# Comparative Evaluation of Anomaly Detection Methods for Fraud Detection in Online Credit Card Payments

Hugo Thimonier[1], Fabrice Popineau[1], Arpad Rimmel[1], Bich-Liên Doan[1], and Fabrice Daniel[2]

[1] Université Paris-Saclay, CNRS, CentraleSupélec,
Laboratoire Interdisciplinaire des Sciences du Numérique,
91190, Gif-sur-Yvette, France
`name.surname@lisn.fr`,
[2] LUSIS AI, Paris, France

**Abstract.** This study explores the application of anomaly detection (AD) methods in imbalanced learning tasks, focusing on fraud detection using real online credit card payment data. We assess the performance of several recent AD methods and compare their effectiveness against standard supervised learning methods. Offering evidence of distribution shift within our dataset, we analyze its impact on the tested models' performances. Our findings reveal that LightGBM exhibits significantly superior performance across all evaluated metrics but suffers more from distribution shifts than AD methods. Furthermore, our investigation reveals that LightGBM also captures the majority of frauds detected by AD methods. This observation challenges the potential benefits of ensemble methods to combine supervised, and AD approaches to enhance performance. In summary, this research provides practical insights into the utility of these techniques in real-world scenarios, showing LightGBM's superiority in fraud detection while highlighting challenges related to distribution shifts.

**Keywords:** Imbalanced Learning, Anomaly Detection, Fraud Detection

## 1 Introduction

Detecting fraudulent behaviors has emerged as a critical problem that has garnered significant attention from practitioners and scholars alike. In sectors such as banking, frauds incurred an estimated annual cost of \$28.58 billion in 2021, as highlighted by the Nilson Report 2021[3]. To address the challenge of identifying frauds within regular credit card payments, banks have increasingly turned to machine learning techniques, known for their effectiveness in many classification tasks, particularly with unstructured data.

Two critical features of fraud detection pose challenges to constructing effective and accurate classifiers: highly imbalanced classes and distribution shift.

---

[3] https://nilsonreport.com/

Imbalanced datasets arise due to the significant disparity in the number of genuine transactions compared to fraudulent ones, making it difficult for traditional classification algorithms to generalize accurately. Highly imbalanced datasets are a specific subset of imbalanced datasets in which the positive class represents less than 1% of the samples; this situation is also referred to as rarity [3]. Moreover, distribution shift occurs as fraudsters constantly adapt their strategies, causing a discrepancy between the training and testing data distributions, thereby hampering the performance of machine learning models.

Learning from imbalanced datasets is a critical topic with implications across various real-life applications. Extensive research has highlighted the consequences of imbalanced learning on canonical classifiers, revealing that most standard classifiers are ill-suited for imbalanced settings. For instance, the limitations of standard machine learning techniques when confronted with imbalanced datasets are examined comprehensively by [42]. This study also sheds light on the struggles faced by backpropagation algorithms in converging within imbalanced setups, as the dominant majority class can overwhelm the gradient vector used for weight updates in neural networks. In contrast, Gradient Boosted Decision Trees (GBDT) are often considered more resilient to imbalanced settings due to their focus on particularly challenging examples [11], thus enabling them to prioritize the minority class more effectively. Highly imbalanced datasets are among the most challenging as research [21] has shown how learners display decreasing performance as imbalance becomes more severe.

In addition to imbalanced datasets, distribution shift is another crucial challenge in detecting fraud. Standard machine learning techniques perform well when the training and testing datasets distributions are similar, if not identical. However, domain shift occurs when the distribution of the test dataset deviates from the original training distribution and thus hinders standard classifiers' performance. Fraud detection involves an iterative game between fraudsters and banks. Fraudsters continually strive to produce increasingly inconspicuous fraudulent behaviors, while banks aim to detect frauds as accurately as possible while avoiding false negatives. This dynamic nature of fraud detection presents significant challenges for machine learning algorithms.

These characteristics of fraud detection underscore the need for methodologies capable of effectively handling distribution shift and highly imbalanced datasets. In response to these challenges, researchers have proposed using anomaly detection (AD) methods, which promise to exhibit robustness in case of both distribution shift and extreme class imbalances. Specifically, AD involves the identification of anomalies within a dataset by delineating deviations from a predefined notion of normality. Anomaly detection methods typically characterize the normal distribution solely based on normal samples during training. Consequently, AD has been regarded as particularly well-suited for imbalanced and extremely imbalanced settings. By design, AD methods do not experience performance deterioration when faced with highly skewed class distributions, as the training process solely requires normal samples. Moreover, assuming only fraudulent behaviors change over time, AD models should be more robust to

distribution shift than standard supervised approaches. Indeed, if the normal distribution is well characterized, they should always be able to exclude new types of anomalies.

