

Data Science Career Track

Capstone 1 -

Milestone Report

by Edward Franke

04/17/2019

Credit Card Fraud Detection Project – Milestone Report

EXECUTIVE SUMMARY:

The purpose for this project is to find a correlation connected to fraudulent credit card transactions that can separate them in real time compared to legitimate transactions. The dataset is a cleaned dataset from Kaggle.com. It has been discovered there is a correlation connected to fraudulent transactions, but the code required to separate it from the legitimate transactions has not yet occurred. The most promising results comes from pandas profiling function which gives the following data:

Fraudulent Transactions (Class 1 indicated fraud)

Warnings

- Amount has 27 / 5.5% zeros Zeros
- Class has constant value 1 Rejected
- Row_Mean is highly correlated with Total ($p = 1$) Rejected
- Time is highly correlated with index ($p = 0.99465$) Rejected
- Total is highly correlated with V10 ($p = 0.94572$) Rejected
- V17 is highly correlated with V16 ($p = 0.96015$) Rejected
- V18 is highly correlated with V17 ($p = 0.97149$) Rejected
- V3 is highly correlated with V1 ($p = 0.90788$) Rejected

Legitimate Transactions (Class 0 indicates legit)

Warnings

- Class has constant value 0 Rejected
- Row_Mean is highly correlated with Total ($p = 1$) Rejected
- Time is highly correlated with index ($p = 0.99338$) Rejected

IDEA: A model to detect fraud in credit card transactions. (problem to solve)

CLIENT: Credit Card Companies

REASON: If they don't detect and stop the fraud as it happens, their customer don't pay the cost, they do. This project is intended to reduce their expenses and increase their profits.

DATA: From Kaggle, a cleaned dataset with 284,807 transactions from 2013 with 492 frauds in total.

SOLUTION: Create a model or analysis to discover what makes fraudulent transaction similar to detect this relationship as it happens.

DETAILS: See Detail Section Below.

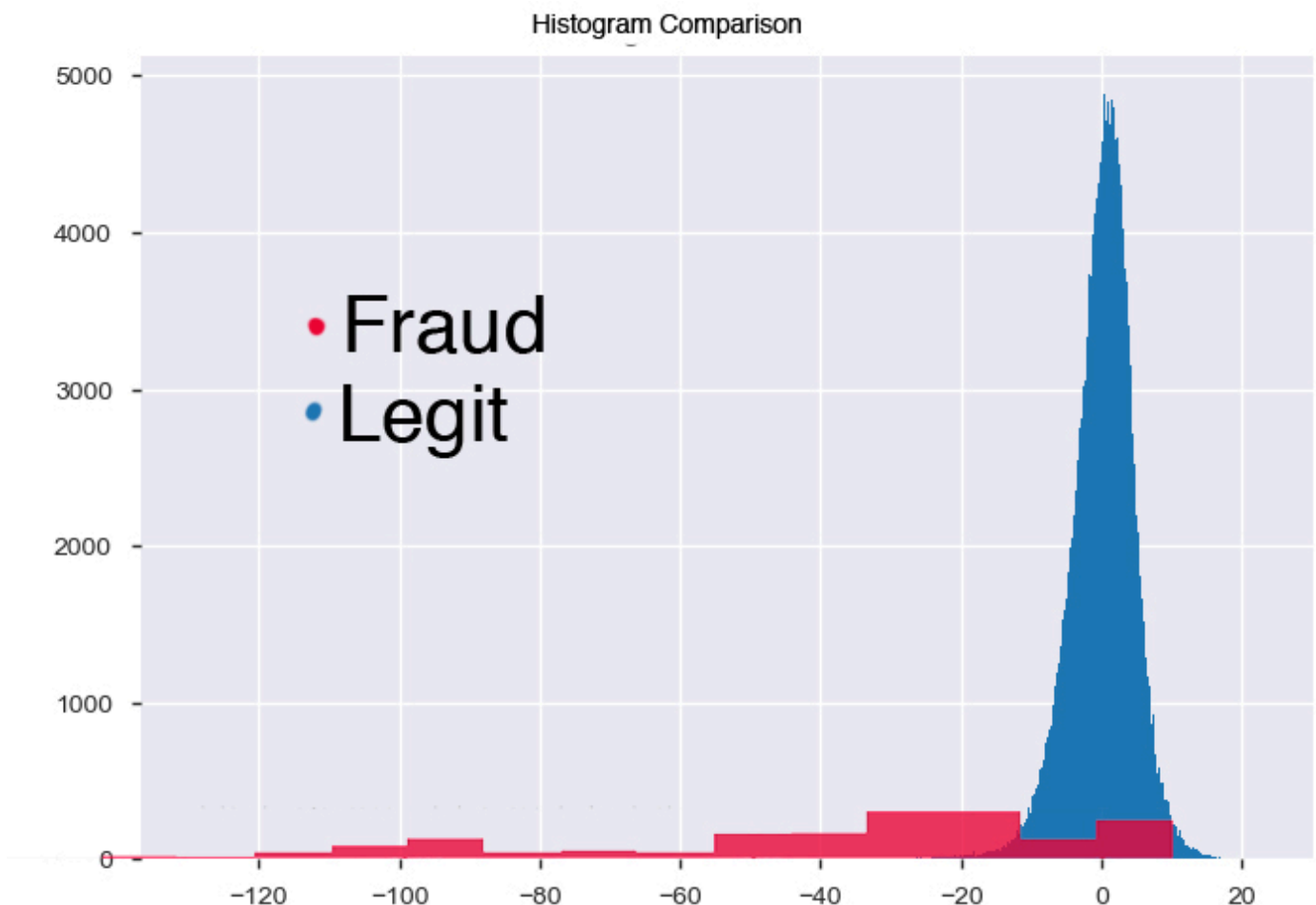
DELIVERABLES: Working code that detect fraud as it receives data and a presentation outlining the discoveries and explain the methods used to reach the discoveries.

