

This version of the article has been accepted for publication, after peer review (when applicable) and is subject to Springer Nature's [AM terms of use](#), but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at: <https://doi.org/10.1007/s10796-022-10346-6>.

# Fraud Detection in Mobile Payment Systems using an XGBoost-based Framework

Petr Hajek

*Science and Research Centre, Faculty of Economics and Administration, University of Pardubice, Studentska 84, Pardubice 532 10, Czech Republic*

Mohammad Zoynul Abedin\*

*Department of Finance, Performance Marketing, Teesside University International Business School, Teesside University, Middlesbrough, TS1 3BX Tees Valley, UK*

Uthayasankar Sivarajah

*Bradford University School of Management, Emm Lane, Heaton, Bradford, UK*

---

## Abstract

Mobile payment systems are becoming more popular due to the increase in the number of smartphones, which, in turn, attracts the interest of fraudsters. Extant research has therefore developed various fraud detection methods using supervised machine learning. However, sufficient labeled data are rarely available and their detection performance is negatively affected by the extreme class imbalance in financial fraud data. The purpose of this study is to propose an XGBoost-based fraud detection framework while considering the financial consequences of fraud detection systems. The framework was empirically validated on a large dataset of more than 6 million mobile transactions. To demonstrate the effectiveness of the proposed framework, we conducted a comparative evaluation of existing machine learning methods designed for modeling imbalanced data and outlier detection. The results suggest that in terms of standard classification measures, the proposed semi-supervised ensemble model integrating multiple unsupervised outlier detection algorithms

---

\*corresponding author

*Email addresses:* petr.hajek@upce.cz (Petr Hajek), m.abedin@tees.ac.uk (Mohammad Zoynul Abedin), u.sivarajah@bradford.ac.uk (Uthayasankar Sivarajah)

and an XGBoost classifier achieves the best results, while the highest cost savings can be achieved by combining random under-sampling and XGBoost methods. This study has therefore financial implications for organizations to make appropriate decisions regarding the implementation of effective fraud detection systems.

*Keywords:* mobile payment, fraud detection, machine learning, imbalanced data, outlier detection

---

## 1. Introduction

Mobile payment transactions are carried out using mobile phone technologies that allow users to deposit, withdraw, spend, transfer and send money. There are nearly three hundred mobile payment services worldwide, which are particularly popular in Sub-Saharan Africa and Asia. In 2020, mobile payment transactions totaled \$767 billion, conducted by approximately 1.2 billion registered users according to Statista. In addition, mobile payments have reportedly enormous potential during the COVID-19 pandemic, as it can greatly increase the promptness and efficiency of money transfers while minimising the necessity of face-to-face contact with bank and government staff (Blumenstock, 2020).

Recent mobile payment case studies (Iman, 2018; Jocevski et al., 2020; Verkijika, 2020) suggest that mobile payment systems have been challenged by several types of factors that have emerged in the context of advances in financial technology. Commercial and technical factors have been identified as particularly important to their future growth. As regards the first group of factors, the need to increase cost efficiency is particularly emphasised because most mobile payment transactions in developing countries are low value but high volume (Franque et al., 2022). Technical factors include, in particular, security concerns, as the legal frameworks and enforcement mechanisms are often inadequate in developing countries (Akanfe et al., 2020; David-West et al., 2022; Pal et al., 2020). To deploy a mobile payment system, it is therefore necessary to minimise fraud in order to increase customer trust and security, as reported in existing mobile payment acceptance models (Chin et al., 2022; Jia et al., 2022; Kar, 2021; Pal et al., 2021).

The increasingly growing use of mobile payments has boosted the chances of criminals committing mobile phone fraud in an illegal effort to circumvent security measures of mobile payment services. There is consequently a lot

29 of pressure to investigate potential security threats that may be exploited,  
30 with the ultimate aim of preventing fraud on a mobile payment service and  
31 developing countermeasures against attacks (Chen et al., 2021; Lopez-Rojas  
32 et al., 2016; Rieke et al., 2013). Early detection of fraudulent transactions  
33 is a key task in this effort. Recent developments in mobile payment services  
34 have therefore heightened the need for automated detection systems that  
35 enable immediate detection and prevention of fraudulent transactions.

36 The main challenges currently facing researchers involved in detecting  
37 fraud in mobile payment transactions include: (1) extreme class imbalance  
38 (only a small proportion of customers have fraudulent intentions); (2) chang-  
39 ing patterns of fraud over time (fraudsters are always looking for new ways  
40 to bypass systems and commit crimes); and (3) inadequate selection of per-  
41 formance metrics. The consequence of the first challenge is a poor user ex-  
42 perience for legitimate customers, as the detection of fraudsters usually also  
43 implies rejecting some legitimate mobile payment transactions. The second  
44 challenge usually leads to a decrease in the performance and efficiency of the  
45 detection model. Therefore, machine learning models must be constantly  
46 updated, otherwise they will not meet their objectives. Regarding the last  
47 challenge, in some cases the providers of mobile payment systems should  
48 prefer a higher false positive rate in exchange for a lower false negative rate  
49 and vice versa. But how to choose the right ratio between these two errors  
50 remains a challenging area in the field of fraud detection in mobile payment  
51 transactions.

52 A relatively high detection accuracy was reported in earlier research by us-  
53 ing both traditional supervised learning methods (Choi and Lee, 2017, 2018)  
54 and deep learning-based methods (Mubalalike and Adali, 2018; Xenopou-  
55 los, 2017). However, a major problem with this kind of application is the  
56 extreme class imbalance of transactions, with a considerable dominance of le-  
57 gitimate transactions in the data. This in turn leads to a poor classification  
58 performance on the minority class of fraudulent transactions. To address  
59 this issue, two approaches have been utilized. The first approach relies on  
60 under-sampling methods used to generate a balanced dataset (Pambudi et  
61 al., 2019). The main limitation of this approach is the loss of potentially im-  
62 portant information stored in discarded legitimate transactions, which can  
63 reduce detection accuracy. Alternatively, an attempt has been made to iso-  
64 late fraudulent transactions in an unsupervised fashion (Buschjäger et al.,  
65 2021), inspired by outlier detection methods. Nevertheless, a comprehensive  
66 evaluation of machine learning methods is not yet available in the literature.

67 Moreover, little is known about how the two approaches can be integrated  
68 to improve the detection performance. To overcome the above problems,  
69 here we propose to enhance the performance of eXtreme Gradient boost-  
70 ing (XGBoost), a state-of-the-art machine learning method, by including a  
71 data sampling component addressing the issue of extreme class imbalance of  
72 mobile payment transactions.

73 In many financial applications it is necessary to filter out unusual ob-  
74 servations to ensure the reliability of the system and prevent attempts to  
75 maliciously use it. This is particularly useful for detecting financial fraud at-  
76 tempts, as their behaviour patterns differ significantly from normal financial  
77 transactions (Bernard et al., 2021). Outlier detection methods are capable of  
78 processing all available data in real time to uncover patterns that evade tra-  
79 ditional supervised learning methods. By doing so, organised crime groups  
80 can be identified with higher accuracy and less false positives. Outlier de-  
81 tection methods have indeed proved effective for detecting credit card fraud  
82 detection (Carcillo et al., 2021), online banking fraud detection (Carminati  
83 et al., 2015), and health insurance fraud detection (Yamanishi et al., 2004).  
84 Overall, however, there has been limited use of these methods to detect fi-  
85 nancial fraud, although some review studies suggest that they deserve more  
86 attention because the detection performance of supervised algorithms is neg-  
87 atively affected by the inherently heavily imbalanced class distribution of  
88 financial fraud data (Ngai et al., 2011). The scarce use of outlier detection  
89 methods can be attributed to the difficulty of detecting fraudulent behaviour  
90 (e.g., abnormal frequency of transactions or spending behaviour) when over-  
91 lapping with legitimate behaviour in datasets contaminated with outliers and  
92 noise. Moreover, several other challenges have been identified that make it  
93 the difficult to detect outliers in the financial domain. First, efficient general  
94 purpose outlier detection methods are lacking because an outlier detection  
95 method in one fraud domain may not be appropriate for other scenarios,  
96 as legitimate and fraudulent behaviour is different from domain to domain  
97 (Ahmed et al., 2016). Second, unsupervised learning is preferred as suffi-  
98 cient labelled data for building models are rarely available. Third, legitimate  
99 behaviour may change over time, and fraudsters try to make their activities  
100 look legitimate. To take advantages of both supervised machine learning and  
101 outlier detection methods, for the first time, we propose a semi-supervised  
102 ensemble fraud detection model combining unsupervised outlier detection  
103 and supervised XGBoost methods that exploit all transactions contained in  
104 a large, highly imbalanced mobile payment transaction dataset.

105 Finally, financial implications of fraud detection methods in mobile pay-  
106 ment transactions have also been neglected in earlier research. Therefore, our  
107 third contribution is to propose a novel performance measure of cost savings  
108 that takes into account the financial implications of false positive and false  
109 negative rates of fraud detection systems. Using the PaySim dataset, our  
110 findings provide evidence for the effectiveness of both XGBoost leveraged by  
111 an under-sampling class-balancing procedure and extreme gradient boost-  
112 ing outlier detection (XGBOD), thus providing important tools to support  
113 operation and management of mobile payment services.

114 In summary, the contributions of this study are threefold:

- 115 1. Developing a novel fraud detection framework for mobile payment sys-  
116 tems by integrating the XGBoost method with class-balancing adjust-  
117 ments and unsupervised outlier detection methods, making it suitable  
118 for detecting fraud in a typical class-imbalanced mobile payment sce-  
119 nario.
- 120 2. Proposing a novel cost savings measure to evaluate the performance of  
121 mobile payment fraud detection systems. Unlike the traditional perfor-  
122 mance measures, the proposed measure considers both the cost savings  
123 from the correct detection of fraudulent transactions and the decrease  
124 in the margin for the transactions incorrectly identified as fraudulent.
- 125 3. Using the benchmark PaySim dataset of more than 6 million mobile  
126 payment transactions, we demonstrate that the proposed fraud detec-  
127 tion framework not only outperforms state-of-the-art fraud detection  
128 methods in terms of detection accuracy but also generates substantial  
129 financial savings to the providers of mobile payment systems.

130 The remainder of this paper is organized as follows. Section 2 reviews  
131 the related work on fraud detection in mobile payment transactions with  
132 respect to data sources, methods used and performance achieved in earlier  
133 studies. Section 3 outlines the proposed fraud detection framework. Section  
134 4 provides the results of the evaluation on the PaySim dataset, robustness  
135 check, and financial implications. Section 5 concludes with providing some  
136 possible directions for future research.

## 137 **2. Fraud Detection in Mobile Payment Systems – Literature Re-** 138 **view**

139 A considerable amount of literature has been published on financial fraud  
140 detection, see West and Bhattacharya (2016) for a review and Hajek and

141 Henriques (2017) for a comprehensive evaluation of financial fraud detection  
142 methods. Risk factors of financial fraud were investigated, indicating that  
143 pressure / incentive to commit fraud is the most important risk factor (Huang  
144 et al., 2017). Related studies can be broadly categorized according to the  
145 financial fraud type as follows (Onwubiko, 2020): (1) account takeover fraud,  
146 (2) payment fraud, and (3) application fraud. Onwubiko (2020) also identi-  
147 fied four main fraud channels, namely physical, web, telephony, and mobile.  
148 Frauds in mobile payment transactions have increasingly been recognized as  
149 a major concern in finance due to recent developments in mobile payment  
150 services (Chen and Sivakumar, 2021). Therefore, security requirements must  
151 be met to address security issues related to mobile payment transactions,  
152 such as mobile malware and SMS-based attacks (Kang, 2018). Heteroge-  
153 neous software and hardware mobile platforms make the security problems  
154 more challenging (Li and Clark, 2013).

