

Fraud Detection with Confidence

A Benchmarking Case Study on Synthetic Data

Daniel Turner-Szymkiewicz, Prof. Ulf Norinder,
Dr. Mirosława Alunowska Figueroa, Dr. Edgar Lopez Rojas

EalaX Ltd
edgarlopez@ealax.com
danielt@ealax.com

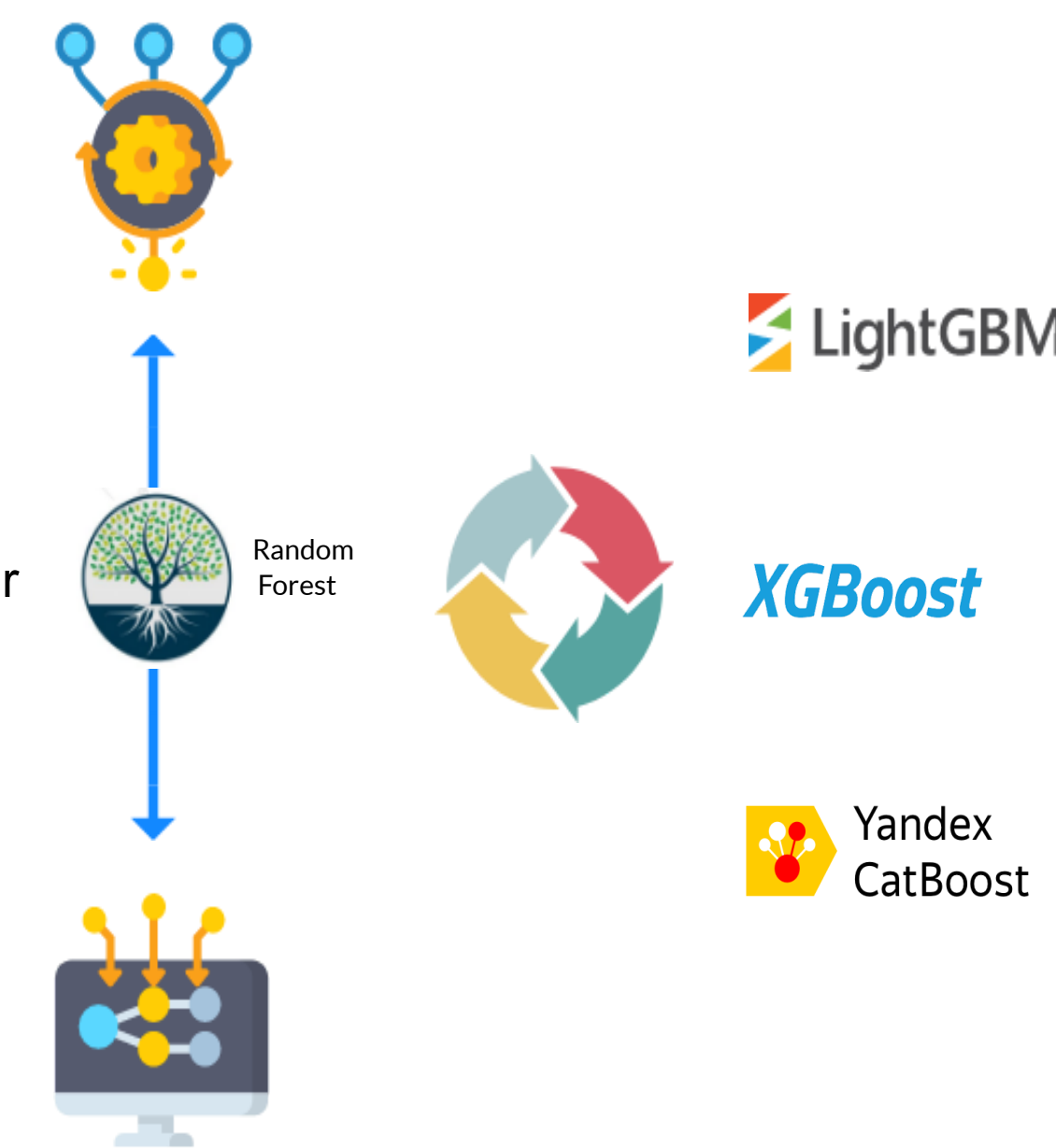
Methodology

Development of a supervised Machine Learning model for fraud detection required an experimental design which would be clear and explainable when running diagnostics & replicated a real-time environment of fraud detection

End Goal Criteria

- Maximize number of found fraudulent transactions
- Maximize amounts of found fraudulent transactions
- Minimize false positives (normal transactions identified as fraud)

