



॥ सा विद्या या विमुक्तये ॥
भारतीय प्रौद्योगिकी संस्थान धारवाड
Indian Institute of Technology Dharwad

Credit Card Fraud Detection

K. Sai Anuroop, Mandeep Bawa, Sushma Biradar, Aniruddha Joshi

Indian Institute of Technology Dharwad

Computer Science and Engineering

[170030035, 170030038, 170010032, 170020004] @iitdh.ac.in



॥ सा विद्या या विमुक्तये ॥
भारतीय प्रौद्योगिकी संस्थान धारवाड
Indian Institute of Technology Dharwad

Abstract

Credit card fraud can be defined as, *'Unauthorized account activity involving a payment card, by a person for which the account is not intended'*.

These frauds cost consumers and banks millions of dollars worldwide, as a response to which several modern fraud-detection techniques are in place today.

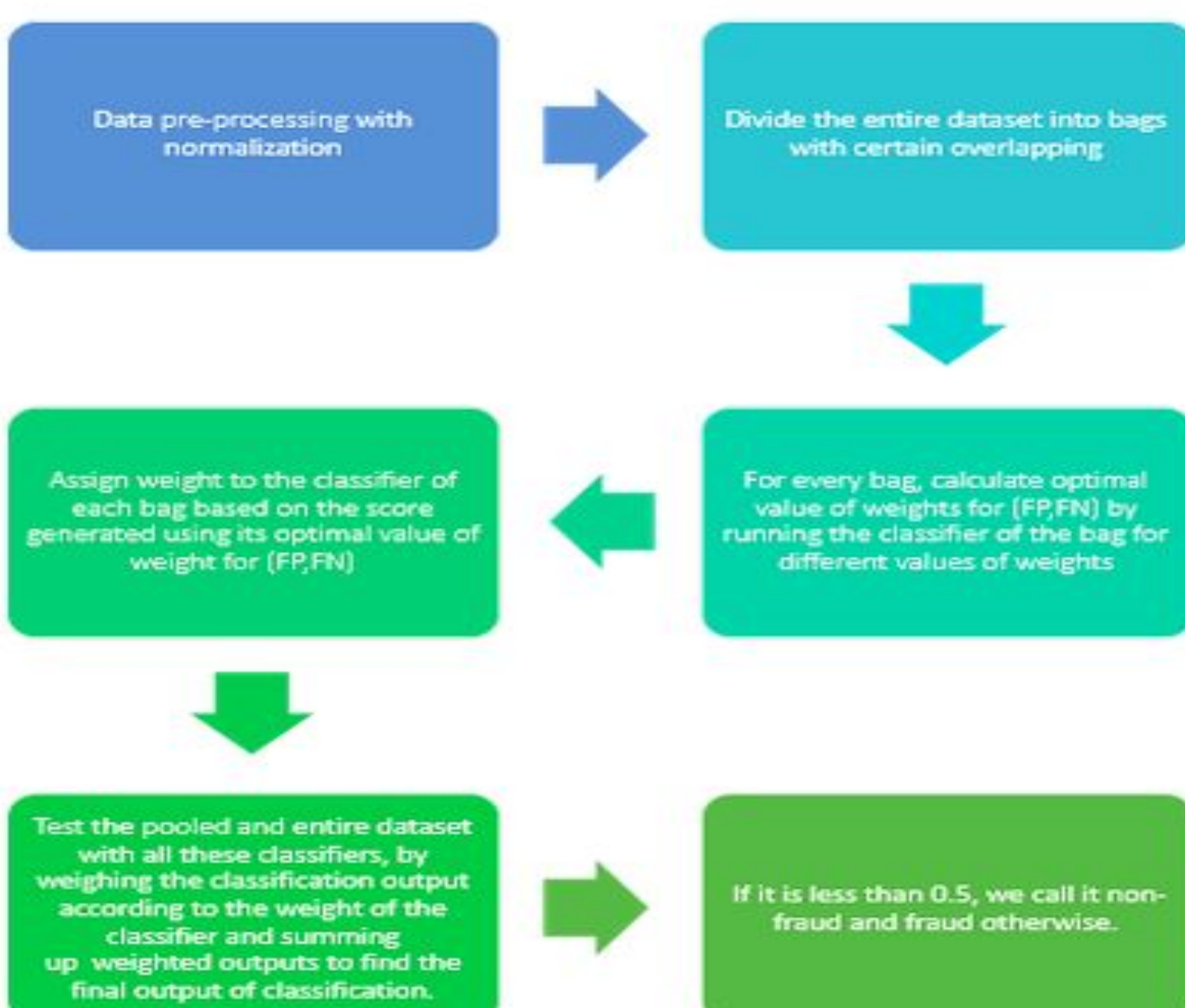
Here, we brief about the algorithm we implemented for detecting fraudulent transactions.