In this work, we empirically investigated AD for fraud detection tasks, exploring their capabilities and limitations. By empirically evaluating various AD methods on a real-world dataset characterized by distribution shifts and extreme class imbalances, we aim to provide insights into the suitability and effectiveness of these techniques for addressing the challenges inherent to fraud detection. In addition to evaluating the performance of AD methods, we conduct a comparative analysis with Gradient Boosted Decision Trees (GBDT), the prevalent choice for machine learning tasks on tabular data [16], to gauge the added value of AD approaches. We rely on the LightGBM implementation [23] and show that GBDT suffer significantly more from distribution shift than AD methods while displaying substantially better fraud detection performance than all tested AD methods. Our paper is structured as follows: in the next section, we discuss works related to our application; in section 3 we discuss in detail the experiments conducted on our dataset; in section 4 we present the obtained results; in section 5 we discuss the results and finally in section 6 we conclude.

## 2 Related Works

Anomaly detection encompasses two types of algorithms: supervised and unsupervised. In the case of supervised AD, one disposes of the label and indirectly uses it in the training process by building a training set solely composed of samples belonging to the *normal* class[4]. Unsupervised AD methods involve situations where the label is unavailable, and anomalies must be directly identified within a dataset containing both *normal* samples and anomalies. While supervised approaches can be used in situations where the imbalance is too severe for standard supervised approaches to work, unsupervised approaches are usually confined to applications that consist in removing samples that may hinder another models' performance on a particular task, e.g. mislabeled samples or outliers. Since we dispose of labels in the context of fraud detection, we will focus on supervised anomaly detection.

Supervised anomaly detection methods differ from standard supervised approaches because labels are only used indirectly. Indeed, standard supervised approaches consist in training a classifier using a dataset

$$\mathcal{D}_{train} = \{(x_i, y_i) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}_{i=1}^n, \tag{1}$$

using both sample features $x_i$ and labels $y_i$. Moreover, $\mathcal{D}_{train}$ contains samples from each class in $\mathcal{Y}$. On the contrary, supervised anomaly detection methods use the label to build a training set solely composed of a single class, referred to as the *normal* class. In the case of binary classification, the *normal* class is the

---

[4]Throughout this paper, the term *normal* relates to the notion of normality.

majority class, e.g. the legit payments in the case of fraud detection, and the training set can then be constructed as

$$\mathcal{D}_{train}^{AD} = \{x_i : y_i = 0\}_{i=1}^n. \tag{2}$$

In this anomaly detection framework, the overall goal is to characterize the normal distribution, $p(x \mid y = 0)$. In inference, this characterization is used to determine whether a sample belongs to the normal distribution or should be seen as an anomaly.

As a field of research, anomaly detection can be divided into several non-exhaustive categories: one-class classification (OCC), reconstruction-based methods and self-supervised methods.

*One-Class Classification* In contrast to traditional machine learning classification problems, one-class classification (OCC) approaches aim to identify samples that do not belong to a specific class by characterizing the distribution of that class. These discriminative models learn a decision boundary using only samples from the designated *normal* class, thereby circumventing the direct estimation of the class distribution. During the inference phase, samples are classified as either belonging to the *normal* class or not, without making any assumptions about the *anomaly* class. One-class support vector machines (OCSVM) [35] and support vector data description (SVDD) [39] are popular OCC methods that rely on kernels to map the data space to a Hilbert space, where a decision boundary is learned. Recently, [33, 34] have introduced extensions to OCC methods that incorporate deep neural networks to alleviate the computational complexity associated with kernels. Other OCC tree-based approaches can be found in the OCC such as isolation forest (IForest) [27], extended isolation forest [19], Robust Random Cut Forest (RRCF) [18] and PIDForest [13]. Other methods have relied on sample-sample dependencies to identify anomalies such as TracInAD [41] relying on influence measures or approaches based on k-nearest neighbors (KNN). In the latter, anomalies are identified by measuring the distance of each sample to its k-nearest neighbors [2, 32]: higher distance indicating abnormality.

*Reconstruction-based methods* Reconstruction-based anomaly detection methods rely on the assumption that different distributions generate normal samples and anomalies. Consequently, training a model to reconstruct samples from the *normal* distribution aims to achieve low reconstruction error for any sample belonging to this distribution. Conversely, anomalies that are believed to stem from a distinct distribution should exhibit significantly higher reconstruction errors.

One of the most prevalent shallow reconstruction-based anomaly detection methods employ Principal Component Analysis (PCA) or Bayesian PCA [9, 20]. Autoencoders [10], regularized autoencoders like Variational Autoencoders (VAEs) [30] and memory-augmented deep autoencoders [12] have also been leveraged for anomaly detection. Recently, [24] proposed a novel methodology for anomaly detection using autoencoders that incorporates the hidden representations of the original and reconstructed samples. Instead of solely comparing

the reconstructed sample and the original sample, the authors suggest comparing the hidden representations of both samples by passing the reconstructed sample through the autoencoder. In addition, recent approaches have explored attention-based architectures for reconstructing masked features of samples, as exemplified by NPT-AD [40]. Other related methods do not directly compute a reconstruction error and only focus on estimating either the entire *normal* distribution such as ECOD [26] or local *normal* distributions as proposed in local outlier factor (LOF) [6].