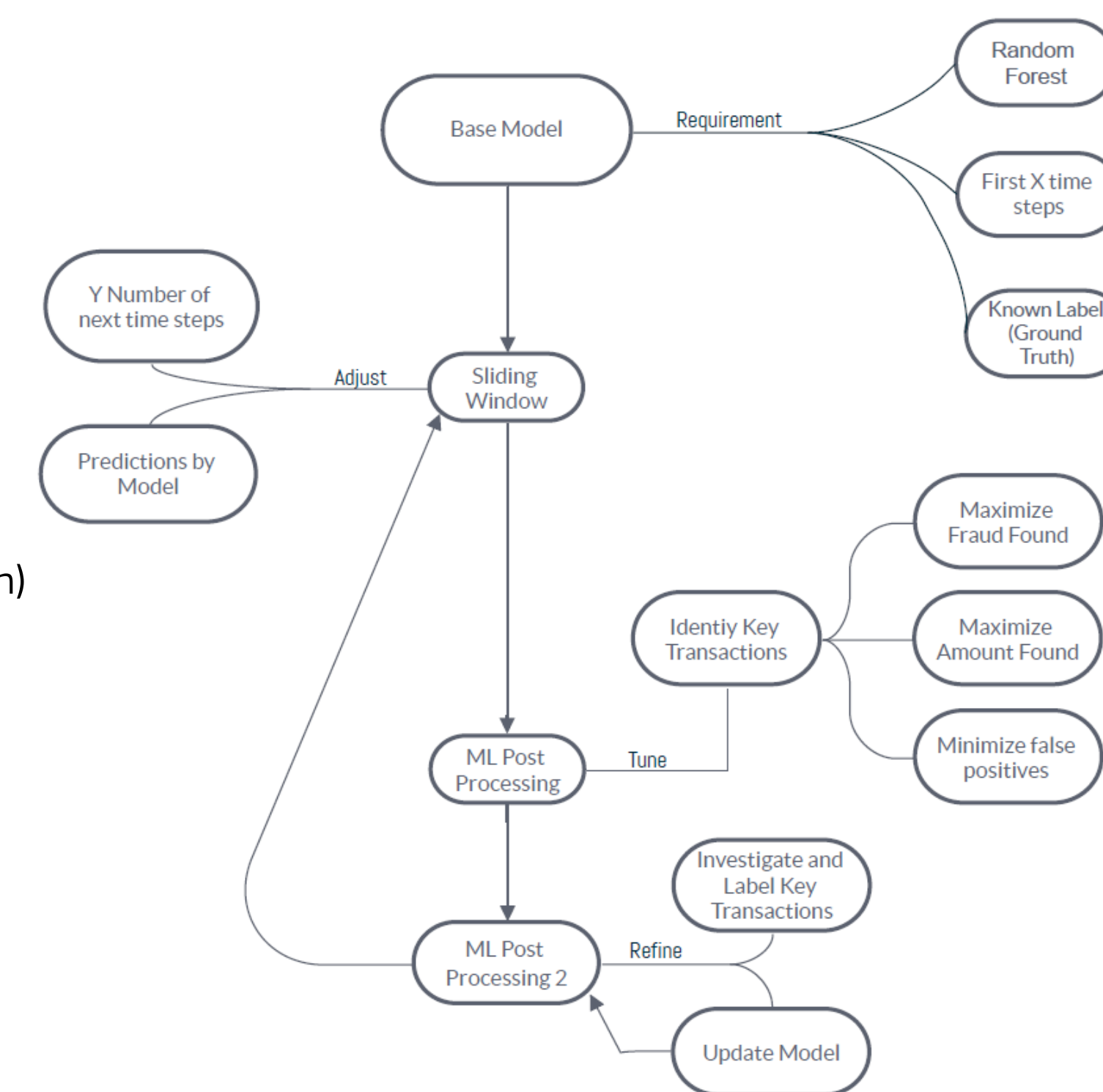
Machine Learning using a base classifier



- Random Forest is Robust and Explainable
- Good performance throughout different domains
- Seldom need for hyper-parameter optimization
- Similar classifiers can be used for comparison (i.e. Boost / LGBM)

Figure: Case-Study Design

Time series analysis



- Base model from first X number of time steps with known labels (ground truth)
- Sliding window (Y number of next time steps) where "identified" transactions are investigated and labels assigned.
- Updating of model using "identified" transactions
- Prediction of next sliding window (Y number of next time steps) where "identified" transactions are investigated and labels assigned.
- Continue Steps 2 and 3

Introduction

We believe that the use of adequately bench-marked synthetic data can help mitigate a great part of the current problems in the financial crime domain, and thus reduce concerns for accuracy, bias and privacy in the realm of compliance (van Driel, 2019). The Conformal Predictions (CP) framework is a recent development in machine learning to associate reliable measures of confidence with pattern recognition settings and is an excellent solution to bridge the gap between machine learning and the validity of synthetic data (Vovk et al., 2005).

- Synthetic data in finance solves the problem posed by the difficulty to access financial datasets due to privacy regulation (Barse, Kvarnstrom, Jonsson, 2004).
- In finance its use has gained significant interest; synthetic data has been adopted to improve financial crime detection and compliance (Lopez-Rojas, Axelsson, and Baca, 2018).
- Criminals' strategies modify constantly over time, and control systems could potentially adapt to this reality by using synthetic data that explores plausible criminal behaviour (Karpoff, 2020).
- A key consideration is if the synthetic data is effective a ML model built from synthetic data performs as well as models built from real data.