Introduction

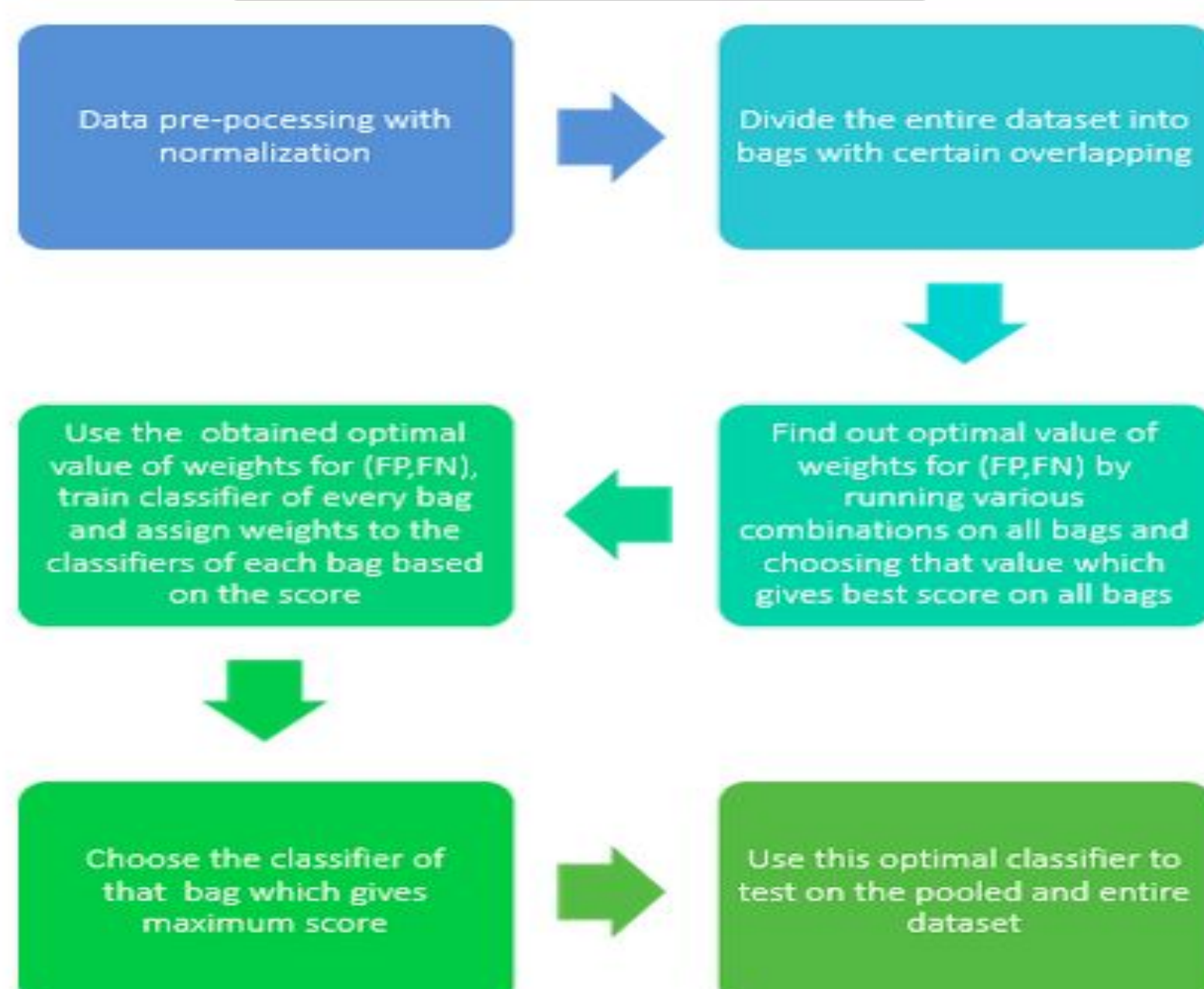
We based our work on the dataset provided in Kaggle under Open Database license. This dataset is already normalized but is highly biased towards Non-Fraud transactions, as is expected of a legal nature of transactions across the world. Feeding this raw data to our algorithm will highly affect its results. So, data pre-processing is a must to ensure correctness and validity of results. For addressing this issue, we have implemented the bagging phase. Since we are concerned with classification of feature vectors, we implemented SVM with Gaussian and Polynomial kernel functions. Also, we have explored two different strategies for training our classifier.