*Self-supervised methods* The literature also features self-supervised approaches employing pretext tasks for anomaly detection [4, 31, 38]. In GOAD [4], several affine transformations are applied to each sample in the training set, while a classifier is trained to predict the specific transformation applied to a transformed sample. During testing, since the classifier was exclusively trained on *normal* samples, it is expected to struggle in correctly predicting the transformation for anomaly samples. Similarly, [31] propose NeuTraL-AD, a contrastive framework in which they transform samples using neural mappings instead of affine transformations. The objective is to learn transformations that maintain similarities in a semantic space between transformed samples and their untransformed counterparts while different transformations are easily distinguishable. In inference, the contrastive loss utilized to optimize the parameters serves as the anomaly score. More recently, [36] introduced a self-supervised methodology for anomaly detection that maximizes the mutual information among the elements of a sample's features using contrastive learning. By maximizing the mutual information, the method effectively captures the underlying structure of normal samples and identifies deviations indicative of anomalies.

*Supervised classification on tabular data* Although deep-learning models have become ubiquitous for a broad range of tasks involving natural language processing (NLP) and computer vision (CV), applying these models to tabular data remains very challenging. Some recent methods [15, 22, 37] have shown promising results when applying deep learning models tailored for tabular data. However, in recent work [14, 16], authors discuss how neural networks tend to struggle with this data type in comparison with other methods based on gradient-boosted decision trees (GBDT). In most scenarios, approaches such as XGBoost [7] or LightGBM [23] have been shown to surpass deep learning algorithms. This type of approach remains the go-to method for practitioners due to its strong classification performance and its simplicity to train in comparison with deep methods. Moreover, GBDT models such as LightGBM and XGBoost are often considered particularly suited for imbalanced and extremely imbalanced set-ups since these models focus on particularly hard-to-classify samples [11], generally the minority class, and thus offer strong performance in comparison to other standard machine learning models.

## 3 Experiments and Datasets

### 3.1 Dataset

We dispose of a labeled dataset of online credit card payments made available by a major French bank. Our dataset contains 145 features describing raw characteristics of payments (e.g. amount, currency) as well as features that we computed, such as rolling sums or rolling means. Our dataset contains 480 million online transactions from the first day of 2018 until the last day of 2021. The dataset is comprised of two independent datasets merged, one from 2018 to 2019 and the other from 2020 to 2021. This dataset displays the characteristic of highly imbalanced classes[5] discussed in section 1 since the proportion of legit payments vastly outnumbers the proportion of frauds. We remove cards with less than 50 payments, cards for which the proportion of frauds exceeds 50%, and cards with too few payments since they would not have derived features (e.g. rolling means) with meaningful values and would risk hindering the models' performance. Similarly, we also omit cards with too many frauds in their payment history since they would also be problematic as they might pollute the fraud distribution. Overall, this preprocessing reduces the dataset to 192 million payments. Most methods we wish to test would be intractable with such dataset size and require further dimension reduction. Thus, we restrict our analysis to two countries with 3 million payments (1.5% of total dataset size) and 20 million payments (10.3% of total dataset size), respectively. Moreover, restricting our analysis to one country at a time should help models learn since payment distributions likely differ between countries.

**Distribution shift** We argue that our datasets undergo a distribution shift and that the distributions of legit and fraudulent payments differ between 2018 and the end of 2021. For instance, as supported by [17] Covid-19 has caused online payment behaviors to change over time drastically. To further support our statement, we display in figure 1 the t-sne [28] and UMAP [29] representations for years 2018-2019 and 2020-2021 for both countries. We observe significant dissimilarities between datasets. For country A, we observe an entire subsample of payments made in 2020-2021 at the bottom left of the UMAP graph, which does not exist for the 2018-2019 period. Similarly, for the t-sne representation of country A, we observe a similar pattern with few data sample overlays between periods. Moreover, for country B, we also observe a very scarce overlay on the graph between periods, especially for the UMAP representation. To further investigate whether distribution shift is present in our datasets, we rely on the Optimal Transport Dataset Distance method (OTDD) [1] to measure the distance between datasets from each period. This method relies on optimal transport, which measures the distance between distributions. For each country, we created two subsamples of 5000 observations for each period and compared

---

[5]Due to confidentiality, we cannot discuss the exact proportion of frauds within our dataset.

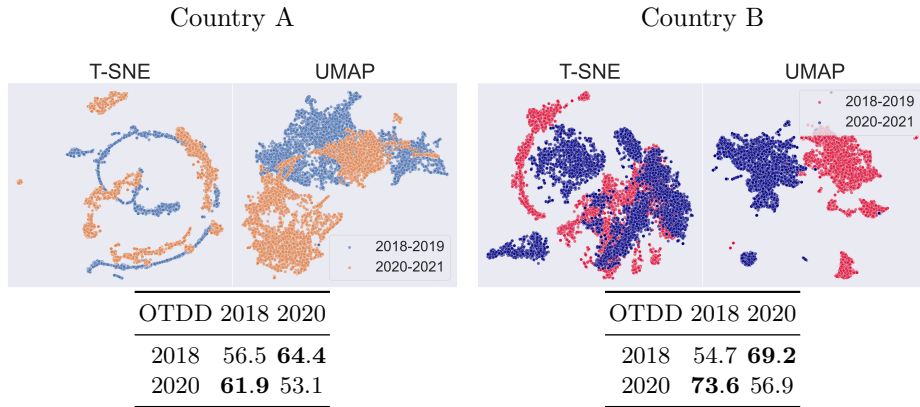|  | Country A |  |  |  | Country B |  |  |
|---|---|---|---|---|---|---|---|
|  | T-SNE | UMAP |  |  | T-SNE | UMAP |  |