155 Regarding the data used in previous studies and summarized in Table  
156 1, the lack of real-world datasets has been identified as a major problem in  
157 the application domain. Therefore, most earlier research tended to generate  
158 simulated synthetic data based on features captured from real-world fraud  
159 and legitimate transactions. To do so, Rieke et al. (2013) extracted pay-  
160 ment laundering patterns from real-world events. However, the number of  
161 instances was insufficient for efficient fraud detection, as indicated by rela-  
162 tively low false negative (legitimate) rates in early studies (Coppolino, 2015;  
163 Rieke et al., 2013). Considerable progress has been made by introducing the  
164 PaySim financial simulator (Lopez-Rojas et al., 2016, 2018) that resembles  
165 normal mobile transactions and injects fraudulent behaviour to produce a  
166 larger number of financial frauds. Agent-based simulations and statistical  
167 analysis confirmed that the simulated data are as prudent as the original  
168 aggregated anonymized real data, thus, representing an optimal control en-  
169 vironment for fraud detection in mobile payment transactions. By leverag-  
170 ing the PaySim data, Lopez-Rojas and Barneaud (2019) demonstrated their  
171 advantages over the relatively small real-world dataset. In addition, the sim-  
172 ulated data retained the transactions and causal dynamics of the original  
173 data. It should be however noted that by preserving the statistical proper-  
174 ties of the real-world data, the high class imbalance in favour of legitimate  
175 transactions is also maintained in the simulated dataset.

176 Traditional machine learning methods with supervised or unsupervised  
177 learning are not effective in handling extreme class imbalance in the data.  
178 Although a relatively high overall accuracy was reported in several studies,

179 these methods performed well only in terms of majority (legitimate) class  
180 accuracy (Choi and Lee, 2017, 2018; Du et al., 2018; Zhou et al., 2018).  
181 This holds also for more recent deep learning models, such as deep belief  
182 networks (Xenopoulos, 2017) and restricted Boltzman machines (Mubalake  
183 and Adali, 2018). To overcome this major limitation, class imbalance was  
184 first approached by using under-sampling methods and then machine learn-  
185 ing methods were trained on the balanced dataset (Pambudi et al. 2019).  
186 Similarly, Xenopoulos (2017) used under-sampling to produce balanced boot-  
187 straps for ensemble learning, and Misra et al. (2020) and Schlör et al. (2021)  
188 applied it to generate balanced training data for deep learning-based detec-  
189 tion models. The main drawback of the under-sampling approach is that  
190 potentially useful instances are often excluded from the training data, which  
191 can significantly degrade the detection accuracy. Alternatively, isolation-  
192 based approaches were used to approximate the data distribution and build  
193 a generative model using mixture components. This outlier detection method  
194 was successfully applied to fraud detection by Buschjäger et al. (2021).

195 However, a comprehensive evaluation of state-of-the-art machine learning-  
196 based approaches exploiting under-sampling methods for handling class im-  
197 balance problem, is lacking in the literature. Hybrid semi-supervised meth-  
198 ods taking advantage of supervised learning and unsupervised outlier detec-  
199 tion methods have also been overlooked. Finally, only standard performance  
200 measures have been used to evaluate fraud detection performance in mobile  
201 payment systems, thus neglecting the financial implications of fraud detec-  
202 tion.

### 203 **3. Fraud Detection Framework**

204 The proposed framework for fraud detection in mobile payment systems is  
205 presented in Fig. 1. The proposed fraud detection models are aimed to take  
206 advantage of XGBoost while overcoming the problem of extremely imbal-  
207 anced classes in mobile payment transaction data. We will demonstrate that  
208 this approach is not only more accurate than supervised machine learning  
209 and outlier detection methods used in existing studies but that our approach  
210 is also more profitable in terms of the proposed cost savings measure.

#### 211 *3.1. Proposed Fraud Detection Models*

212 This section outlines two fraud detection models proposed in this study.  
213 First, the eXtreme Gradient boosting (XGBoost) method, augmented with

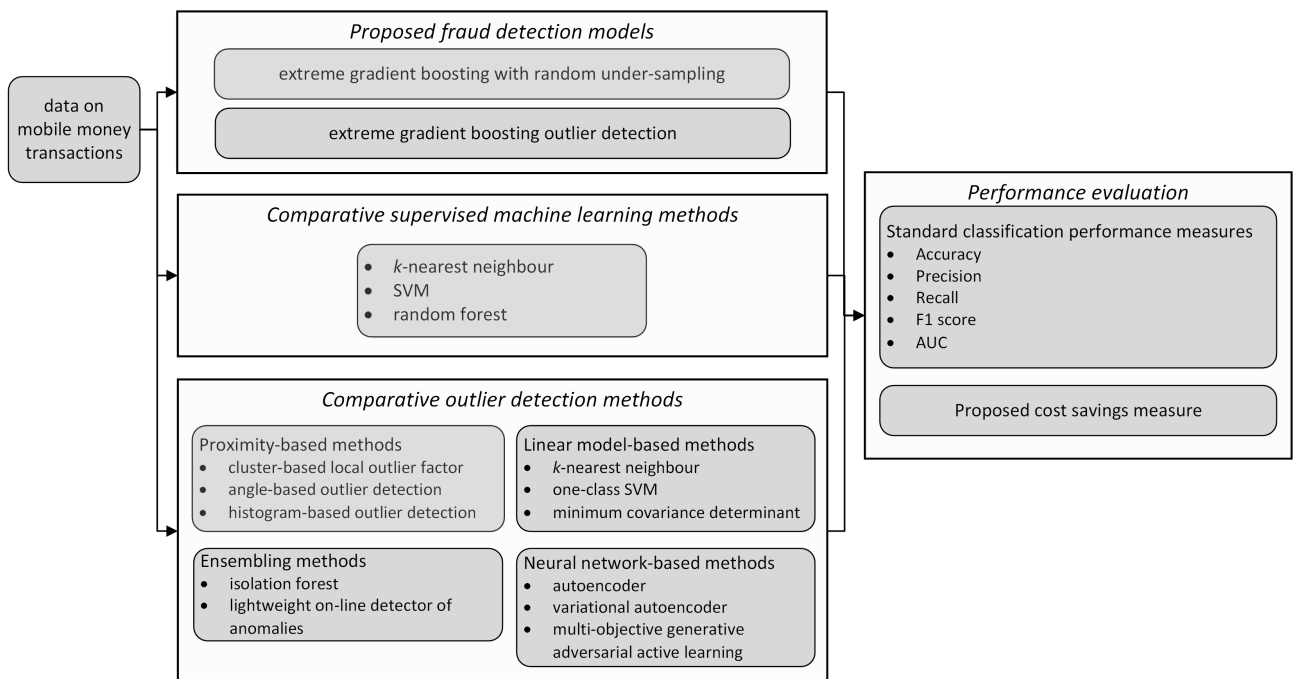


Figure 1: Fraud detection framework

Table 1: Summary of data and methods used in previous studies

Study	Data (# fraud / legitimate)	Method	Performance
Rieke et al. (2013)	synthetic logs (20/5,297)	predictive security analyser	$FNR=0.550$
Coppolino et al. (2015)	synthetic logs	Dempster-Shafer theory	$FNR=0.240$
Xenopoulos (2017)	PaySim (492/284,315)	ensemble of deep belief networks	$Acc=89.05$ , $AUC=0.961$
Choi and Lee (2017, 2018)	Korean payment data (2,402/274,670)	unsupervised (EM, K-means, Farthest-First, X-means, MakeDensity), supervised (NB, SVM, LR, OneR, C4.5, RF)	$Acc=99.97$
Mubalaike and Adali (2018)	PaySim (8,213/6M)	restricted Boltzman machines	$Acc=91.53$
Du et al. (2018)	PaySim (8,213/6M)	SVM with LogDet regularization	$Acc=97.57$ , $AUC=0.978$
Zhou et al. (2018)	Chinese bankcard enrollment (5,753/~52M)	GB DT, LR, RF, rule-based expert	$Precision=50.83$ , $Recall=0.25$
Pambudi et al. (2019)	PaySim (4,093/246,033)	RUS+SVM	$F1=0.900$ , $AUC=0.880$
Misra et al. (2020)	PaySim (492/284,315)	Autoencoder+MLP	$Acc=0.999$ , $F1=0.827$
Mendelson and Lerner (2020)	PaySim (8,213/6M)	cluster drift detection	$AUC=0.898$
Turner et al. (2021)	Bitcoin blockchain transactions	DeepWalk network analysis	-
Schlör et al. (2021)	PaySim (8,213/6M)	deep MLP with ReLU and iNALU	$F1=0.880$ , $AUC=0.960$
Buschjager et al. (2021)	PaySim (269/572K)	generalized Isolation Forest	$AUC=0.821$
This study	PaySim (8,213/6M)	RUS+XGBoost, XGBOD	

Legend:  $Acc$  – accuracy,  $AUC$  – area under ROC curve, DT – decision tree, EM – expectation-maximization,  $F1$  –  $F1$ -score (average of precision and recall),  $FNR$  – false negative (legitimate) rate, GB – gradient boosting, LR – logistic regression, MLP – multilayer perceptron, NB – Naïve Bayes, PNN – probabilistic neural network, RF – random forest, SVM – support vector machine, XGBOD – extreme gradient boosting outlier detection, and XGBoost – eXtreme Gradient boosting.

214 random under-sampling, is introduced to leverage both the supervised learning  
215 capability and robustness of XGBoost, a state-of-the-art machine learning  
216 method, and the data sampling component to overcome the class imbalance  
217 problem inherent in mobile payment transaction data. The second model  
218 exploits the extreme gradient boosting outlier detection (XGBOD) method,  
219 a semi-supervised algorithm that improves the performance of the XGBoost  
220 method on highly imbalanced mobile payment transaction data by intro-  
221 ducing outlier scores obtained from multiple unsupervised outlier detection  
222 methods.

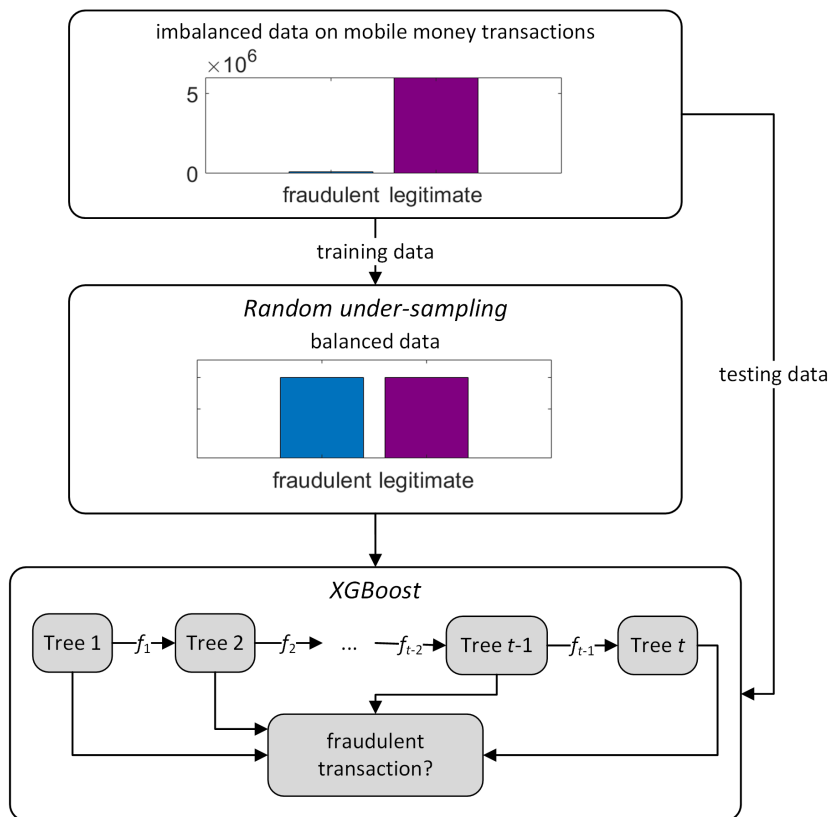


Figure 2: Flowchart of RUS-XGBoost for fraud detection

223 *3.1.1. Extreme gradient boosting with random under-sampling*

224 The proposed RUS+XGBoost integrates the random under-sampling (RUS)  
 225 method with XGBoost, as depicted in Fig. 2. The RUS component is first  
 226 used to generate balanced training samples, and XGBoost then generates ad-  
 227 ditive models to produce the final prediction on whether the mobile payment  
 228 transaction is fraudulent or not.