Block Diagram

Strategy #1



Strategy #2



Proposed Algorithm

We decompose the proposed algorithm into the following parts:

- Bagging of feature vectors
- Assigning weights to the classifier associated with every bag
- Obtaining the model for classification
- Testing the dataset with obtained classifier

Proposed Modeling Scheme

Bagging Phase

Random selections tend to retain data proportions, hence the constituent imbalance levels in training data are carried forward to the base learners. This leads to data imbalance affecting the training process to a large extent. The proposed model enables balanced data selection such that the effects of data imbalance are considerably reduced during model training. Let, T_{min} to be the set of Fraudulent Transactions and T_{major} to be Legitimate Transactions. T_{b_i} be the b^{th} bag containing feature vectors on which base learner is trained. Then define,

$$T_{b_i} = T_{major} \cup T_{min}$$

where T_{major} defines sampled instances of T_{major} , and is created by performing n overlapping divisions of T_{major} . Instances for T_{major} are obtained by sampling the data from T_{major} within the interval $[(b-1)NT+1, (b)NT]$ where $1 \leq b \leq n$ is the bag identifier.

$$NT = |T_{major}| = (|T_{major}|/n) + \theta * |T_{major}|$$

$0 \leq \theta \leq 1$ is a hyper-parameter that defines the degree of accepted overlap among majority classes.

Consecutive bags contain certain levels of overlaps to make sure that the temporal distribution change is gradual and hence does not exhibit sudden changes in predictions between consecutive bags.

Training Phase

Strategy #1

We vary the values of weights given to FPs and FNs to find the optimal values of these at which a bag gives maximum sum of accuracy, precision and recall. We do this for all the bags and obtain the classifier of the bag at the optimal value. Next, we assign weights to the classifiers obtained, on the basis of the sum of accuracy, precision and recall at their optimal FP and FN weights. We then test the pooled and entire datasets with all these classifiers, by weighing the classification output according to the weight of the classifier and summing up weighted outputs to find the final output of classification. If it is less than 0.5, we call it non-fraud and fraud otherwise.

Strategy #2

We vary the values of weights given to FPs and FNs to find the optimal values of these at which the sum of accuracy, precision and recall of all the bags is maximum. Next, at the optimal value of weights for FPs and FNs obtained, we train the classifiers of every bag and find that classifier which gives the maximum sum of accuracy, precision and recall. We then use this classifier to test on the entire dataset. Using polynomial kernel performed better than Gaussian kernel.

Testing Phase

We merged data used for testing different bags to create one big pooled testing dataset. This was achieved by pooling 20% of the testing dataset of each bag. Also, we tested on the entire dataset.