Initial findings from exploratory analysis

```
1 data_fraud_df = data_df[data_df.Class == 1]
2 data_legit_df = data_df[data_df.Class == 0]
```

This is the code that was used to separate the initial dataframe into separate dataframes only containing fraud or legit transactions.

While comparing the transactions, I noticed fraudulent transactions have a greater histogram area that legitimate transactions. I was under the impression that I could create code to determine any transaction with a total V set as below -17 as fraudulent (see Appendix for more details). However, as this overlaid histogram comparison shows, that would fail the purpose of the project. A solution is still being sought.



WHAT'S NEXT

Deep Dive

- Discover what connects fraud transactions (V1 and V3, V16 and V17 and V18)

- Write code that detects this correlation as the transactions come in

- Run Statistical Analysis to determine error levels

Final Report and Presentation

APPENDIX – V Numbers compared

Fruadalent Transactions



Distinct count473

Unique (%)96.1%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-1.3133

Minimum-4.7719

Maximum2.1324

Zeros (%)0.0%

V3

Numeric

Distinct count473

Unique (%)96.1%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean4.542

Minimum-1.3133

Maximum12.115

Zeros (%)0.0%

V4

Numeric

Distinct count473

Unique (%)96.1%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-3.1512

Minimum-22.106

Maximum11.095

Zeros (%)0.0%

V5

Numeric

Distinct count473

Unique (%)96.1%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-1.3977

Minimum-6.4063

Maximum6.4741

Zeros (%)0.0%

V6

Numeric

Distinct count473

Unique (%)96.1%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-5.5687

Minimum-43.557

Maximum5.8025

Zeros (%)0.0%

V7

Numeric

Distinct count473

Unique (%)96.1%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean0.57064

Minimum-41.044

Maximum20.007

Zeros (%)0.0%

V8

Numeric

Distinct count473

Unique (%)96.1%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-2.5811

Minimum-13.434

Maximum3.3535

Zeros (%)0.0%

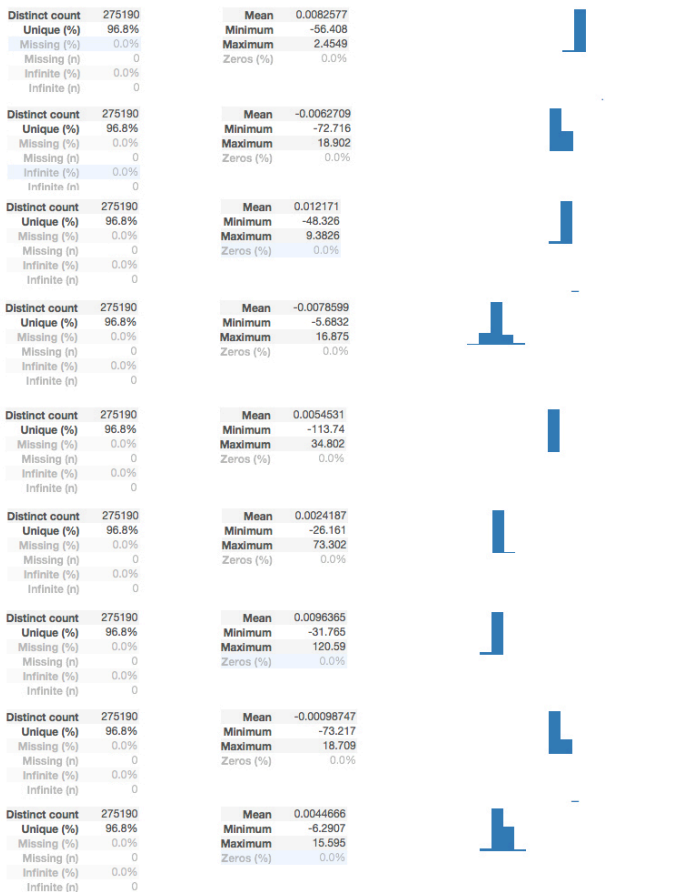
V9

Numeric

This variable is highly correlated with V1 and should be ignored for analysis

Correlation0.90788

Legitimate Transactions



Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-0.0062709

Minimum-72.716

Maximum18.902

Zeros (%)0.0%

V2

Numeric

Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean0.012171

Minimum-48.326

Maximum9.3826