Fig. 1: T-sne [28] and UMAP [29] bi-dimensional representation of payments in countries A and B for each period. These graphs give evidence of a distribution shift for both countries' payment behaviors between 2018 and 2020 since we observe very few sample overlays. Tables give the Optimal Transport Dataset Distance [1] between subsamples for each country. We observe a much higher distance between subsamples from the same country between different periods than for the same period

**Country A**

| OTDD | 2018 | 2020 |
|---|---|---|
| 2018 | 56.5 | **64.4** |
| 2020 | **61.9** | 53.1 |

**Country B**

| OTDD | 2018 | 2020 |
|---|---|---|
| 2018 | 54.7 | **69.2** |
| 2020 | **73.6** | 56.9 |

the distance between the subsample of the same period and between periods. Results of this analysis are shown in the tables in figure 1 and indicate that the distance between the dataset increases across periods. This is especially true for country B.

**Data splits and preprocessing** For a fair comparison between supervised approaches and anomaly detection methods, we split the 2018 datasets of each country into two separate datasets constituted of legit payments and frauds. We take a training set representing 75% of the 2018 dataset for each country and use a test set the 25% remaining. We include the frauds in the training set for LightGBM, while for anomaly detection approaches, the frauds are excluded from the training set. The 2020 dataset of each country serves entirely as test sets. Overall, the considered dataset used for every tested model contains features describing characteristics of the payment (e.g. amount of the transaction, the currency used, duration since the last transaction.), among which eight are categorical. All eight categorical features are encoded using Catboost encoding following [5]. Continuous features are scaled to be in $(0,1)$ through standard normalizing by removing the mean and reducing to unit variance.

### 3.2 Experimental Settings

In order to evaluate the suitability of anomaly detection methods for fraud detection, we conduct a comprehensive analysis using state-of-the-art approaches

on our dataset. We investigate both deep learning-based techniques designed for tabular data, such as the self-supervised approaches of GOAD [4], NeuTraL-AD [31] the contrastive approach proposed in [36], and the reconstruction-based approach of NPT-AD [40]. Additionally, we include non-deep learning methods, namely Isolation Forest [27], ECOD [26], COPOD [25], and the KNN AD approach [2, 32], as they have demonstrated remarkable performance on various tabular datasets. As a baseline, we employ LightGBM [23], a well-established supervised classification method, to assess the added value of anomaly detection compared to standard supervised techniques in this context.

To effectively compare the performance of various models in detecting fraud, we employ three commonly used metrics from the anomaly detection literature and the banking industry for real-life model evaluation. The anomaly detection literature has widely adopted the F1-score and the AUROC as evaluation metrics. While the AUROC is suitable for balanced class distributions, it may not fully account for class proportions, which is crucial in assessing performance in imbalanced set-ups. We include the Area Under the Precision-Recall Curve (AUPRC) to address this limitation, which is better suited for imbalanced datasets. In support of this choice, [8] has demonstrated that a model dominates in the ROC space if and only if it also dominates in the PR space.

For the F1-Score, whose value depends on a threshold, we adhere to the practices of the anomaly detection literature by selecting the threshold for the anomaly score that predicts an equal number of fraud cases as those present in the dataset. This approach ensures a fair and consistent evaluation across all models.

## 4 Results

### 4.1 Models Hyperparameters

We implemented the non-deep models using the PyOD library [43] with default hyperparameter values and LightGBM [23] also with default parameters. For the deep learning approaches, we adopted the hyperparameters suggested in the original papers for the dataset that most resembled our datasets. Specifically, we set the hyperparameters to those used for the KDD dataset for NeuTraL-AD [31], GOAD [4], and NPT-AD [40] using their official implementations available on GitHub. Regarding the approach proposed by [36], we kept the parameters at their default values as specified in their implementation. Deep models were trained on 4 Nvidia GPUs V100 16Go/32Go, while non-deep models were trained on 2 Intel Cascade Lake 6248 processors (20 cores at 2.5 GHz), thus 40 cores.

