

Fraud Detection for E-commerce Transactions by Employing a Prudential Multiple Consensus Model

Salvatore Carta^a, Gianni Fenu^a, Diego Reforgiato Recupero^a, Roberto Saia^a

^a*Department of Mathematics and Computer Science, University of Cagliari
Palazzo delle Scienze, Via Ospedale 72, 09124 Cagliari*

Abstract

More and more financial transactions through different E-commerce platforms have appeared now-days within the big data era bringing plenty of opportunities but also challenges and risks of stealing information for potential frauds that need to be faced. This is due to the massive use of tools such as credit cards for electronic payments which are targeted by attackers to steal sensitive information and perform fraudulent operations. Although intelligent fraud detection systems have been developed to face the problem, they still suffer from some well-known problems due to the imbalance of the used data. Therefore this paper proposes a novel data intelligence technique based on a Prudential Multiple Consensus model which combines the effectiveness of several state-of-the-art classification algorithms by adopting a twofold criterion, probabilistic and majority based. The goal is to maximize the effectiveness of the model in detecting fraudulent transactions regardless the presence of any data imbalance. Our model has been validated with a set of experiments on a large real-world dataset characterized by a high degree of data imbalance and results show how the proposed model outperforms several state-of-the-art solutions, both in terms of ensemble models and classification approaches.

Keywords: Information Security, Credit Card, Fraud Detection, Machine Learning

1. Introduction

Nowadays, the employment of credit cards for financial transactions represent the backbone of the E-commerce dynamics and business, since they allows purchasing in real-time of goods and services all over the world and using any device (smart-phone, 5 tablet, pc) connected to the Internet. As it can be noticed, there are risks associated to this operation that might cause the theft of sensitive information associated to the

*Corresponding author

Email address: diego.reforgiato@unica.it (Diego Reforgiato Recupero)

customers' credit cards. A recent report from the European Payments Council¹ shows that a certain percentage of the Internet electronic payments is related to frauds.

One more analysis has been performed by the *Euromonitor International*² in the
10 Europe, Middle East and Africa (*EMEA*) area, which shows that the number of frauds, and the associated budget in euros, within the *EMEA* area kept growing from 2006 to 2016, the year of the publication of the study. The values showing that are reported in Figure 1. Although these data refer only to the *EMEA* area, they clearly underline the seriousness of the problem. In US, the *American Association of Fraud Examiners*³
15 found out that 15% of all the frauds are somehow connected to credit cards transactions, and this represents the 80% of the whole financial value.

According to the *FBI's* Internal Crime Complaint Center (*IC3*)⁴, the term *credit card fraud* is defined as:

20 **“a wide-ranging term for fraud committed using a credit card or any similar payment mechanism as a fraudulent source of funds in a transaction.”**

It can be carried out in two different ways, *off-line* or *on-line* [1]. If the fraud is off-line, that means the credit card has been previously stolen and then used to perform fraudulent payments, assuming the identity of the legitimate owner. For this case the thief has a limit amount of time which lasts from the time of the theft of the credit card
25 to the time when the owner reports to his/her bank and the bank consequently disables the card. When the fraud is on-line, and this is the most common, the information been stolen is digital and it has been obtained in several ways (e.g., *skimming*, *shimming*, *cloning*, or *phishing*). Once the fraudster obtains this information, he/she can purchase through the Internet, until either the legitimate owner does not notice the problem and
30 blocks the card or the budget in the card ends. The latter (fraud on-line) is the case we take into account within the proposed paper.

Several research institutions and industries have made huge investments with the aim of designing effective methods capable of tackling the problem by employing machine learning, deep learning, big data, and computational intelligence technologies.

35 The efforts in this context have led to a large number of solutions that are able to automatically distinguish legitimate credit card transactions from the fraudulent ones.

However, regardless of the used approach, there are some common problems that reduce its performance. The most common is represented by the unbalanced distribution nature of the training data characterizing the past transactions which generates
40 different problems of over-fitting and leads to low performances of the adopted classifiers. In other words, such a problem arises because the number of available *fraudulent* samples is usually much lower than the *legitimate* ones and this high grade of unbalance does not allow the definition of a reliable model of evaluation [2, 3, 4].