Zeros (%)0.0%

V3

Numeric

Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-0.0078599

Minimum-5.6832

Maximum34.802

Zeros (%)0.0%

V4

Numeric

Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean0.0054531

Minimum-113.74

Maximum34.802

Zeros (%)0.0%

V5

Numeric

Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean0.0024187

Minimum-26.161

Maximum73.302

Zeros (%)0.0%

V6

Numeric

Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean0.0096365

Minimum-31.765

Maximum120.59

Zeros (%)0.0%

V7

Numeric

Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean-0.00098747

Minimum-73.217

Maximum18.709

Zeros (%)0.0%

V8

Numeric

Distinct count275190

Unique (%)96.8%

Missing (%)0.0%

Missing (n)0

Infinite (%)0.0%

Infinite (n)0

Mean0.0044666

Minimum-6.2907









Maximum15.595

Zeros (%)0.0%











V9

Numeric

Fruadalent Transactions

<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-5.6769</div><div>-24.588</div><div>4.0314</div><div>0.0%</div></div>	<div><div></div><div>V10</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>3.8002</div><div>-1.7022</div><div>12.019</div><div>0.0%</div></div>	<div><div></div><div>V11</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-6.2594</div><div>-18.684</div><div>1.3759</div><div>0.0%</div></div>	<div><div></div><div>V12</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-0.10933</div><div>-3.1278</div><div>2.8154</div><div>0.0%</div></div>	<div><div></div><div>V13</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-6.9717</div><div>-19.214</div><div>3.4424</div><div>0.0%</div></div>	<div><div></div><div>V14</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-0.092929</div><div>-4.4989</div><div>2.4714</div><div>0.0%</div></div>	<div><div></div><div>V15</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-4.1399</div><div>-14.13</div><div>3.1397</div><div>0.0%</div></div>	<div><div></div><div>V16</div><div>Numeric</div></div>
<div><div>This variable is highly correlated with</div><div>Correlation</div><div>0.96015</div></div> <div>V16 and should be ignored for analysis</div>		
<div><div>This variable is highly correlated with</div><div>Correlation</div><div>0.97149</div></div> <div>V17 and should be ignored for analysis</div>		
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>473</div><div>96.1%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.68066</div><div>-3.6819</div><div>5.2283</div><div>0.0%</div></div>	<div><div></div><div>V19</div><div>Numeric</div></div>