### 4.2 Fraud Detection Performances

Metrics reported in table 1 are averaged over 10 runs and we performed t-tests on the highest metrics to asses whether models obtained significantly different results. Among the AD approaches, we observe that the non-deep methods demonstrate the best performance, specifically ECOD, COPOD, and KNN. However,

Table 1: Performance metrics of AD models in comparison with LightGBM [23]. Results are averaged over 10 runs for 10 different splits of the data, the standard deviation are displayed below the metrics. We report the F1-Score in terms of percentage. The highest metric over all models is highlighted in bold, while the highest metrics among the AD method are underlined. We perform 5% t-test between the highest metrics to measure whether they are statistically different

| Model | Country A | | | | | | Country B | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | F1 (↑) | | AUROC (↑) | | AUPRC (↑) | | F1 (↑) | | AUROC (↑) | | AUPRC (↑) | |
| | 2018 | 2020 | 2018 | 2020 | 2018 | 2020 | 2018 | 2020 | 2018 | 2020 | 2018 | 2020 |
| LightGBM | **21.52** | **0.7** | **89.98** | **66.49** | **18.74** | **0.31** | **17.15** | 0.48 | **93.5** | **75.51** | **34.84** | **2.68** |
| | 1.6 | 0.33 | 0.41 | 2.38 | 1.78 | 0.05 | 1.93 | 0.39 | 0.31 | 1.24 | 1.48 | 0.29 |
| ECOD | 0.48 | 0.2 | 62.2 | <u>62.49</u> | 0.4 | <u>0.25</u> | 0.57 | 1.04 | 54.02 | 51.59 | 1.09 | 0.76 |
| | 0.2 | 0.22 | 0.64 | 1.02 | 0.02 | 0.01 | 0.26 | 0.38 | 0.54 | 0.87 | 0.06 | 0.04 |
| COPOD | 0.34 | 0.16 | 64.77 | <u>62.64</u> | 0.43 | <u>0.25</u> | 0.5 | 1.12 | 51.7 | 50.15 | 1.0 | 0.73 |
| | 0.21 | 0.15 | 0.51 | 0.98 | 0.02 | 0.01 | 0.2 | 0.44 | 0.59 | 0.91 | 0.05 | 0.04 |
| Isolation Forest | 0.16 | 0.19 | 64.14 | 60.55 | 0.43 | 0.23 | 0.71 | **<u>1.3</u>** | 60.52 | 46.86 | 1.35 | 0.67 |
| | 0.12 | 0.19 | 0.75 | 0.76 | 0.02 | 0.01 | 0.23 | 0.34 | 0.52 | 0.84 | 0.07 | 0.03 |
| KNN | 0.34 | 0.01 | <u>68.87</u> | 55.6 | 0.51 | 0.18 | 0.78 | 0.38 | <u>65.92</u> | 49.28 | <u>1.58</u> | 0.63 |
| | 0.12 | 0.04 | 0.76 | 0.85 | 0.02 | 0.01 | 0.21 | 0.27 | 0.64 | 0.67 | 0.08 | 0.02 |
| GOAD | 0.14 | 0.19 | 53.72 | 52.73 | 0.17 | 0.17 | 0.7 | 0.69 | 50.36 | <u>64.45</u> | 0.67 | <u>1.03</u> |
| | 0.09 | 0.13 | 1.41 | 1.69 | 0.01 | 0.01 | 0.36 | 0.34 | 2.35 | 1.25 | 0.05 | 0.06 |
| NeuTraL-AD | 0.6 | 0.02 | 59.12 | 51.52 | 0.35 | 0.15 | 1.45 | 0.38 | 53.23 | 45.19 | 1.08 | 0.58 |
| | 0.22 | 0.08 | 3.56 | 1.15 | 0.05 | 0.01 | 0.44 | 0.21 | 1.9 | 1.75 | 0.07 | 0.03 |
| Internal Cont. | 0.64 | 0.0 | 39.43 | 46.7 | 0.18 | 0.13 | 1.21 | 0.23 | 45.63 | 50.66 | 0.87 | 0.68 |
| | 0.08 | 0.0 | 1.05 | 0.17 | 0.0 | 0.0 | 0.16 | 0.07 | 2.46 | 0.9 | 0.1 | 0.03 |
| NPT-AD | <u>0.97</u> | **<u>0.66</u>** | 67.21 | 53.2 | <u>0.81</u> | 0.18 | <u>1.67</u> | 0.58 | <u>66.21</u> | 53.45 | 1.28 | 0.61 |
| | 0.07 | 0.06 | 1.25 | 0.65 | 0.01 | 0.03 | 0.11 | 0.03 | 1.14 | 0.43 | 0.11 | 0.06 |

it is worth noting that the deep learning approach NPT-AD also yields comparable results. While the AD methods achieve satisfactory results regarding the AUROC, their performances are consistently poor for both the F1-Score and AUPRC metrics across both countries and periods. In contrast, LightGBM exhibits significantly better performance across all metrics, consistently achieving the highest values, except for the F1-Score on the 2020 dataset of Country B. The overall poor performance of all models, including LightGBM, for the F1-Score and AUPRC, highlights the inherent challenges associated with fraud detection. However, a noteworthy observation is the substantial performance gap between LightGBM and the anomaly detection methods.