229 *Under-sampling for handling class imbalance problem*

230 The extremely high imbalance between legitimate and fraud classes makes  
 231 detecting financial fraud a challenge (Du et al., 2018). Considering the impor-  
 232 tance of class imbalance in financial fraud detection, numerous methods have  
 233 been used to improve the classification performance of supervised learning  
 234 methods. In the related literature (Pambudi et al., 2019), data-level solutions  
 235 have been particularly successful because they allow to address the imbalance

236 problem before training machine learning methods. In addition, data-level  
 237 methods integrated into classifier ensembles appear to be particularly effec-  
 238 tive (Galar et al., 2012). From the data-level methods, over-sampling meth-  
 239 ods create artificial instances in the minority class to balance the training  
 240 data. However, this can lead to problems of overfitting and overgeneral-  
 241 ization as instances of the majority class are ignored. Moreover, given the  
 242 gradual increase in data on financial fraud, under-sampling methods should  
 243 be a better choice than their over-sampling counterpart.

244 The RUS method used in this study enables controlling for the number  
 245 of samples selected from the original data. RUS is a non-heuristic method  
 246 that randomly selects a data subset from the majority class, which is com-  
 247 putationally effective and enables sampling heterogeneous data (Haixian et  
 248 al., 2017).

249 *Extreme gradient boosting*

250 XGBoost is a computationally efficient and scalable implementation of  
 251 gradient boosted decision trees that build additive models in a stepwise fash-  
 252 ion. The overall error is minimized incrementally by introducing additive  
 253 models based on the errors obtained in the previous steps. This results in  
 254 an ensemble of base learners with better prediction ability than the individ-  
 255 ual classifiers. This is achieved by gradually improving the accuracy, low tree  
 256 depth and equal contribution of the base learners to the final combined model.  
 257 To further improve robustness to noise and overfitting, gradient boosting was  
 258 augmented with a random sampling scheme (stochastic gradient boosting).  
 259 XGBoost is an enhanced implementation with a more regularized model to  
 260 control overfitting. The objective function of XGBoost to be minimized is  
 261 given as follows (Chen and Guestrin, 2016):

$$obj^{(t)} = \sum_{i=1}^n (y_i - (\hat{y}_i^{(t-1)} + f_t(x_i)))^2 + \sum_{t=1}^T \Omega(f_t), \quad (1)$$

262 where  $y_i$  is the target value of the  $i$ -th instance,  $\hat{y}_i^{(t)}$  is its predicted value at  
 263 the  $t$ -th iteration,  $f_t(x_i)$  is the additive decision tree model greedily added  
 264 to improve the model performance, and  $\Omega(f_t)$  is a regularization term penal-  
 265 izing the model complexity. The goal of this regularization procedure is to  
 266 compress the weights for many features to zero to perform feature selection,  
 267 which is advantageous when dealing with high-dimensional data. Therefore,  
 268 XGBoost is currently one of the best performing classifiers across domains  
 269 and has been successfully applied to insurance fraud detection (Dhieb et al.,

270 2019).

### 271 3.1.2. Extreme gradient boosting outlier detection model

272 The XGBOD method (Zhao and Hryniewicki, 2018) is a semi-supervised  
273 ensemble algorithm integrating multiple unsupervised outlier detection algo-  
274 rithms and an XGBoost classifier, as illustrated in Fig. 3. First, unsupervised  
275 methods are used to obtain data representations in terms of transformed out-  
276 lier scores (TOS). Second, a feature selection method is used to reduce the  
277 TOS feature space so that only relevant TOS are retained. Then, the outlier  
278 score matrix is combined with the original features to produce a combined  
279 feature space. An improved feature space is thus generated, and the XGBoost  
280 classifier is used in this feature space to produce the final outlier scores for  
281 each mobile payment transaction. The advantage of this approach is its  
282 good predictive ability, which is due to its robustness to overfitting and data  
283 imbalance.

284 In the proposed XGBOD-based fraud detection model, a variety of unsu-  
285 pervised outlier detection methods (presented in subsection 3.2.2) are used  
286 to produce the TOS features. To maintain the balance between their diver-  
287 sity and accuracy, the balance selection algorithm (Zhao and Hryniewicki,  
288 2018) is used to perform TOS selection. This algorithm applies a discounted  
289 accuracy function  $\Psi(TOS_i)$  to pick the subset of  $p$  most relevant TOS. The  
290 function is defined as follows:

$$\Psi(TOS_i) = \frac{AUC_i}{\sum_{i,j=1}^k |\rho(TOS_i, TOS_j)|}, \quad (2)$$

291 where  $AUC_i$  is the  $AUC$  performance of the  $i$ -th outlier detection method,  
292 and  $\rho(TOS_i, TOS_j)$  denotes the Pearson correlation coefficient between a  
293 pair of TOS.

### 294 3.2. Machine Learning Methods for Comparative Evaluation

295 In this section, we present the machine learning methods used for com-  
296 parative evaluation in detecting fraud in mobile payment transactions. The  
297 methods can be broadly divided into (1) machine learning methods with  
298 supervised learning that address the class imbalance problem typical for fi-  
299 nancial fraud detection data, and (2) outlier detection methods.

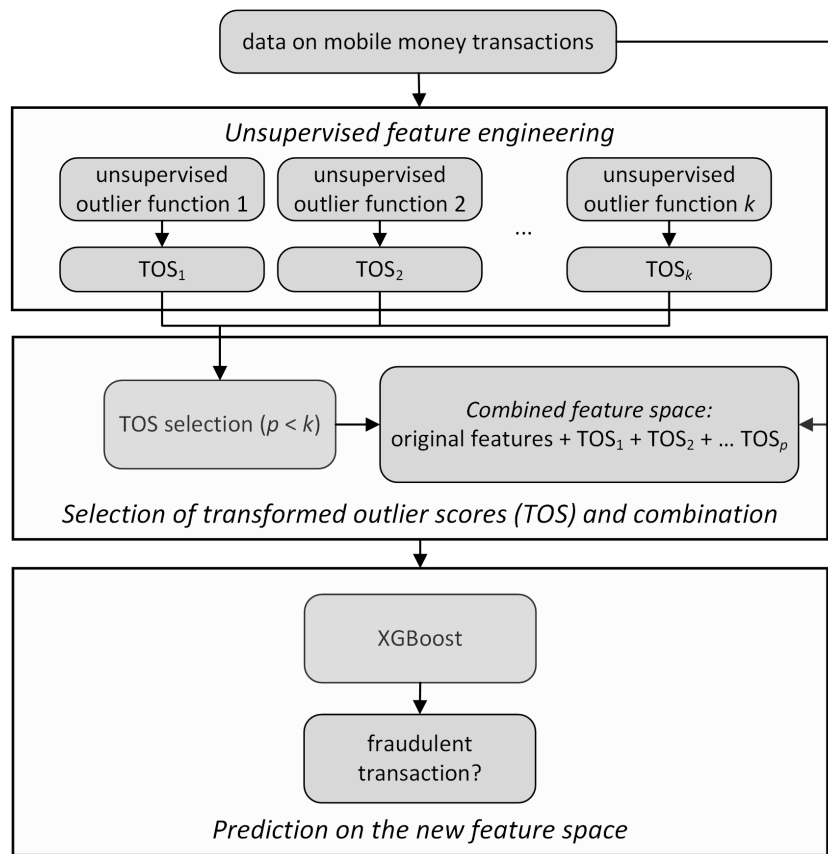


Figure 3: Flowchart of XGBOD for fraud detection

300 *3.2.1. Supervised learning methods for imbalanced data*

301 *k*-nearest neighbour classifier

302 The *k*-nearest neighbour (*k*-NN) method is an instance-based non-parametric  
303 classifier that uses training instances for comparison purpose. An instance is  
304 classified considering its *k* most-similar instances (typically in terms of Eu-  
305 clidean distance) using a majority vote. This simple approach proved to be  
306 accurate in a comparative analysis of machine learning methods for highly  
307 imbalanced credit card fraud detection (Awoyemi et al., 2017). In financial  
308 fraud detection, it is assumed that fraud instances are far from the sam-  
309 ples of the legitimate class. Therefore, *k*-NN can be effectively used even in  
310 unsupervised outlier detection mode (Ramaswamy et al., 2000).

311 *Support vector machine*

312 SVM is a particularly effective classifier for financial fraud detection due  
313 to its capacity to deal with high-dimensional data (Du et al., 2018; Pambudi  
314 et al., 2019; Seera et al., 2021). The SVM algorithm aims to find the optimal  
315 separating hyperplane that maximizes the margin between instances from  
316 different classes. The decision boundary is represented by a subset of the  
317 data known as support vectors. Finding the parameters of the hyperplane is  
318 an optimization problem that takes into consideration both, minimizing the  
319 training error and maximizing the margin. To handle nonlinear relationships  
320 in the data, kernel functions (e.g., linear, polynomial or radial basis functions)  
321 are employed to map the classification problem from the original feature space  
322 to a new feature space of higher dimension where linear separation is possible.

323 *Random forest*

324 Random forest (RF) integrates multiple decision tree predictors trained  
325 independently on different data samples. This allows to generate a number  
326 of trees, ensuring that the generalization error converges to a certain limit.  
327 Another major advantage of RF is its non-differentiable decision boundary.  
328 In addition, random feature selection is used to split the nodes in each tree,  
329 making the RF classifier more robust to noise. The application of RF in  
330 financial fraud detection is particularly effective when the class distribution is  
331 imbalanced because its hierarchical structure enables learning patterns from  
332 both classes (Nami and Shajari, 2018). These advantages explain the good  
333 performance of RF on financial fraud detection tasks (Zhou et al., 2018).

334 *3.2.2. Outlier detection methods*

335 Outlier detection is typically conducted using unsupervised machine learn-  
336 ing methods. The methods presented in this section are trained to represent

337 the legitimate data using clusters of similar data observations. Then, an un-  
338 seen instance is assigned a score that is compared to a threshold representing  
339 the decision boundary separating legitimate instances from outliers.

340 The evaluation conducted in this study contains four types of outlier  
341 detection methods, namely (1) proximity-based methods, (2) linear model-  
342 based methods, (3) ensembling methods, and (4) neural network-based meth-  
343 ods.

#### 344 *Proximity-based methods*

345 To detect outliers, proximity-based methods investigate the neighbour-  
346 hood of each data instance. For example, the local outlier factor (LOF)  
347 method (Breunig et al., 2000) uses the Euclidean distance between the data  
348 instance and its closest neighbour to obtain an outlier score. In the  $k$ -NN  
349 method (KNN) (Ramaswamy et al., 2000), a partition-based algorithm is  
350 first used to identify candidate partitions containing outliers, and then the  
351 distances of instances from these partitions are calculated to detect outliers.  
352 An important advantage of proximity-based methods is their independence  
353 of the data distribution. In other words, no a priori knowledge about the  
354 data distribution is required. However, these methods usually do not scale  
355 well for high-dimensional data. To reduce the sensitivity of LOF to the curse  
356 of dimensionality, the cluster-based local outlier factor (CBLOF) method  
357 (He et al., 2003) replaces closest neighbours with closest clusters, and the  
358 angle-based outlier detection (ABOD) method (Kriegel et al., 2008) replaces  
359 distances with the angular radius and variance of each data vector. The  
360 histogram-based outlier detection (HBOS) method assumes independence of  
361 features to score instances in linear time and is thus computationally more  
362 efficient compared to nearest-neighbour-based methods. However, HBOS  
363 fails in detecting local outliers because the density estimation produced by  
364 histograms does not allow modelling local outliers.

#### 365 *Linear model-based methods*

366 Linear model-based methods rely on the construction of decision bound-  
367 ary separating instances in the legitimate class from the rest of the input  
368 data space. The one-class SVM (OCSVM) method (Schölkopf et al., 2000)  
369 constructs a separating hyperplane in high-dimensional space by minimiz-  
370 ing the structural risk to capture regions of data belonging to the legitimate  
371 class. To prevent overfitting, this method allows a certain percentage of data  
372 instances (regularization parameter) to fall outside the separation bound-  
373 ary. The minimum covariance determinant (MCD) method (Hardin and  
374 Rocke, 2004) combine a multivariate location and scale estimator with a ro-

375 bust clustering algorithm so that the determinant of the covariance matrix  
376 is minimized for each cluster. This method is first trained to fit a minimum  
377 covariance determinant model and then the outlier score is calculated us-  
378 ing the Mahalanobis distance. However, problems can arise when clusters  
379 overlap significantly, leading to poor convergence of the algorithm.