- 1 Secure Data**
Establish source of original data.
- 2 Generate Synthetic Data**
Data is generated via historical notation or fraud typology injection.
- 3 Benchmark Model Training**
Multiple secure synthetic datasets are generated, and benchmarked.
- 4 Model Deployment**
An ideal candidate can be selected with confidence.

Conformal Approach

For a binary classification task, such as fraud detection, the CP framework operates as follows:

- The Training data set is randomly split into a proper training set (70 %) and a calibration set (30 %).
- The model is built on the proper training set and the calibration set as well as the new test data set are predicted from the model.
- The test set predictions are compared to the calibration set predictions of the two classes separately and the CP p-values, one value for each of the two classes, are calculated for every test set example.
- Comparing these p-values for each of the two classes (fraud, normal) separately against the acceptable error rate determined by the user provides the final classification (label, labels or no labels) for every new example (Vovk et al., 2005)

A Comparative Analysis of the Conformal Framework

Our case study focuses on looking at the viability of our framework when combined with tried and tested gold standard algorithms for the analysis of financial crime data.

This includes observing different time windows and comparing these to classical approaches some financial institutions still implements (such as dynamic thresholding) and seeing how our conformal approach compares.

Below is the sum output of each algorithms performance in the conformal approach within a 5 hour sliding-window and its respective comparison with the thresholding approach (th) that utilizes the same dynamic Time window but does not use the conformal approach.

| PaySim Comparison of Algorithms | | | | | |
|---------------------------------|---------------|------------|-------------|-------------|------------|
| Type | Random Forest | XGBoost | LGBM | CatBoost | Threshold |
| Fraud Found | 7499 | 7102 | 7252 | 6723 | 5465 |
| Total Fraud | 8169 | 8169 | 8169 | 8169 | 8169 |
| Missed Fraud | 670 | 1067 | 917 | 1446 | 2704 |
| Missclassified Fraud | 16179 | 7410 | 26197 | 107631 | 1541793 |
| Missed Fraud (M) | 643294435 | 2160258449 | 1923317584 | 795973769 | 379836668 |
| Missclassified Fraud (M) | 5684855488 | 2719956851 | 10107092423 | 28011401527 | 6.9315E+11 |
| F1 | 0.8675 | 0.8985 | 0.8137 | 0.4667 | 0.222512 |
| Recall | 0.9191 | 0.8707 | 0.8847 | 0.8309 | 0.66771 |
| Precision | 0.8869 | 0.9474 | 0.8541 | 0.4403 | 0.227023 |
| K | 0.8503 | 0.899 | 0.7926 | 0.3886 | 0.093066 |
| Specificity | 0.9965 | 0.9989 | 0.9962 | 0.9814 | 0.818034 |
| MCC | 0.8672 | 0.9079 | 0.817 | 0.4555 | 0.138441 |

Figure: An overview of classifier performance on PaySim



A visualisation of the conformal approach with Random Forest versus the classical thresholding approach

Conclusion

- Within this data our framework really shines in being able to isolate fraudulent signals and avoids the pitfalls of false positives.
- The conformal approach has shown great promise in its application on financial fraud data.
- Throughout this case study we have observed fascinating results from the conformal approach, yet it is important to mention that it can serve one other key function, and that is to bridge the gap of data explainability.
- Therefore, we would also like to explore the use-case of synthetic data playing a major role in creating more robust ML models that may avoid the pitfalls of bias by incorporating better information.
- In close, we believe that the conformal approach will be essential in the validation and improvement of synthetic data generation

Acknowledgements

This work was supported by UK's innovation agency, Innovate UK, under granted projects to EalaX Ltd: FraudSim 82929 and CP-Mark 89039 - 2020/2021.



Barse E.L., Kvarnstrom H., Jonsson E., 2004. Synthesizing test data for fraud detection systems Conference Paper Conference: Computer Security Applications Conference, 2003. Proceedings, 19th Annual IEEE Xplore

Karpoff Jonathan M. 2020 "The future of financial fraud." Journal of Corporate Finance, pp. 10169

Lopez-Rojas E. A., Axelsson S., and Baca D. 2018. "Analysis of Fraud Controls Using the PaySim Financial Simulator." International Journal of Simulation and Process Modelling, 13 (4), pp. 377-386, ISBN: 1740-2131.

Lopez-Rojas, Edgar Alonso & Elmi, Ahmad & Axelsson, Stefan. (2016). PAYSIM: A FINANCIAL MOBILE MONEY SIMULATOR FOR FRAUD DETECTION.

van Driel Hugo. 2019. "Financial fraud, scandals, and regulation: A conceptual framework and literature review." Business History, 61(8):1259-1299, 2019. doi: 10.1080/00076791.2018.1519026.

V. Vovk, A. Gammerman, G. Shafer, Algorithmic Learning in a Random World, Springer, New York, 2005, pp. 1-324.