## 5   Discussion

### 5.1   Distribution Shift

The obtained results for both countries support the hypothesis that a distribution shift occurred between 2018 and 2020 in our dataset. Across most tested

approaches, we observe a significant decrease in all metrics between the 2018 and 2020 datasets for both countries. Notably, the distribution shift appears more pronounced in country B, with metrics experiencing a more significant decline. The findings presented in Figure 1 further corroborate this statement as the dataset distance is higher between periods than for country A, and the bi-dimensional representations also show less overlap between the periods than for country A. Although the AD methods display poor overall performance, they exhibit more resilience in the face of distribution shift than LightGBM. This trend is particularly evident for ECOD and COPOD, which maintained relatively similar metrics between the 2018 and 2020 datasets for both countries. Conversely, KNN and NPT-AD experienced a significant decline in performance across both periods and datasets compared to ECOD and COPOD. While Light-GBM still achieves the highest metrics for most of the 2020 dataset, it suffers a substantial drop in performance between the two periods. This drop is especially pronounced in the AUPRC and F1-Score metrics. Notably, most AD methods outperform LightGBM by a significant margin in terms of the F1-Score for the 2020 dataset of country B.

Based on these findings on our dataset, it appears that LightGBM is a favorable choice in the fraud detection framework when labels are available and no distribution shift occurs. However, in the presence of a distribution shift, retraining LightGBM on an updated dataset becomes crucial to prevent a significant performance decline.

### 5.2 Anomaly Detection Methods for Ensembling

One potential advantage of anomaly detection (AD) methods is their ability to identify fraud cases that differ from those flagged by supervised approaches. If AD models can successfully detect fraud instances supervised models cannot identify, resorting to ensembling techniques could enhance overall fraud detection performance. To investigate this further, we focused on ECOD, one of the top-performing AD methods on the dataset consisting of payments in *country A*. We examined whether the fraud cases detected by ECOD differed from those identified by LightGBM. Across the 10 iterations, we observed that, on average, 3.28% of the fraud cases in the test set were detected by ECOD but not by LightGBM. Conversely, LightGBM detected 20.41% of the fraud cases in the test set that ECOD did not flag. Notably, the 3.28% represents 10.98% of the fraud cases detected by ECOD. In other words, 89.02% of the fraud cases detected by ECOD were also detected by LightGBM. As a result, ensembling models by combining AD methods with LightGBM to enhance fraud detection performance may prove ineffective.

## 6 Conclusion

In conclusion, our study highlights several key findings concerning applying machine learning techniques for fraud detection. Our results demonstrate that

LightGBM consistently outperforms the tested AD methods across various evaluation metrics, emphasizing its efficacy in fraud detection tasks compared to other methods. However, we also observed that LightGBM's performances are susceptible to degradation due to distribution shift. This finding underscores the importance of retraining LightGBM on updated datasets when there is suspicion or evidence of a distribution shift. By adapting the model to the changing data distribution, it is possible to mitigate the drop in performance and maintain its effectiveness in fraud detection. Furthermore, our investigation revealed that ensembling techniques with AD methods would not significantly improve overall fraud detection performance. Despite the potential for AD methods to detect frauds that may elude supervised approaches, our analysis showed that LightGBM also detected most of the frauds identified by AD methods. This finding suggests limited benefits in combining AD methods with LightGBM in our specific fraud detection framework. We believe these insights may contribute to advancing the field of fraud detection and inform practitioners in selecting appropriate models and strategies for robust and accurate fraud detection systems.

Future work may involve replicating our analysis on other credit card payment datasets to determine whether our obtained results can be generalized. Furthermore, enhancing the robustness of GBDT models against distribution shifts emerges as a critical direction for further exploration. Addressing this challenge is paramount for financial institutions, as it enables them to embrace machine learning techniques in fraud detection systems confidently.

## References

1. Alvarez-Melis, D., Fusi, N.: Geometric dataset distances via optimal transport. In: H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, H. Lin (eds.) Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual (2020). URL https://proceedings.neurips.cc/paper/2020/hash/f52a7b2610fb4d3f74b4106fb80b233d-Abstract.html
2. Angiulli, F., Pizzuti, C.: Fast outlier detection in high dimensional spaces. In: European Conference on Principles of Data Mining and Knowledge Discovery, pp. 15–27. Springer (2002)
3. Bauder, R.A., Khoshgoftaar, T.M., Hasanin, T.: An empirical study on class rarity in big data. In: 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA), pp. 785–790 (2018). DOI 10.1109/ICMLA.2018.00125
4. Bergman, L., Hoshen, Y.: Classification-based anomaly detection for general data. In: International Conference on Learning Representations (2020). URL https://openreview.net/forum?id=H1lK_lBtvS

5. Bourdonnaye, F.D.L., Daniel, F.: Evaluating categorical encoding methods on a real credit card fraud detection database. CoRR **abs/2112.12024** (2021). URL https://arxiv.org/abs/2112.12024

6. Breunig, M.M., Kriegel, H., Ng, R.T., Sander, J.: LOF: identifying density-based local outliers. In: W. Chen, J.F. Naughton, P.A. Bernstein (eds.) Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data, May 16-18, 2000, Dallas, Texas, USA, pp. 93–104. ACM (2000). DOI 10.1145/342009.335388. URL https://doi.org/10.1145/342009.335388