#### 380 *Ensembling methods*

381 Isolation Forest (Liu et al., 2008) aims to separate outliers from the rest of  
382 the data samples. To calculate an isolation score for the data instances, ran-  
383 dom forest is employed. The method assumes that outliers are susceptible to  
384 isolation and, therefore, can be isolated closer to the root of the tree. Specif-  
385 ically, the average path length from the root of the trees can be used obtain  
386 the isolation score. Isolation trees are thus able to build sub-models on dif-  
387 ferent data samples while maintaining low computational complexity and the  
388 ability to scale to handle large volumes of data and high-dimensional prob-  
389 lems. Similarly, lightweight on-line detector of anomalies (LODA) comprises  
390 a collection of weak learners represented by one-dimensional histograms ap-  
391 proximating probabilities of random data projections. The use of sparse pro-  
392 jections makes LODA robust to both the large number of samples and missing  
393 data, allowing the detection of anomalous samples in real-time (Pevny, 2016).

#### 394 *Neural network-based methods*

395 Neural network-based methods utilize feature learning to reduce dimen-  
396 sionality. An autoencoder is an unsupervised neural network capable of non-  
397 linear dimensionality reduction and reproducing input data vectors. Saku-  
398 rada and Yairi (2014) showed that autoencoder (AE) can be successfully  
399 applied to outlier detection. To detect outliers in financial fraud, AEs can be  
400 trained to learn legitimate behaviour and compute a reconstruction error rep-  
401 resenting the outlier score (Sakurada and Yairi, 2014). To achieve robustness  
402 in learning disentangled representations, variational autoencoder (VAE) was  
403 proposed that utilizes both the joint data distribution and their latent gener-  
404 ative factors (Burgess et al., 2018). VAE represents a probabilistic graphical  
405 model whose posterior distribution is estimated using a neural network. The  
406 outlier score of VAE is calculated as the reconstruction probability. Recently,  
407 generative adversarial networks (GANs) have been deployed to unsupervised  
408 outlier detection. Specifically, multi-objective generative adversarial active  
409 learning (MO-GAAL) uses GANs to sample informative potential outliers  
410 following a mini-max game between a discriminator and a generator (Liu  
411 et al., 2019). Thus, GANs assist the discriminative algorithm in finding a  
412 boundary that can effectively separate fraudulent outliers from legitimate

413 normal data. This has been exploited in several studies on financial fraud  
414 (Sethia et al., 2018; Delecourt and Guo, 2019).

### 415 3.3. Performance Evaluation

416 In many related studies (Du et al., 2018; Misra et al., 2020; Mubalaike and  
417 Adali, 2018), the ratio of correctly classified transactions to the total number  
418 of transactions (i.e., accuracy) has been used as the evaluation measure.  
419 However, in the scenario of class-imbalanced data, this measure fails to detect  
420 well the model performance for the minority (fraud) class.

421 As noted in previous research (Lopez-Rojas and Barneaud, 2019), an in-  
422 herent problem in detecting financial fraud that needs to be addressed is the  
423 unknown distribution and impact of all fraudulent transactions. In the ab-  
424 sence of an adequate measure of fraud detection performance, existing fraud  
425 detection approaches rely on traditional measures of classification perfor-  
426 mance. The most desirable performance measure is the ability to correctly  
427 identify fraudulent transactions (true positive rate). In addition, minimizing  
428 false positive and false negative transaction rates (see confusion matrix in  
429 Table 2) is also a desirable quality of fraud detection systems, especially in a  
430 changing fraudulent environment. Here, we use these standard classification  
431 measures to evaluate the performance of fraud detection models. The true  
432 positive rate (*Recall*) is defined as the number of transactions correctly iden-  
433 tified as fraudulent as a percentage of all fraudulent transactions as follows:

$$Recall = \frac{TP}{TP + FN}, \quad (3)$$

434 where  $TP$  and  $FN$  are the numbers of true positive and false negative trans-  
435 actions. The false positive rate ( $FPR$ ) is the number of transactions incor-  
436 rectly identified as fraudulent as a percentage of all legitimate transactions:

$$FPR = \frac{FP}{FP + TN}, \quad (4)$$

437 where  $FP$  and  $TN$  are the numbers of false positive and true negative trans-  
438 actions. The false negative rate ( $FNR$ ) is the number of transactions incor-  
439 rectly identified as legitimate as a percentage of all fraudulent transactions:

$$FNR = \frac{FN}{TP + FN} = 1 - Recall. \quad (5)$$

Table 2: Confusion matrix for fraud detection

Prediction/Target	Positive	Negative
Positive (fraudulent transaction)	$TP$	$FP$
Negative (legitimate transaction)	$FN$	$TN$

440 In reality, financial institutions try to reduce the risk of fraud while trying  
 441 to comply with regulations, but *Recall* is difficult to estimate in the real world  
 442 because *FN* is unknown (hidden fraud). Therefore, financial institutions can  
 443 only calculate *Precision* (i.e., the number of transactions correctly identified  
 444 as fraudulent as a percentage of all transactions that are expected to be  
 445 fraudulent) (Lopez-Rojas and Barneaud, 2019):

$$Precision = \frac{TP}{TP + FP}. \quad (6)$$

446 Previous studies have also considered the *F1* measure (Pambudi et al.,  
 447 2019; Schlör et al., 2021), defined as the harmonic mean of precision and  
 448 recall:

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}. \quad (7)$$

449 The area under the receiver operating characteristic curve (*AUC*) has  
 450 also been used as a more appropriate measure for fraud detection in mobile  
 451 payment transactions due to its robustness to imbalanced data (Buschjäger  
 452 et al., 2021; Mendelson and Lerner, 2020). *AUC* can be defined as the prob-  
 453 ability that a fraud detection model ranks a randomly selected fraudulent  
 454 transaction higher than a randomly selected legitimate transaction, as fol-  
 455 lows:

$$AUC = \int_0^1 Recall(T) \times \frac{d}{dT} FPR(T) dT, \quad (8)$$

456 where  $T$  is the cut-off point.

### 457 3.4. Cost savings measure

458 In addition to the traditional performance measures above, here we pro-  
 459 pose a measure of cost savings measure to account for the financial impli-  
 460 cations of fraud detection models. The proposed cost savings measure was  
 461 inspired by profit-based loan default prediction systems, considering poten-  
 462 tial returns and losses (Papouškova and Hajek, 2019; Ye et al., 2018). On the

463 one hand, correct detection of a fraudulent transaction leads to the following  
 464 cost savings:

$$CS_{TP} = \sum_{i=1}^n (TP_i \times A_i \times 3.36) - (TP_B \times A_F \times 3.36), \quad (9)$$

465 where  $CS_{TP}$  are cost savings from  $TP$  transactions,  $TP_i$  is the  $i$ -th transac-  
 466 tion of  $TP$ ,  $A_i$  is the amount of the  $i$ -th transaction,  $TP_B$  is the number of  
 467  $TP$  transactions detected by the reference fraud detection system, and  $A_F$   
 468 is the average amount of fraudulent transactions. We also took into account  
 469 that fraud now costs financial institutions \$3.36 for every dollar lost to fraud  
 470 and that the current average percentage of successful fraud attempts is 48%  
 471 (i.e.,  $TP_B=0.52$ )<sup>1</sup>.

472 On the other hand, mobile transactions generate a revenue margin of 3.5%  
 473 on average [4]. Therefore, we also considered the cost of  $FP$  transactions,  
 474 estimated as the decrease in the margin for these transactions:

$$Cost_{FP} = (TN \times A_L \times 0.035) - \sum_{j=1}^m (FP_j \times AT_j \times 0.035), \quad (10)$$

475 where  $Cost_{FP}$  is cost of  $FP$  transactions,  $FP_j$  is the  $j$ -th  $FP$  transaction,  
 476  $A_L$  is the average amount of legitimate transactions, and  $AT_j$  is the amount  
 477 of the  $j$ -th transaction. The total cost savings  $CS_{total}$  is then calculated as:

$$CS_{total} = CS_{TP} - Cost_{FP}. \quad (11)$$

478 Note that the proposed measure is expressed in financial terms and is  
 479 instance-dependent (with respect to the amount of each transaction), allow-  
 480 ing for a direct interpretation by financial institutions.

## 481 4. Experimental Results and Analysis

### 482 4.1. Data

483 Consistent with most previous studies (Buschjäger et al., 2021; Du et al.,  
 484 2018; Xenopoulos, 2017), we used the PaySim dataset<sup>2</sup> in this study. The  
 485 main objective of the simulations performed by Lopez-Rojas and his research

<sup>1</sup><https://chainstoreage.com/study-fraud-costs-increased-73-year-over-year-us-retailers>

<sup>2</sup><https://www.kaggle.com/ealaxi/paysim1>

486 team (Lopez-Rojas et al., 2016; Lopez-Rojas et al., 2018; Lopez-Rojas and  
487 Barneaud, 2019) was to replicate typical fraud scenarios that have similar  
488 statistical characteristics to the original mobile payment transaction data. To  
489 this end, different types of fraudulent transactions were injected, including  
490 cash-in (increasing account balance), cash-out (withdrawing cash), payment  
491 (paying for goods or services), transfer (to another user) and debit (sending  
492 money to a bank account). PaySim simulated 743 time steps, representing  
493 thirty days of real-time data. To introduce fraudulent behaviour into the  
494 system, 1,000 fraudsters were included with a 3% probability of committing  
495 fraud at any time step. A total of 6,362,620 mobile transactions were involved  
496 in the dataset, of which 8,213 were fraudulent. Table 3 provides descriptive  
497 statistics of the dataset, and Fig. 4 shows the numbers and amounts of  
498 transactions in time steps.

499 We opted for this dataset for several reasons (Lopez-Rojas and Barneaud,  
500 2019). First, real-time historical data do not include enough fraudulent trans-  
501 actions. Therefore, some previous studies have considered all abnormal trans-  
502 actions to be fraudulent (Choi and Lee, 2017). Second, privacy protections  
503 prevent companies from making datasets public. Third, fraudulent behaviour  
504 is adaptive, making it difficult to create sufficiently diverse real-world fraud  
505 data. In addition, a similar approach based on typical real attack scenarios  
506 was taken in studies related to online banking fraud detection (Carminati et  
507 al., 2015).

