recall metric in the testing dataset: ", recall)
36
37        # Plot non-normalized confusion matrix
38        class_names = [0, 1]
39        plt.figure()
40        plot_confusion_matrix(cnf_matrix, classes=class_names, title='Confusion matrix')
41        plt.show()

```

Recall metric in the testing dataset: 0.6190476190476191



Random Forest

```
In [31]: 1 import numpy as np
2 import pandas as pd
3 from sklearn.ensemble import RandomForestClassifier
4
5 # Create a random forest classifier
6 clf = RandomForestClassifier(n_estimators=100, random_state=42)
7
8 # Fit the classifier to the training data
9 clf.fit(X_train_undersample, y_train_undersample)
10
11 # Predict the target variable for the test data
12 y_pred = clf.predict(X_test)
13
14 # Convert the predicted target variable into a Pandas Series
15 y_pred = pd.Series(y_pred)
16
17 # Convert the test target variable into a Pandas Series
18 y_test = pd.Series(y_test.values.ravel())
19
20 # Compare the two Series element-wise
21 accuracy = np.mean(y_pred == y_test)
22 print('Accuracy:', accuracy)
23
```

Accuracy: 0.9746146553842913

This is so much better.

```
In [34]: 1 import xgboost as xgb
2
3 # Create an XGBoost model
4 clf = xgb.XGBClassifier(random_state=42)
5
6 # Fit the model to the training data
7 clf.fit(X_train, y_train)
8
9 # Predict the target variable for the test data
10 y_pred = clf.predict(X_test)
11
12 # Calculate the accuracy of the model
13 accuracy = clf.score(X_test, y_test)
14 print('Accuracy:', accuracy)
```

Accuracy: 0.9995318516437859

And This is even better!

Conclusion:

Our undersampled dataset performed well in our Logistic Regression models however, it will likely misclassify nonfraudulent transactions as fraudulent which will make our customers unhappy. Our Random Forest model did so much better. But the best model was our XGBoost. It seems to do much better with skewed data.

```
In [ ]: 1
```