7. Chen, T., Guestrin, C.: Xgboost: A scalable tree boosting system. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, p. 785–794. Association for Computing Machinery, New York, NY, USA (2016). DOI 10.1145/2939672.2939785. URL https://doi.org/10.1145/2939672.2939785

8. Davis, J., Goadrich, M.: The relationship between precision-recall and roc curves. In: Proceedings of the 23rd International Conference on Machine Learning, ICML '06, p. 233–240. Association for Computing Machinery, New York, NY, USA (2006). DOI 10.1145/1143844.1143874. URL https://doi.org/10.1145/1143844.1143874

9. Dutta, H., Giannella, C., Borne, K.D., Kargupta, H.: Distributed top-k outlier detection from astronomy catalogs using the DEMAC system. In: Proceedings of the Seventh SIAM International Conference on Data Mining, April 26-28, 2007, Minneapolis, Minnesota, USA, pp. 473–478. SIAM (2007). DOI 10.1137/1.9781611972771.47. URL https://doi.org/10.1137/1.9781611972771.47

10. Finke, T., Krämer, M., Morandini, A., Mück, A., Oleksiyuk, I.: Autoencoders for unsupervised anomaly detection in high energy physics. Journal of High Energy Physics **2021**(6) (2021). DOI 10.1007/jhep06(2021)161. URL http://dx.doi.org/10.1007/JHEP06(2021)161

11. Frery, J.: Ensemble Learning for Extremely Imbalced Data Flows. Theses, Université de Lyon (2019). URL https://tel.archives-ouvertes.fr/tel-02899943

12. Gong, D., Liu, L., Le, V., Saha, B., Mansour, M.R., Venkatesh, S., van den Hengel, A.: Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection. 2019 IEEE/CVF International Conference on Computer Vision (ICCV) pp. 1705–1714 (2019)

13. Gopalan, P., Sharan, V., Wieder, U.: Pidforest: Anomaly detection and certification via partial identification. In: Neural Information Processing Systems (2019). URL https://api.semanticscholar.org/CorpusID:202766416

14. Gorishniy, Y., Rubachev, I., Kartashev, N., Shlenskii, D., Kotelnikov, A., Babenko, A.: Tabr: Unlocking the power of retrieval-augmented tabular deep learning (2023)

15. Gorishniy, Y., Rubachev, I., Khrulkov, V., Babenko, A.: Revisiting deep learning models for tabular data. In: A. Beygelzimer, Y. Dauphin, P. Liang, J.W. Vaughan (eds.) Advances in Neural Information Processing Systems (2021). URL https://openreview.net/forum?id=i_Q1yrOegLY

16. Grinsztajn, L., Oyallon, E., Varoquaux, G.: Why do tree-based models still outperform deep learning on typical tabular data? In: Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (2022). URL https://openreview.net/forum?id=Fp7__phQszn

17. Gu, S., Ślusarczyk, B., Hajizada, S., Kovalyova, I., Sakhbieva, A.: Impact of the covid-19 pandemic on online consumer purchasing behavior. Journal of Theoretical and Applied Electronic Commerce Research **16**(6), 2263–2281 (2021). DOI 10.3390/jtaer16060125. URL https://www.mdpi.com/0718-1876/16/6/125

18. Guha, S., Mishra, N., Roy, G., Schrijvers, O.: Robust random cut forest based anomaly detection on streams. In: ICML (2016). URL https://www.amazon.science/publications/robust-random-cut-forest-based-anomaly-detection-on-streams

19. Hariri, S., Kind, M.C., Brunner, R.J.: Extended isolation forest. IEEE Transactions on Knowledge and Data Engineering **33**(4), 1479–1489 (2021). DOI 10.1109/TKDE.2019.2947676

20. Huang, L., Nguyen, X., Garofalakis, M.N., Jordan, M.I., Joseph, A.D., Taft, N.: In-network PCA and anomaly detection. In: B. Schölkopf, J.C. Platt, T. Hofmann (eds.) Advances in Neural Information Processing Systems 19, Proceedings of the Twentieth Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 4-7, 2006, pp. 617–624. MIT Press (2006). URL https://proceedings.neurips.cc/paper/2006/hash/2227d753dc18505031869d44673728e2-Abstract.html

21. Japkowicz, N., Stephen, S.: The class imbalance problem: A systematic study. Intell. Data Anal. **6**(5), 429–449 (2002)

22. Kadra, A., Lindauer, M., Hutter, F., Grabocka, J.: Well-tuned simple nets excel on tabular datasets. In: Thirty-Fifth Conference on Neural Information Processing Systems (2021)

23. Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., Liu, T.Y.: Lightgbm: A highly efficient gradient boosting decision tree. Advances in neural information processing systems **30**, 3146–3154 (2017)

24. Kim, K.H., Shim, S., Lim, Y., Jeon, J., Choi, J., Kim, B., Yoon, A.S.: Rapp: Novelty detection with reconstruction along projection pathway. In: International Conference on Learning Representations (2020). URL https://openreview.net/forum?id=HkgeGeBYDB