Table 3: Attributes in the PaySim dataset

Attribute	Mean value / Range
Step	1-743
Type of transaction	cash-out (35%), cash-in (34%), transfer and debit (31%)
Amount of transaction	180K
Customer name	6.35M unique values
Initial balance	834K
New balance	855K
Recipient name	2.72M unique values
Initial balance of the recipient	1.1M
New balance of the recipient	1.22M
Fraud	0 (legitimate 6.36M) / 1 (fraud 8.2K)

#### 508 4.2. Experimental Setup

509 For data partitioning, we randomly created training and testing data with  
510 a 3:1 ratio (75% training data, 25% testing data). To ensure reliable perfor-  
511 mance evaluation, we repeated this process five times. Since the performance

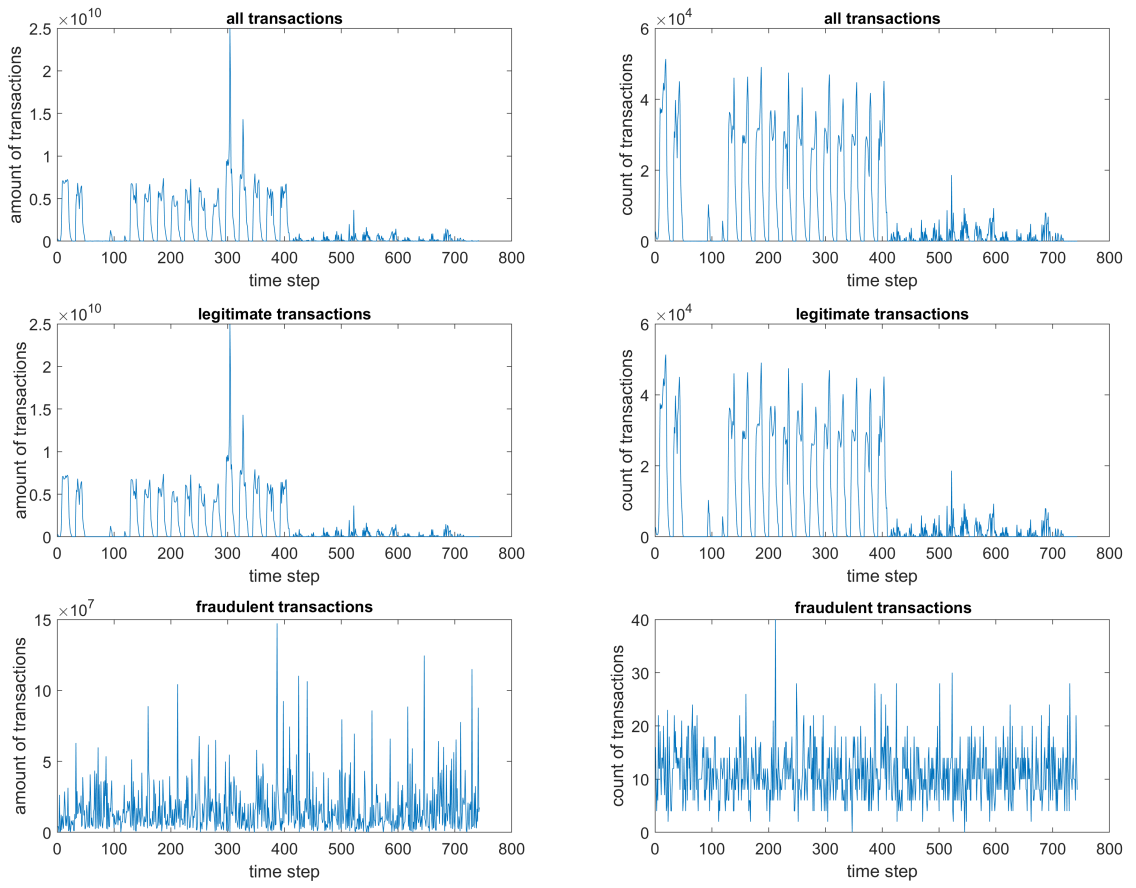


Figure 4: Visualization of amounts and counts of transactions in the PaySim dataset

512 of the fraud detection methods strongly relies on their hyperparameter selec-  
513 tion, we then conducted their optimal selection using 5-fold cross-validation  
514 on the training data (for the list of hyperparameters and their values, see Ap-  
515 pendix 1). Then, we performed fraud detection in mobile payment transac-  
516 tions using the above supervised learning and outlier detection methods. For  
517 experiments, we used the following implementations: (1) supervised learning  
518 methods in the Python library Scikit-Learn 0.23.0, (2) the RUS algorithm  
519 available in the library Imbalanced-Learn 0.6.2, and (3) the outlier detection  
520 methods available in the library PyOD (Zhao et al., 2019). The performance  
521 of the methods was evaluated using the measures defined in the following  
522 subsection.

### 523 4.3. Empirical Results

524 We performed empirical experiments using the PaySim dataset. This  
525 section consists of four subsections. First, we investigate the performance  
526 of supervised learning methods and the effect of random under-sampling on  
527 their effectiveness. Second, the performance of outlier detection methods is  
528 evaluated. Third, the financial consequences of the fraud detection models  
529 are evaluated. Finally, the robustness of the models is tested using a credit  
530 card fraud dataset.

#### 531 4.3.1. Supervised learning methods

532 In the first set of experiments, we compared the performance of four su-  
533 pervised learning methods (XGBoost,  $k$ -NN, SVM, and RF), without using  
534 RUS, to obtain baseline performance. Table 4 shows the testing results of  
535 overall accuracy  $Acc$ ,  $AUC$ ,  $F1$ ,  $Precision$  and  $Recall$ . The values of perfor-  
536 mance measures were obtained as the average of five experiments. For each  
537 performance measure, the number in bold represents the best value among  
538 the tested methods. The non-parametric Wilcoxon test was performed on  
539 the performance measure values obtained in the five experiments to statisti-  
540 cally compare the performance between the best performing method and the  
541 remaining methods. Significantly similar results at the 5% level with respect  
542 to  $AUC$  and  $F1$  are marked with an asterisk.

543 In terms of accuracy, all the supervised learning methods used performed  
544 well. However, as noted above, the extreme class imbalance suggests that this  
545 evaluation measure is not as relevant in this case. As for the  $AUC$  measure,  
546 XGBoost was superior to the other methods, indicating a well-balanced per-  
547 formance for both legitimate and fraud classes. The good balance between

548 *Precision* and *Recall* caused XGBoost to achieve the best results also in  
 549 terms of *F1* measure. By contrast, SVM and *k*-NN performed well only with  
 550 respect to *Precision* and *Recall*, respectively, making them unsuitable meth-  
 551 ods for fraud detection in mobile payment transactions. Overall, these results  
 552 indicate that only XGBoost without class-balancing adjustment is suitable  
 553 for detecting fraud in such a class-imbalanced scenario.

554 Then, we investigated the effect of the RUS under-sampling procedure  
 555 on the performance of the supervised learning methods. On the one hand,  
 556 Table 4 shows that RUS greatly improved the values of *AUC* for SVM, RF  
 557 and XGBoost. On the other hand, there was a considerable deterioration in  
 558 *F1*, which can be attributed to the lower *Precision* achieved at the cost of  
 559 higher *Recall*. In other words, RUS caused almost all fraudulent transac-  
 560 tions to be detected, but this was accompanied by a substantial increase in  
 561 the number of *FP* transactions. This resulted in a bias for the minority class  
 562 while reducing the accuracy for the majority class. It is worth noting that we  
 563 also experimented with other heuristic-based under-sampling methods, such  
 564 as edited nearest neighbour and Tomek links, to address the class imbal-  
 565 ance problem but without improvement in detection performance. Finally,  
 566 it should be noted that the execution time (training time + testing time)  
 567 was substantially reduced by using RUS. For example, RUS+XGBoost was  
 568 computationally most efficient with 2.38 seconds compared to 207.02 seconds  
 569 required for XGBoost without using RUS.

Table 4: Fraud detection performance of supervised learning methods

Method	<i>AUC</i>	<i>F1</i>	<i>Acc</i>	<i>Precision</i>	<i>Recall</i>	<i>Execution time [s]</i>
<i>k</i> -NN	0.9313	0.1588	0.9881	0.0873	0.8744	4,581.4
SVM	0.6543	0.4655	0.9991	<b>0.9474</b>	0.3086	12,082.9
RF	0.8961	0.8394*	0.9996	0.9146	0.7756	1,196.2
XGBoost	0.9350	<b>0.8410*</b>	<b>0.9998</b>	0.8794	0.8059	207.0
RUS+ <i>k</i> -NN	0.8996	0.0405	0.9475	0.0207	0.8516	145.3
RUS+SVM	0.8344	0.0321	0.9431	0.0164	0.7255	1,041.5
RUS+RF	0.9933*	0.2305	0.9914	0.1303	0.9947	12.6
RUS+XGBoost	<b>0.9955*</b>	0.2812	0.9934	0.1637	<b>0.9976</b>	2.4

Notes: The best results are in bold, \* statistically similar at 5% as the best performer in bold. The experiments were performed on Intel® Core™ i5-8400 CPU @ 2.8GHz, 32 GB RAM with six cores on a Windows 10 oper. system in the Python libraries Scikit-Learn 0.23.0 and Imbalanced-Learn 0.6.2.

### 570 4.3.2. Outlier detection methods

571 In the seconds run of experiments, the performance of XGBOD was evalu-  
 572 ated compared with other outlier detection methods. Table 5 shows that XG-

573 BOD significantly outperformed the remaining methods in terms of *AUC* and  
 574 *F1*. In addition, XGBOD was also dominant with respect to both *Precision*  
 575 and *Recall*, indicating excellent performance on both classes.

Table 5: Fraud detection performance of outlier detection methods

Method	<i>AUC</i>	<i>F1</i>	<i>Acc</i>	<i>Precision</i>	<i>Recall</i>	<i>Execution time [s]</i>
ABOD	0.8353	0.0680	0.9953	0.0675	0.0685	2,646.5
CBLOF	0.8593	0.0822	0.9954	0.0829	0.0822	41.3
HBOS	0.7731	0.0077	0.9951	0.0078	0.0076	4.1
LODA	0.6818	0.1060	0.9954	0.1026	0.1096	14.8
Isolation Forest	0.8358	0.0189	0.9964	0.0307	0.0137	189.9
KNN	0.8618	0.1260	0.9957	0.1288	0.1233	1,948.5
MCD	0.7705	0.1084	0.9956	0.1087	0.1081	127.4
OCSVM	0.6732	0.0273	0.9951	0.0272	0.0274	802.9
AE#	0.8050	0.0869	0.9954	0.0870	0.0868	931.1
VAE#	0.8050	0.0869	0.9954	0.0870	0.0868	2,922.9
MO-GAAL	0.9071	0.6059	0.9980	0.5902	0.6225	13,184.4
XGBOD	<b>0.9958</b>	<b>0.8737</b>	<b>0.9994</b>	<b>0.9942</b>	<b>0.7793</b>	4,256.3

Notes: The best results are in bold, # The experiments were performed on GPU NVIDIA GeForce GTX 1060 6GB, 1280 cores on a Windows 10 oper. system in the Python library PyOD.

576 These results can be explained by the semi-supervised learning approach  
 577 used in the XGBOD method. This is because, unlike other outlier detection  
 578 methods, XGBOD leverages the labels assigned to mobile transactions. In  
 579 addition, the transactions contained in the majority class of legitimate trans-  
 580 actions are fully utilized by the multiple unsupervised outlier detection meth-  
 581 ods that produce outlier scores in XGBOD. The XGBoost algorithm applied  
 582 in the improved XGBOD feature space exhibits good robustness to overfit-  
 583 ting and data imbalance, and outperforms the supervised learning methods  
 584 reported in Table 4 in terms of *AUC* and *F1*. However, it should be admit-  
 585 ted that the drawback of XGBOD is the longer execution time, on average  
 586 4,256.25 seconds.

#### 587 4.4. Financial Impact of Fraud Detection

588 To investigate the financial consequences of the evaluated fraud detection  
 589 systems, we used the performance measures defined in Eqs. (9)-(11). Table  
 590 6 shows the average financial performance of all methods in terms of cost  
 591 savings from *TP* transactions, cost of *FP* transactions and total cost savings.  
 592 To calculate these results, we used the average amounts of fraudulent and  
 593 legitimate transactions in the data, i.e.,  $A_F = 1,468,000$  and  $A_L = 178,200$ .

594 In general, supervised learning methods outperformed outlier detection  
 595 methods in terms of overall cost savings, which can be attributed to the high

Table 6: Financial impact of fraud detection methods

Method	$CS_{TP}^*$	$Cost_{FP}$	$CS_{total}$
$k$ -NN	3,576.4	120.3	3,456.1
SVM	-2,135.4	218.3	-2,135.6
RF	2,575.1	923.1	2,574.2
XGBoost	3,630.7	380.5	3,630.3
RUS+ $k$ -NN	3,443.3	519.3	2,924.0
RUS+SVM	2,155.9	561.4	1,594.5
RUS+RF	4,903.3	85.8	4,817.5
RUS+XGBoost	<b>4,932.9</b>	65.9	<b>4,866.9</b>
ABOD	-4,556.9	23.5	-4,580.4
CBLOF	-4,418.7	22.6	-4,441.3
HBOS	-5,171.0	24.2	-5,195.2
LODA	-4,142.3	23.8	-4,166.2
Isolation Forest	-5,109.6	10.7	-5,120.3
KNN	-4,004.2	20.7	-4,024.9
MCD	-4,157.7	22.2	-4,179.7
OCSVM	-4,971.4	24.3	-4,995.7
AE	-4,372.6	22.6	-4,395.3
VAE	-4,372.6	22.6	-4,395.3
MO-GAAL	1,031.6	10.7	1,020.9
XGBOD	2,612.9	<b>0.1</b>	2,612.8

Notes: \* amounts are given in mil. units of an African currency that could not be disclosed by data providers.