Legitimate Transactions

<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.0098237</div><div>-14.741</div><div>23.745</div><div>0.0%</div></div>	<div><div></div><div>V10</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-0.0065761</div><div>-4.7975</div><div>10.002</div><div>0.0%</div></div>	<div><div></div><div>V11</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.010832</div><div>-15.145</div><div>7.8484</div><div>0.0%</div></div>	<div><div></div><div>V12</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.0001892</div><div>-5.7919</div><div>7.1269</div><div>0.0%</div></div>	<div><div></div><div>V13</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.012064</div><div>-18.392</div><div>10.527</div><div>0.0%</div></div>	<div><div></div><div>V14</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.00016081</div><div>-4.3913</div><div>8.8777</div><div>0.0%</div></div>	<div><div></div><div>V15</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.0071641</div><div>-10.116</div><div>17.315</div><div>0.0%</div></div>	<div><div></div><div>V16</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.011535</div><div>-17.098</div><div>9.2535</div><div>0.0%</div></div>	<div><div></div><div>V17</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>0.0038872</div><div>-5.3667</div><div>5.0411</div><div>0.0%</div></div>	<div><div></div><div>V18</div><div>Numeric</div></div>
<div><div>Distinct count</div><div>Unique (%)</div><div>Missing (%)</div><div>Missing (n)</div><div>Infinite (%)</div><div>Infinite (n)</div></div> <div><div>275190</div><div>96.8%</div><div>0.0%</div><div>0</div><div>0.0%</div><div>0</div></div>	<div><div>Mean</div><div>Minimum</div><div>Maximum</div><div>Zeros (%)</div></div> <div><div>-0.0011779</div><div>-7.2135</div><div>5.592</div><div>0.0%</div></div>	<div><div></div><div>V19</div><div>Numeric</div></div>

Fruadalent Transactions



Legitimate Transactions



APPENDIX – Kaggle Details about the dataset.

Context

It is important that credit card companies are able to recognize fraudulent credit card transactions so that customers are not charged for items that they did not purchase.

Content

The datasets contains transactions made by credit cards in September 2013 by european cardholders. This dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions. The dataset is highly unbalanced, the positive class (frauds) account for 0.172% of all transactions.

It contains only numerical input variables which are the result of a PCA transformation. Unfortunately, due to confidentiality issues, we cannot provide the original features and more background information about the data. Features V1, V2, ... V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'. Feature 'Time' contains the seconds elapsed between each transaction and the first transaction in the dataset. The feature 'Amount' is the transaction Amount, this feature can be used for example-dependent cost-sensitive learning. Feature 'Class' is the response variable and it takes value 1 in case of fraud and 0 otherwise.

Inspiration

Identify fraudulent credit card transactions.

Given the class imbalance ratio, we recommend measuring the accuracy using the Area Under the Precision-Recall Curve (AUPRC). Confusion matrix accuracy is not meaningful for unbalanced classification.

Acknowledgements

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. More details on current and past projects on related topics are available on <http://mlg.ulb.ac.be/BruFence> and <http://mlg.ulb.ac.be/ARTML>

Please cite: Andrea Dal Pozzolo, Olivier Caelen, Reid A. Johnson and Gianluca Bontempi. Calibrating Probability with Undersampling for Unbalanced Classification. In Symposium on Computational Intelligence and Data Mining (CIDM), IEEE, 2015