Machine Learning Application Report

Overview

This section describes the various Machine Learning algorithms used and applied to the "Credit Card Fraud Detection" dataset.

Summary Files

I took a creditcards.csv file where Class is the target variable in there already prelabeled as 1 and 0 to easier classify as Fraud or Non-Fraud transaction and do not have to find them first. Dataset is ready for training a model but very imbalanced to make accurate classifications

The dataset that is used for credit card fraud detection is derived from the Kaggle https://www.kaggle.com/mlg-ulb/creditcardfraud The dataset has been collected and analysed during a research collaboration of Worldline and the Machine Learning Group (http://mlg.ulb.ac.be) of ULB (Université Libre de Bruxelles) on big data mining and fraud detection.

Feature Engineering

In this section I had to prepare all the data to make it readable and standardized for computations. I had data to be scaled and labeled.

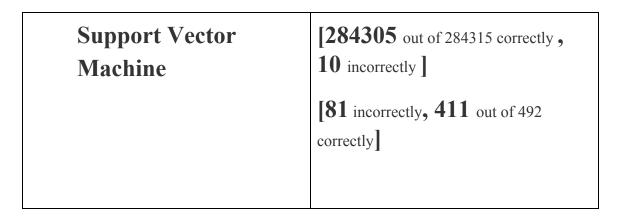
Training a model

There are many algorithms invented to train a classification and anomaly detection model. This will take a lot of time to test them all. That's why I tried only four, the most common:IsolationForest,Support Vector Machine,KNeighborsClassifier, Autoencoder.

Findings and conclusions

Model Algorithm	Confusion Matrix
IsolationForest:	[401, 283914]
	[175, 317]
Support Vector Machine	[284305, 10]
	[81,411
KNeighborsClassifier	[284315,0]
	[98,394]
Autoencoder	[284038,277]
	[352, 140

Let me explain the confusion matrix. Confusion Matrix tells me how many a Classifier classified correct values and incorrect out of all counts.



I justify this using this metric in my problem because my problem has very imbalanced data to detect and secondly, my problem is to detect correct fraud values, not to classify in general. Accuracy metric will not be trustworthy here. According to the confusion matrix, I have the best results in Support Vector Machine. I never know which one is the best for all cases. In one case, Isolation Forest will win, in another case, Autoencoder will win. There is no one Algorithm that will work perfectly for all cases. It depends on a dataset features it is made of and its properties. Obviously, the Support Vector Machine performed the best. Because it detected the highest number of all fraud transactions that I have in my dataset in total. If I have 492 fraud transactions in my dataset, SVM detects almost all of them 411. Other algorithms did a good job too but SVM did the best. Even if KNeighborsClassifier classified all Normal transactions correctly, it did worse on

Fraud. My goal is to detect all Fraudulent transactions instead. Of course, I want to use the Support Vector Machine algorithm in my case and problem. Support Vector Machine is designed for classification and regression problems. SVM is based on the idea of finding a hyperplane that best separates the features into different domains. So in my problem it works properly and the best of all. There are techniques that will fix imbalanced data problems: Undersampling and Oversampling. I can balance data artificially just for training models properly and achieve better results in classification accuracy. This is how I can improve my project.