596 *Recall* values of supervised learning methods. Note that cost savings from  
597 *TP* transactions were considered to have a stronger financial impact on total  
598 cost savings compared to *FP* transactions. In contrast, XGBOD delivered  
599 the lowest costs associated with *FP* transactions, which is related to its  
600 high *Precision* performance. Surprisingly, SVM and unsupervised outlier  
601 detection methods used in previous studies (Buschjäger et al., 2021; Du et  
602 al., 2018) did not perform well in terms of financial impact and provided  
603 negative overall cost savings due to their low *Recall* values.

#### 604 4.5. Comparison with State-of-the-art Methods

605 To further show the effectiveness of the proposed fraud detection model,  
606 the obtained *AUC* was compared with that of previous studies that examined  
607 the same dataset (Table 7). The best *AUC* performance thus far reported  
608 was achieved using SVM with LogDet regularization (Du et al., 2018). Our  
609 result in Table 4 obtained for SVM confirm that equipping SVM with LogDet  
610 regularization improves the *AUC* performance. Indeed, the traditional SVM  
611 method is reportedly sensitive to outliers and noisy data (Shajalal et al.,  
612 2021). Table 7 also shows that deep neural networks performed well in  
613 previous studies (Schlör et al., 2021; Xenopoulos, 2017). However, their

614 performance is limited by the relatively low number of fraudulent transac-  
615 tions in the dataset. By contrast, the worst performance was reported for  
616 the Isolation Forest method (Buschjäger et al., 2021). Note that the results  
617 for Isolation Forest obtained here (Table 5) are consistent with those from  
618 Buschjäger et al. (2021). The results in Table 7 suggest that the proposed  
619 XGBoost-based models perform better than those used in previous studies in  
620 terms of  $AUC$ , which can be attributed to their good scalability and efficient  
621 processing of sparse data.

Table 7: Comparison of fraud detection performance of the proposed XGBoost-based models with the results of previous studies

Study	Method	AUC
Xenopoulos (2017)	ensemble of deep belief networks	0.961
Du et al. (2018)	SVM with LogDet regularization	0.978
Pambudi et al. (2019)	RUS+SVM	0.880
Mendelson and Lerner (2020)	cluster drift detection	0.898
Schlör et al. (2021)	deep MLP with ReLU and iNALU	0.960
Buschjäger et al. (2021)	generalized Isolation Forest	0.821
This study	RUS+XGBoost	<b>0.996</b>
	XGBOD	<b>0.996</b>

Notes: The best results are in bold.

#### 622 4.6. Robustness Check on Bank Payment Datasets

623 To confirm the obtained performance evaluation, we checked the robust-  
624 ness of the considered fraud detection methods using a bank payment dataset.  
625 The BankSim dataset<sup>3</sup> (Lopez-Rojas and Axelsson, 2014) was generated us-  
626 ing a multi-agent simulation based on a sample of transactional data from a  
627 Spanish bank. The dataset was validated using statistical techniques and so-  
628 cial network analysis of customer-merchant relationships, thus approximating  
629 key features of real bank payment frauds. Each transaction was characterized  
630 by payment amount (in EUR), customer and merchant zip codes, customer  
631 gender and age, and merchant category (e.g., fashion, technology, transport,  
632 and travel). A total of 594,643 transaction records were included, of which  
633 7,200 were fraudulent transaction. The simulation was run for 180 steps  
634 representing months. Thieves were injected to steal or clone an average of  
635 three credit cards at each step and conduct approximately two fraudulent  
636 transactions per day. The result of the simulation is depicted in Fig. 5.

<sup>3</sup><https://www.kaggle.com/ealaxi/banksim1>

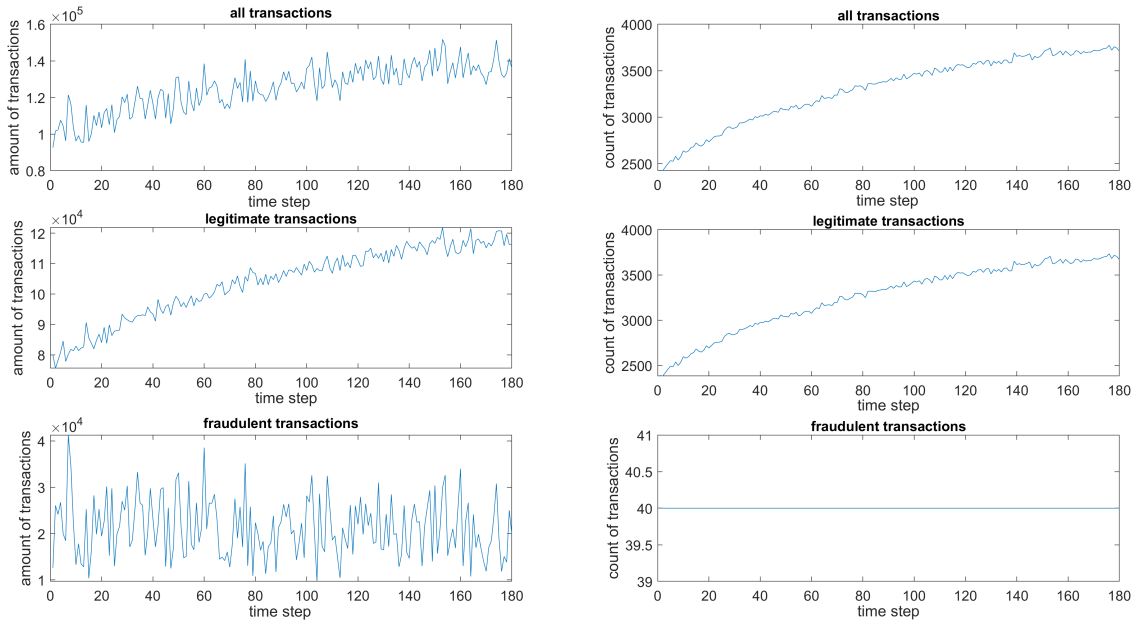


Figure 5: Visualization of amounts and counts of transactions in the BankSim dataset

637 The BankSim dataset provides a benchmark for detecting fraud in bank  
638 payment transactions, as several recent studies have shown (Cui et al., 2021;  
639 Vaughan, 2020). As a robustness check, we trained the evaluated models on  
640 the BankSim dataset using the same experimental setup as for the PaySim  
641 dataset. Note that the sampling process and data collection system was  
642 unique and heterogeneous for both dataset, which allowed us to verify the  
643 robustness of the tested fraud detection models. The results in Table 8 sug-  
644 gest that the under-sampling procedure is not as effective for smaller financial  
645 fraud datasets, improving the performance of supervised learning methods  
646 only in terms of  $AUC$ . In contrast, the performance of unsupervised out-  
647 lier detection methods substantially improved compared to the large PaySim  
648 dataset, suggesting their poor scalability. Overall, XGBoost and XGBOD  
649 performed well in terms of both  $AUC$  and  $F1$  measures, indicating their  
650 good robustness to data size.

#### 651 4.7. Discussion

652 Prior studies (Buschjäger et al., 2021; Pambudi et al., 2019) have noted  
653 the importance of addressing the problem of extreme class imbalance in mo-

Table 8: Fraud detection performance on the BankSim dataset

Method	<i>AUC</i>	<i>F1</i>	<i>Acc</i>	<i>Precision</i>	<i>Recall</i>
<i>k</i> -NN	0.9466	0.2029	0.9089	0.1131	<b>0.9851</b>
SVM	0.7723	0.6849	0.9941	<b>0.9208</b>	0.5451
RF	0.8973	0.8145	0.9957	0.8332	0.7966
XGBoost	0.9112	<b>0.8391</b>	<b>0.9963</b>	0.8240	0.8240
RUS+ <i>k</i> -NN	0.9454	0.3216	0.9535	0.1941	0.9371
RUS+SVM	0.9433	0.3379	0.9571	0.2065	0.9291
RUS+RF	0.9746	0.4674	0.9738	0.3073	0.9754
RUS+XGBoost	0.9774	0.4898	0.9760	0.3266	0.9789
ABOD	0.9852	0.5039	0.9877	0.5032	0.4995
CBLOF	0.9688	0.6072	0.9902	0.6069	0.6074
HBOS	0.9340	0.1490	0.9787	0.1488	0.1488
LODA	0.7272	0.0743	0.9770	0.0736	0.0736
Isolation Forest	0.9647	0.3932	0.9851	0.3974	0.3867
KNN	0.9894	0.5957	0.9899	0.5956	0.5961
MCD	0.9695	0.6922	0.9923	0.6928	0.6944
OCSVM	0.4431	0.0077	0.9758	0.0075	0.0075
AE	0.9350	0.3861	0.9848	0.3829	0.3813
VAE	0.9351	0.3863	0.9848	0.3829	0.3813
MO-GAAL	0.9367	0.3510	0.9807	0.3029	0.4173
XGBOD	<b>0.9968</b>	0.7893	0.9953	0.8018	0.7084

654 bile payment transactions. Therefore, our first set of experiments was de-  
655 signed to investigate the effect of under-sampling the majority class of le-  
656 gitimate transactions on the performance of supervised learning methods.  
657 Consistent with Pambudi et al. (2019), we observed that the detection per-  
658 formance improved for most of the machine learning methods, especially for  
659 the proposed RUS+XGBoost fraud detection model. In contrast to earlier  
660 findings (Buschjäger et al., 2021), however, the second set of experiments  
661 did not detect any evidence for the effectiveness of outlier detection meth-  
662 ods. However, when conducted in a semi-supervised manner, the proposed  
663 XGBOD detection model was found to be superior even to the supervised  
664 learning methods. Finally, the financial consequences of the fraud detection  
665 models were examined to provide guidance on how to set up the right per-  
666 formance metrics for fraud detection in mobile payment transactions. This  
667 experiment addressed the need for an adequate measure of fraud detection  
668 performance as raised in recent research (Lopez-Rojas and Barneaud, 2019).  
669 We found that RUS+XGBoost performed best in terms of cost savings from  
670 correctly detecting fraudulent transactions, while XGBOD minimized the  
671 cost of false positive transactions.

672 Based on the experimental results of this study, we propose the following  
673 suggestions for mobile payment systems.

674 Firstly, the providers of mobile payment systems should pay more at-

675 tention to recent developments in the machine learning research. Specifi-  
676 cally, XGBoost enhanced with class-balancing or outlier detection methods  
677 should be applied to effectively handle the extreme class imbalance prob-  
678 lem in the data and accurately detect fraud in mobile payment transactions.  
679 RUS+XGBoost is particularly recommended for its low execution time, in-  
680 dicating its capability for real-time fraud detection.

681 Secondly, cost savings and transaction costs should be considered when  
682 implementing fraud detection systems in mobile payment systems. For fraud  
683 detection models in mobile payment systems, these evaluation metrics are  
684 critical due to the high costs associated with mobile payment default. The  
685 proposed cost savings measure can be used for this purpose as it offers  
686 providers appropriate guidance for making decisions on the selection cost-  
687 effective fraud detection systems.

688 The importance of accurate and cost-effective fraud detection systems has  
689 dramatically increased during the COVID-19 pandemic because many emerg-  
690 ing and developing countries used mobile money transfer to provide COVID-  
691 19 aid (Blumenstock, 2020). Indeed, mobile payment systems provide a fast  
692 and scalable solution while complying with social distancing measures, which  
693 encouraged government-to-person mobile payments. To enable sustainable  
694 solutions for mobile money transfer, fraud detection technologies represent a  
695 critical component of the frameworks for sustainable government-to-person  
696 mobile money transfers proposed in response to COVID-19 (Davidovich et  
697 al., 2020).