25. Li, Z., Zhao, Y., Botta, N., Ionescu, C., Hu, X.: Copod: Copula-based outlier detection. 2020 IEEE International Conference on Data Mining (ICDM) (2020). DOI 10.1109/icdm50108.2020.00135. URL http://dx.doi.org/10.1109/ICDM50108.2020.00135

26. Li, Z., Zhao, Y., Hu, X., Botta, N., Ionescu, C., Chen, G.: Ecod: Unsupervised outlier detection using empirical cumulative distribution functions. IEEE Transactions on Knowledge and Data Engineering pp. 1–1 (2022). DOI 10.1109/TKDE.2022.3159580

27. Liu, F.T., Ting, K.M., Zhou, Z.H.: Isolation forest. In: 2008 Eighth IEEE International Conference on Data Mining, pp. 413–422 (2008). DOI 10.1109/ICDM.2008.17

28. van der Maaten, L., Hinton, G.: Visualizing data using t-SNE. Journal of Machine Learning Research **9**, 2579–2605 (2008). URL http://www.jmlr.org/papers/v9/vandermaaten08a.html

29. McInnes, L., Healy, J., Melville, J.: Umap: Uniform manifold approximation and projection for dimension reduction (2020)

30. Pol, A.A., Berger, V., Cerminara, G., Germain, C., Pierini, M.: Anomaly Detection With Conditional Variational Autoencoders. In: ICMLA 2019 - 18th IEEE International Conference on Machine Learning and Applications, 18th International Conference on Machine Learning Applications. Boca Raton, United States (2019). URL https://hal.inria.fr/hal-02396279

31. Qiu, C., Pfrommer, T., Kloft, M., Mandt, S., Rudolph, M.: Neural transformation learning for deep anomaly detection beyond images. In: International Conference on Machine Learning, pp. 8703–8714. PMLR (2021)

32. Ramaswamy, S., Rastogi, R., Shim, K.: Efficient algorithms for mining outliers from large data sets. In: ACM Sigmod Record, vol. 29, pp. 427–438. ACM (2000)

33. Ruff, L., Vandermeulen, R., Goernitz, N., Deecke, L., Siddiqui, S.A., Binder, A., Müller, E., Kloft, M.: Deep one-class classification. In: J. Dy, A. Krause (eds.) Proceedings of the 35th International Conference on Machine Learning, *Proceedings of Machine Learning Research*, vol. 80, pp. 4393–4402. PMLR (2018). URL https://proceedings.mlr.press/v80/ruff18a.html

34. Ruff, L., Vandermeulen, R.A., Görnitz, N., Binder, A., Müller, E., Müller, K.R., Kloft, M.: Deep semi-supervised anomaly detection. In: International Conference on Learning Representations (2020). URL https://openreview.net/forum?id=HkgH0TEYwH

35. Schölkopf, B., Williamson, R., Smola, A., Shawe-Taylor, J., Platt, J.: Support vector method for novelty detection. In: Proceedings of the 12th International Conference on Neural Information Processing Systems, NIPS'99, p. 582–588. MIT Press, Cambridge, MA, USA (1999)

36. Shenkar, T., Wolf, L.: Anomaly detection for tabular data with internal contrastive learning. In: International Conference on Learning Representations (2022). URL https://openreview.net/forum?id=_hszZbt46bT

37. Somepalli, G., Goldblum, M., Schwarzschild, A., Bruss, C.B., Goldstein, T.: SAINT: improved neural networks for tabular data via row attention and contrastive pre-training. CoRR **abs/2106.01342** (2021). URL https://arxiv.org/abs/2106.01342

38. Tack, J., Mo, S., Jeong, J., Shin, J.: Csi: Novelty detection via contrastive learning on distributionally shifted instances. In: H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, H. Lin (eds.) NeurIPS (2020). URL https://proceedings.neurips.cc/paper/2020

39. Tax, D., Duin, R.: Support vector data description. Machine Learning **54**, 45–66 (2004). DOI 10.1023/B:MACH.0000008084.60811.49

40. Thimonier, H., Popineau, F., Rimmel, A., DOAN, B.L.: Beyond individual input for deep anomaly detection on tabular data. In: NeurIPS 2023 Second Table Representation Learning Workshop (2023). URL https://openreview.net/forum?id=lsn7ehxAdt

41. Thimonier, H., Popineau, F., Rimmel, A., Doan, B.L., Daniel, F.: TracInAD: Measuring influence for anomaly detection. In: 2022 International Joint Conference on Neural Networks (IJCNN), pp. 1–6 (2022). DOI 10.1109/IJCNN55064.2022.9892058

42. Yanmin, S., Wong, A., Kamel, M.S.: Classification of imbalanced data: a review. International Journal of Pattern Recognition and Artificial Intelligence **23**, 687–719 (2011). DOI 10.1142/S0218001409007326

43. Zhao, Y., Nasrullah, Z., Li, Z.: Pyod: A python toolbox for scalable outlier detection. Journal of Machine Learning Research **20**(96), 1–7 (2019). URL http://jmlr.org/papers/v20/19-011.html