Results

Strategy #1

Pooled Dataset

Accuracy : 97.32%
Precision : 97.93%
Recall : 96.19%

Entire Dataset

Accuracy : 97.33%
Precision : 97.93%
Recall : 96.14%

Strategy #2

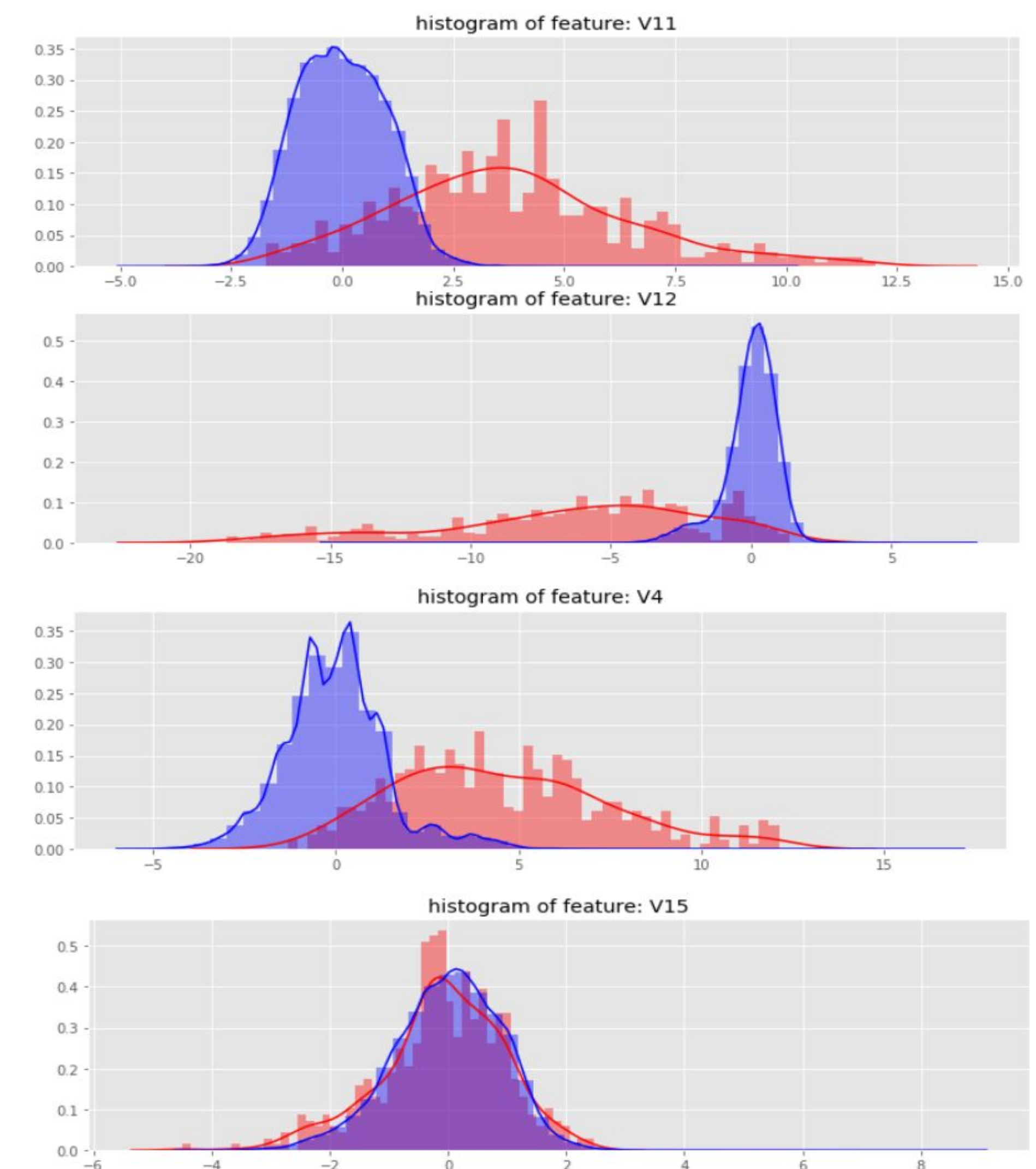
Pooled Dataset

Accuracy : 87.18%
Precision : 81.20%
Recall : 93.36%

Entire Dataset

Accuracy : 99.92%
Precision : 87.34%
Recall : 69.69%

Distribution of FVs



Conclusion

We reviewed different algorithms for credit card fraud detection starting with logistic regression, SVM and then came up with our own algorithm, fusing different techniques employed in the literature reviewed.

In our tests, Strategy #1 performed better on the pooled dataset when compared to Strategy #2. Also, it was noticed that Strategy #2 gave better performance when tested over entire dataset.

It is observed that adding weights based on false positives and false negatives to SVM classifier improves performance.

While training, it is observed that the method of bagging enabled balanced data selection such that the effects of data imbalance are considerably reduced during model training.

We thus proposed and implemented an algorithm for credit card fraud detection which yields good classification results.

Future Work

Future work:

For real-time fraud detection, bags containing recent transactions could be given more weight when compared to bags containing older transactions.

Acknowledgements

We would like to thank Prof. SRM Prasanna for his guidance all throughout the project.

References

- [1] Parallel and incremental credit card fraud detection model to handle concept drift and data imbalance
Somasundaram, A. & Reddy, S.
Neural Computing and Applications, Springer, (2019) 31(Suppl 1): 3
<https://doi.org/10.1007/s00521-018-3633-8>
- [2] Fraud Detection using Machine Learning
Aditya Oza, Stanford University <http://cs229.stanford.edu/proj2018/report/261.pdf>