698 Finally, our results suggest that unsupervised outlier detection methods  
699 are not appropriate for fraud detection in mobile payment transactions. The  
700 current study was unable to evaluate the use of the fraud detection system  
701 in a real environment because the number of labelled instances is insuffi-  
702 cient in existing real-world data. Instead, we experimented with a controlled  
703 environment with fraudulent behaviour injected into the data to obtain a  
704 well-performing fraud detection system. However, we believe that the accu-  
705 racy of the proposed fraud detection system would not deteriorate in real-  
706 world applications as the data used in this study are based on the real-world  
707 anonymized data. To further improve the detection accuracy and to assist  
708 the providers of mobile payment systems with the development of fraud de-  
709 tection systems, large labelled real-world data should be collected and made  
710 available to enable effective training of state-of-the-art supervised learning  
711 methods.

## 712 5. Conclusion

713 In this paper, we have proposed an XGBoost-based fraud detection frame-  
714 work while considering the financial impact of fraud detection. The findings  
715 from this study make several noteworthy contributions to the current lit-  
716 erature. First, the XGBoost model was combined with under-sampling to  
717 effectively address the problem of extreme class imbalance and avoid over-  
718 fitting. Second, to fully exploit the large amount of underlying data, unsu-  
719 pervised outlier detection methods were integrated into the XGBoost-based  
720 model. The comparison of the XGBoost-based fraud detection performance  
721 with various state-of-the-art machine learning methods confirmed that we  
722 have found a cutting-edge solution for fraud detection in mobile payment  
723 systems. Our findings also suggest a role for the proposed model in pro-  
724 moting cost savings of fraud detection systems. Taken together, our results  
725 strongly argue against a major role of single machine learning methods and  
726 unsupervised outlier detection methods in fraud detection of mobile payment  
727 transactions, implying that ensemble XGBoost-based methods are preferable.

728 In the future research, ensemble methods combined with alternative under-  
729 sampling and unsupervised outlier detection methods should be further in-  
730 vestigated, including automatic optimization of outlier detection ensembles  
731 (Reunanen et al., 2020) and the XGBoost method enhanced with weighted  
732 and focal losses (Wang et al., 2020). Unfortunately, it was not possible to  
733 investigate our model’s robustness against different mobile payment trans-  
734 action data distributions due to privacy concerns and other limitations of  
735 existing datasets. Therefore, further data would be needed to evaluate model  
736 robustness, including testing the feasibility of transfer learning across mul-  
737 tiple datasets. The proposed fraud detection models should also be applied  
738 to solving related fraud detection problems, such as credit card and loan  
739 frauds, which also exhibit class imbalance characteristics and large real-world  
740 datasets are available for these problems (West and Bhattacharya, 2016).  
741 Other possible application fields of the proposed model include credit scor-  
742 ing (default prediction) (Mahbobi et al., 2021), direct marketing (Wong et  
743 al., 2020), and customer churn prediction (Wong et al., 2020). An issue that  
744 was not addressed in this study was the interpretability property of the fraud  
745 detection models. Therefore, further research might explore the tradeoff be-  
746 tween achieving a high detection accuracy while maintaining interpretability  
747 (Hajek, 2019). Finally, the current investigation was limited by the use of  
748 the cost savings measure only to evaluate the trained model, and thus not in

749 the objective function of the fraud detection model. Future research should  
750 therefore examine the performance of fraud detection models using the cost  
751 savings measure as the objective function. This could lead to our model  
752 delivering even greater cost savings to the end user.

### 753 **Acknowledgements**

754 This article was supported by the scientific research project of the Czech  
755 Sciences Foundation Grant No. 19-15498S.

### 756 **Conflict of Interest**

757 The authors have no competing interests to declare that are relevant to  
758 the content of this article.

### 759 **References**

- 760 [1] Ahmed, M., Mahmood, A.N., & Islam, M.R. (2016). A survey of anomaly  
761 detection techniques in financial domain. *Future Generation Computer  
762 Systems*, 55, 278–288.
- 763 [2] Akanfe, O., Valecha, R., & Rao, H.R. (2020). Assessing country-level pri-  
764 vacy risk for digital payment systems. *Computers & Security*, 99, 102065
- 765 [3] Awoyemi, J.O., Adetunmbi, A.O., & Oluwadare, S.A. (2017). Credit card  
766 fraud detection using machine learning techniques: A comparative anal-  
767 ysis. *IEEE International Conference on Computing, Networking and In-  
768 formatics*, ICCNI 2017, IEEE, p. 1–9.
- 769 [4] Bansal, S., Bruno, P., Denecker, O., & Niederkorn, M. (2019). *Global  
770 Payments Report 2019: Amid Sustained Growth, Accelerating Challenges  
771 Demand Bold Actions*.
- 772 [5] Bernard, P., De Freitas, N.E.M., & Maillet, B.B. (2021). A financial  
773 fraud detection indicator for investors: an IDeA. *Annals of Operations  
774 Research*, 1–24.
- 775 [6] Blumenstock, J. (2020). Machine learning can help get COVID-19 aid to  
776 those who need it most. *Nature* 13.7.2020, 1–3.

- 777 [7] Breunig, M.M., Kriegel, H.P., Ng, R.T., & Sander, J. (2020). LOF: Identifying density-based local outliers. In: *2000 ACM SIGMOD International Conference on Management of Data - SIGMOD '00*, New York, New York, USA, p. 93–104.
- 781 [8] Burgess, C.P., Higgins, I., Pal, A., Matthey, L., Watters, N., Desjardins, G., & Lerchner, A. (2018). Understanding disentangling in  $\beta$ -VAE. In: *Proc. of the 31st Conference on Neural Information Processing Systems*, p. 1–11.
- 785 [9] Buschjäger, S., Honysz, P.J., & Morik, K. (2021). Randomized outlier detection with trees. *International Journal of Data Science and Analytics*, 1–14.
- 788 [10] Carcillo, F., Le Borgne, Y.A., Caelen, O., Kessaci, Y., Oblé, F., & Bontempi, G. (2021). Combining unsupervised and supervised learning in credit card fraud detection. *Information Sciences*, 557, 317–331.
- 791 [11] Carminati, M., Caron, R., Maggi, F., Epifani, I., & Zanero, S. (2015). BankSealer: A decision support system for online banking fraud analysis and investigation. *Computers & Security*, 53, 175–186.
- 794 [12] Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In: *Proc. of the 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, p. 785–794.
- 797 [13] Chen, Y., & Sivakumar, V. (2021). Investigation of finance industry on risk awareness model and digital economic growth. *Annals of Operations Research*, 1–22.
- 800 [14] Chen, S., Yuan, Y., Luo, X.R., Jian, J., & Wang, Y. (2021). Discovering group-based transnational cyber fraud actives: A polymethodological view. *Computers & Security*, 104, 102217.
- 803 [15] Chin, A. G., Harris, M. A., & Brookshire, R. (2022). An empirical investigation of intent to adopt mobile payment systems using a trust-based extended valence framework. *Information Systems Frontiers*, 24, 329–347.

- 807 [16] Choi, D., & Lee, K. (2017). Machine learning based approach to finan-  
808 cial fraud detection process in mobile payment system. *IT CoNvergence*  
809 *PRActice (INPRA)*, 5(4), 12–24.
- 810 [17] Choi, D., & Lee, K. (2018). An artificial intelligence approach to finan-  
811 cial fraud detection under IoT environment: A survey and implementa-  
812 tion. *Security and Communication Networks*, 2018, 5483472.
- 813 [18] Coppolino, L., D’Antonio, S., Formicola, V., Massei, C., & Romano, L.  
814 (2015). Use of the Dempster–Shafer theory to detect account takeovers  
815 in mobile money transfer services. *Journal of Ambient Intelligence and*  
816 *Humanized Computing*, 6(6), 753–762.
- 817 [19] Cui, J., Yan, C., & Wang, C. (2021). ReMEMBeR: Ranking metric  
818 embedding-based multicontextual behavior profiling for online banking  
819 fraud detection. *IEEE Transactions on Computational Social Systems*,  
820 8(3), 643–654.
- 821 [20] David-West, O., Oni, O., & Ashiru, F. (2022) Diffusion of innovations:  
822 Mobile money utility and financial inclusion in Nigeria. Insights from  
823 agents and unbanked poor end users. *Information Systems Frontiers*, 1–  
824 21.
- 825 [21] Davidovic, S., Nunhuck, S., Prady, D., Tourpe, H., & Anderson,  
826 E. (2020). Beyond the COVID-19 crisis: a framework for sustainable  
827 government-to-person mobile money transfers. *IMF Working Papers*,  
828 198, 1–38.
- 829 [22] Delecourt, S., & Guo, L. (2019). Building a robust mobile payment  
830 fraud detection system with adversarial examples. In: *2019 IEEE 2nd*  
831 *Int. Conf. on Artificial Intelligence and Knowledge Engineering (AIKE)*,  
832 IEEE, p. 103–106.
- 833 [23] Dhieb, N., Ghazzai, H., Besbes, H., & Massoud, Y. (2019). Extreme gra-  
834 dient boosting machine learning algorithm for safe auto insurance opera-  
835 tions. In: *2019 IEEE International Conference on Vehicular Electronics*  
836 *and Safety, ICVES 2019, IEEE*, p. 1–5.
- 837 [24] Du, J.Z., Lu, W.G., Wu, X.H., Dong, J.Y., & Zuo, W.M. (2018). L-  
838 SVM: A radius-margin-based SVM algorithm with LogDet regularization.  
839 *Expert Systems with Applications*, 102, 113–125.

- 840 [25] Franque, F. B., Oliveira, T., & Tam, C. (2022). Continuance intention  
841 of mobile payment: TTF model with Trust in an African context. *Information Systems Frontiers*, 1–19.  
842
- 843 [26] Galar, M., Fernandez, A., Barrenechea, E., Bustince, H., & Herrera, F.  
844 (2012). A review on ensembles for the class imbalance problem: Bagging-,  
845 boosting-, and hybrid-based approaches. *IEEE Transactions on Systems,  
846 Man and Cybernetics Part C: Applications and Reviews*, 42(4), 463–484.
- 847 [27] Haixiang, G., Yijing, L., Shang, J., Mingyun, G., Yuanyue, H., & Bing,  
848 G. (2017). Learning from class-imbalanced data: Review of methods and  
849 applications. *Expert Systems with Applications*, 73, 220–239.
- 850 [28] Hajek, P., & Henriques, R. (2017). Mining corporate annual reports for  
851 intelligent detection of financial statement fraud – A comparative study  
852 of machine learning methods. *Knowledge-Based Systems*, 128, 139–152.
- 853 [29] Hajek, P. (2019). Interpretable fuzzy rule-based systems for detecting  
854 financial statement fraud. In: *IFIP International Conference on Artificial  
855 Intelligence Applications and Innovations*, AIAI 2019, Springer, p. 425–  
856 436.
- 857 [30] Hardin, J., & Rocke, D.M. (2004). Outlier detection in the multiple cluster  
858 setting using the minimum covariance determinant estimator. *Computational Statistics and Data Analysis*, 44(4), 625–638.  
859
- 860 [31] He, Z., Xu, X., & Deng, S. (2003). Discovering cluster-based local out-  
861 liers. *Pattern Recognition Letters*, 24(9–10), 1641–1650.
- 862 [32] Huang, S. Y., Lin, C. C., Chiu, A. A., & Yen, D. C. (2017). Fraud  
863 detection using fraud triangle risk factors. *Information Systems Frontiers*,  
864 19(6), 1343–1356.
- 865 [33] Iman, N. (2018). Is mobile payment still relevant in the fintech era?.  
866 *Electronic Commerce Research and Applications*, 30, 72–82.
- 867 [34] Jia, L., Song, X., & Hall, D. (2022). Influence of habits on mobile pay-  
868 ment acceptance: An ecosystem perspective. *Information Systems Fron-  
869 tiers*, 24, 247–266.

- 870 [35] Jocevski, M., Ghezzi, A., & Arvidsson, N. (2020). Exploring the growth  
871 challenge of mobile payment platforms: A business model perspective.  
872 *Electronic Commerce Research and Applications*, 40, 100908.
- 873 [36] Kang, J. (2018). Mobile payment in Fintech environment: trends, secu-  
874 rity challenges, and services. *Human-Centric Computing and Information*  
875 *Sciences*, 8(1), 1–16.
- 876 [37] Kar, A. K. (2021). What affects usage satisfaction in mobile payments?  
877 Modelling user generated content to develop the “digital service usage  
878 satisfaction model”. *Information Systems Frontiers*, 23(5), 1341–1361.
- 879 [38] Kriegel, H.P., Schubert, M., & Zimek, A. (2008). Angle-based outlier  
880 detection in high-dimensional data. In: *Proc. of the 14th ACM SIGKDD*  
881 *International Conference on Knowledge Discovery and Data Mining*, p.  
882 444–452.
- 883 [39] Li, Q., & Clark, G. (2013). Mobile security: A look ahead. *IEEE Security*  
884 *and Privacy*, 11(1), 78–81.
- 885 [40] Liu, F.T., & Ting, K.M., Zhou, Z.H. (2008). Isolation forest. In: *IEEE*  
886 *Int. Conf. on Data Mining*, ICDM, IEEE, p. 413–422.
- 887 [41] Liu, Y., Li, Z., Zhou, C., Jiang, Y., Sun, J., Wang, M., & He, X. (2019).  
888 Generative adversarial active learning for unsupervised outlier detection.  
889 *IEEE Transactions on Knowledge and Data Engineering*, 32(8), 1517–  
890 1528.
- 891 [42] Lopez-Rojas, E., Elmir, A., & Axelsson, S. (2016). Paysim: A financial  
892 mobile money simulator for fraud detection. In: *28th European Model-*  
893 *ing and Simulation Symposium*, EMSS 2016, Dime University of Genoa,  
894 Larnaca, p. 249–255.
- 895 [43] Lopez-Rojas, E.A., & Axelsson, S. (2014). Banksim: A bank payments  
896 simulator for fraud detection research. In: *the 26th European Modeling*  
897 *and Simulation Symposium (EMSS)*, p. 144–152.
- 898 [44] Lopez-Rojas, E.A., Axelsson, S., & Baca, D. (2018). Analysis of fraud  
899 controls using the PaySim financial simulator. *International Journal of*  
900 *Simulation and Process Modelling*, 13(4), 377–386.

- 901 [45] Lopez-Rojas, E.A., & Barneaud, C. (2019). Advantages of the PaySim  
902 simulator for improving financial fraud controls. *Advances in Intelligent*  
903 *Systems and Computing*, 998, 727–736.
- 904 [46] Mahbobi, M., Kimiagari, S., & Vasudevan, M. (2021). Credit risk classi-  
905 fication: an integrated predictive accuracy algorithm using artificial and  
906 deep neural networks. *Annals of Operations Research*, 1–29.
- 907 [47] Mendelson, S., & Lerner, B. (2020). Online cluster drift detection for  
908 novelty detection in data streams. In: *Proc. of the 19th IEEE Interna-*  
909 *tional Conference on Machine Learning and Applications*, ICMLA 2020,  
910 p. 171–178.
- 911 [48] Misra, S., Thakur, S., Ghosh, M., & Saha, S.K. (2020). An autoencoder  
912 based model for detecting fraudulent credit card transaction. *Procedia*  
913 *Computer Science*, 167, 254–262.
- 914 [49] Mubalalike, A.M., & Adali, E. (2018). Deep learning approach for intel-  
915 ligent financial fraud detection system. In: *UBMK 2018 3rd Int. Conf.*  
916 *on Computer Science and Engineering*, p. 598–603.
- 917 [50] Nami, S., & Shajari, M. (2018). Cost-sensitive payment card fraud de-  
918 tection based on dynamic random forest and k-nearest neighbors. *Expert*  
919 *Systems with Applications*, 110, 381–392.
- 920 [51] Ngai, E.W.T., Hu, Y., Wong, Y.H., Chen, Y., & Sun, X. (2011). The  
921 application of data mining techniques in financial fraud detection: A  
922 classification framework and an academic review of literature. *Decision*  
923 *Support Systems*, 50(3), 559–569.
- 924 [52] Onwubiko, C. (2020). Fraud matrix: a morphological and analysis-based  
925 classification and taxonomy of fraud. *Computers & Security*, 96, 101900.
- 926 [53] Pal, A., De, R., & Herath, T. (2020). The role of mobile payment tech-  
927 nology in sustainable and human-centric development: evidence from  
928 the post-demonetization period in India. *Information Systems Frontiers*,  
929 22(3), 607–631.
- 930 [54] Pal, A., Herath, T., De, R., & Rao, H. R. (2021). Is the convenience  
931 worth the risk? An investigation of mobile payment usage. *Information*  
932 *Systems Frontiers*, 23(4), 941–961.

- 933 [55] Pambudi, B.N., Hidayah, I., & Fauziati, S. (2019). Improving money  
934 laundering detection using optimized support vector machine. In: *2019*  
935 *2nd International Seminar on Research of Information Technology and*  
936 *Intelligent Systems*, ISRITI 2019, p. 273–278.
- 937 [56] Papouskova, M., & Hajek, P. (2019). Two-stage consumer credit risk  
938 modelling using heterogeneous ensemble learning. *Decision Support Sys-*  
939 *tems*, 118, 33–45.
- 940 [57] Pevny, T. (2016). Loda: Lightweight on-line detector of anomalies. *Ma-*  
941 *chine Learning*, 102(2), 275–304.
- 942 [58] Ramaswamy, S., Rastogi, R., & Shim, K. (2000). Efficient algorithms for  
943 mining outliers from large data sets. In: *Proc. of the 2000 ACM SIGMOD*  
944 *Int. Conf. on Management of Data*, p. 427–438.
- 945 [59] Reunanen, N., Rätty, T., & Lintonen, T. (2020). Automatic optimization  
946 of outlier detection ensembles using a limited number of outlier examples.  
947 *International Journal of Data Science and Analytics*, 10, 377–394.
- 948 [60] Rieke, R., Zhdanova, M., Repp, J., Giot, R., & Gaber, C. (2013). Fraud  
949 detection in mobile payments utilizing process behavior analysis. In: *2013*  
950 *Int. Conf. on Availability, Reliability and Security*, ARES 2013, p. 662–  
951 669.
- 952 [61] Sakurada, M., & Yairi, T. (2014). Anomaly detection using autoencoders  
953 with nonlinear dimensionality reduction. In: *Proc. of the MLSDA 2014*  
954 *2nd Workshop on Machine Learning for Sensory Data Analysis*, p. 4–11.
- 955 [62] Seera, M., Lim, C.P., Kumar, A., Dhamotharan, L., & Tan, K.H. (2021).  
956 An intelligent payment card fraud detection system. *Annals of Operations*  
957 *Research*, 1–23.
- 958 [63] Sethia, A., Patel, R., & Raut, P. (2018). Data augmentation using gener-  
959 ative models for credit card fraud detection. In: *4th International Con-*  
960 *ference on Computing Communication and Automation (ICCCA)*, IEEE,  
961 p. 1–6.
- 962 [64] Schlör, D., Ring, M., Krause, A., & Hotho, A. (2021). Financial fraud  
963 detection with improved neural arithmetic logic units. *Lecture Notes in*  
964 *Computer Science*, 12591, 40–54.

- 965 [65] Schölkopf, B., Williamson, R., Smola, A., Shawe-Taylor, J., Platt, J.,  
966 & Holloway, R. (2000). Support vector method for novelty detection. In:  
967 *Advances in Neural Information Processing Systems*, MIT Press, p. 582–  
968 588.
- 969 [66] Shajalal, M., Hajek, P., & Abedin, M. Z. (2021). Product backorder  
970 prediction using deep neural network on imbalanced data. *International*  
971 *Journal of Production Research*, 1–18.
- 972 [67] Turner, A., McCombie, S., & Uhlmann, A. (2021). Follow the money:  
973 Revealing risky nodes in a Ransomware-Bitcoin network. In: *Proc. of the*  
974 *54th Hawaii Int. Conf. on System Sciences*, p. 1560–1572.
- 975 [68] Vaughan, G. (2020). Efficient big data model selection with applications  
976 to fraud detection. *International Journal of Forecasting*, 36(3), 1116–  
977 1127.
- 978 [69] Verkijika, S.F. (2020). An affective response model for understanding the  
979 acceptance of mobile payment systems. *Electronic Commerce Research*  
980 *and Applications*, 39, 100905.
- 981 [70] Wang, C., Deng, C., & Wang, S. (2020). Imbalance-XGBoost: leveraging  
982 weighted and focal losses for binary label-imbalanced classification with  
983 XGBoost. *Pattern Recognition Letters*, 136, 190–197.
- 984 [71] West, J., & Bhattacharya, M. (2016). Intelligent financial fraud detec-  
985 tion: A comprehensive review. *Computers & Security*, 57, 47–66.
- 986 [72] Wong, M. L., Seng, K., & Wong, P. K. (2020). Cost-sensitive ensemble of  
987 stacked denoising autoencoders for class imbalance problems in business  
988 domain. *Expert Systems with Applications*, 141, 112918.
- 989 [73] Xenopoulos, P. (2017). Introducing DeepBalance: Random deep belief  
990 network ensembles to address class imbalance. In: *2017 IEEE Int. Conf.*  
991 *on Big Data, Big Data 2017*, p. 3684–3689.
- 992 [74] Yamanishi, K., Takeuchi, J.I., Williams, G., & Milne, P. (2004). On-  
993 line unsupervised outlier detection using finite mixtures with discounting  
994 learning algorithms. *Data Mining and Knowledge Discovery*, 8(3), 275–  
995 300.

- 996 [75] Ye, X., Dong, L.A., & Ma, D. (2018). Loan evaluation in P2P lending  
997 based on random forest optimized by genetic algorithm with profit score.  
998 *Electronic Commerce Research and Applications*, 32, 23–36.
- 999 [76] Zhao, Y., & Hryniewicki, M.K. (2018). XGBOD: Improving supervised  
1000 outlier detection with unsupervised representation learning. In: *Proc. of*  
1001 *the Int. Joint Conf. on Neural Networks*, p. 1–8.
- 1002 [77] Zhao, Y., Nasrullah, Z., & Li, Z. (2019). PyOD: A Python toolbox for  
1003 scalable outlier detection. *Journal of Machine Learning Research*, 20(96),  
1004 1–7.
- 1005 [78] Zhou, H., Chai, H.F., & Qiu, M.L. (2018). Fraud detection within  
1006 bankcard enrollment on mobile device based payment using machine  
1007 learning. *Frontiers of Information Technology and Electronic Engineer-*  
1008 *ing*, 19(12), 1537–1545.

## Appendix 1: Settings of machine learning methods

Method	Parameters
ABOD	contamination = the proportion of frauds in the training dataset, neighbours $k = \{5, 10\}$
CBLOF	number of clusters = 8, clustering estimator = $K$ -means, alpha = 0.9
HBOS	alpha = 0.1
LODA	number of bins = 10, number of random cuts = 100
Isolation Forest	number of estimators = $\{100, 200\}$
KNN	neighbours $k = \{2, 3, 5\}$ , radius = 1.0
MCD	contamination = the proportion of frauds in the training dataset
OCSVM	kernel function: {linear, polynomial, RBF with gamma = 0.01}, nu = 0.1
AE	hidden activation = ReLU, optimizer = adam, epochs = 100, dropout rate = 0.2, L2 regularizer = 0.2, hidden neurons = [8, 4, 4, 8]
VAE	hidden activation = ReLU, optimizer = adam, epochs = 100, gamma = 1.0, dropout rate = 0.2, L2 regularizer = 0.1, encoder neurons = [8, 4, 2], decoder neurons = [2, 4, 8]
MO-GAAL	contamination = the proportion of frauds in the training dataset, number of sub generators = 10, learning rate of the discriminator = 0.01, learning rate of the generator = 0.0001, epochs = 20
XGBOD	estimator list = {ABOD, CBLOF, HBOS, LODA, Isolation Forest, KNN, MCD, OCSVM, AE, VAE}, $p = 5$ , learning rate = 0.1
$k$ -NN	$k = 2, 3, 5$
SVM	complexity parameter $C = 1$ , kernel function: {linear, polynomial, RBF with gamma = 0.01}
RF	number of trees = {100, 200}
XGBoost	booster = gbtree, eta = 0.3, gamma = 0, maximum depth of a tree = {3, 6, 9}, sampling method = uniform, lambda = 1, alpha = 0
RUS	sampling strategy = {0.5, 0.75, 1.0}