45 This happens because the fraudulent transactions collected in the past by the fraud detection systems are much less frequent than the legitimate ones. Moreover, (i) the heterogeneity nature of the data and the (ii) presence of overlaps among the data [5]

¹<https://bit.ly/2yQC7G1>

²<http://www.euromonitor.com/>

³<http://www.acfe.com>

⁴<https://www.ic3.gov/>

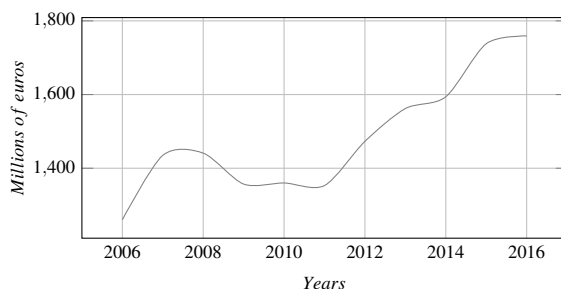


Figure 1: *Fraud Amount in EMEA Area*

are two elements that worsen the problem.

The two elements mentioned above highly affect the effectiveness of any fraud detection systems generating a large number of miss-classifications.

50 The approach we propose in this paper aims at improving the classification performances of several state-of-the-art classification algorithms, by adopting an ensemble approach where the final classification is given by the combination of the different elements of the ensemble through a novel model regulated by a twofold policy. The policy is further defined with probabilistic and prudential criteria in order to maximize the effectiveness of single approaches.

55 More in detail, the main scientific contributions of our proposed approach are the following:

- (i) we introduce a formalization of the *Prudential Multiple Consensus (PMC)* model aimed at combining the classification made by each single approach by adopting both a probabilistic and a prudential criterion;
- 60 (ii) we defined the algorithm used to classify the new transactions as *legitimate* or *fraudulent* depending on the *PMC* model previously formalized, according to the performed ensemble criteria analysis.

65 The remainder of this paper is organized as it follows. Section 2 introduces background and related work of the scenario taken into account. Section 3 includes the adopted formal notation defining formally the problem we face. In Section 4 we describe the implementation of the proposed approach whereas in Section 5 we perform a preliminary study aimed at selecting and ensembling a set of classification algorithms. Section 6 describes the characteristics of the experimental environment as well as the description of the adopted datasets, strategy, and metrics whereas Section 7 shows the obtained results along with a related discussion. Remarks and future work where we are headed are given in Section 8 which also ends the paper.

2. Background and Related Work

75 After an introduction on the most common approaches and methods used to tackle the fraud detection problem, this section underlines the current open problems, introducing the ensemble classification methods and the most suitable performance metrics that are used for evaluation.

2.1. Fraud Detection Approaches

Different approaches in literature which tackle the fraud detection problem exploits
80 the following techniques:

- *Data Mining*: an example is provided by the work of researchers in [6] where the generation of ad-hoc patterns are presented to recognize frauds. In one more example [7], authors investigate several combinations of manual and automatic approaches of classification.
- 85 • *Artificial Intelligence*: as in [8], which uses a technique to obtain a reduction of the number of false alarms, during the evaluation process.
- *Machine Learning*: the work presented in [9] makes use of several types of classification (single and ensemble). Another work, [10], takes into account the combination of unsupervised and supervised strategies.
- 90 • *Genetic Programming*: an example is represented by [11]), where an *evolutionary computation* technique has been implemented to improve the fraud detection process, taking into account the dynamics of the credit card transactions.
- *Reinforcement Learning*: as in [12], which formalizes the interactions between fraudsters and card-issuers as a *Markov Decision Process*.
- 95 • *Transformed-domain-based*: as in [13, 14, 15, 16], where the evaluation process has been performed in a non-canonical domain (e.g., time, frequency, or frequency-time).
- *Combined Criteria*: as in [17, 18], where a multidimensional technique is exploited to improve the classification performance. One more example is given
100 in [17], where the authors introduce several fraud indicators in the classification process.

2.2. Operating Modalities

Under a different point of view, all the state-of-the-art fraud detection approaches can be divided in two categories which depend on the involvement of artificial intelligence technology [19]: there are *supervised approaches* and *unsupervised approaches*.
105 In the following we will show a brief description of each of the two classes.

- *Supervised approaches* define their evaluation model by exploiting the past fraudulent and non-fraudulent transactions collected by the fraud detection system. These methods do not work well without a considerable number of examples
110 (training data) which have annotations in both classes (i.e., *legitimate* and *fraudulent*). As these methods must learn from the training data how to predict future data, and therefore they discover only known patterns, it is important to provide them with a fully consistent and complete training data.
- *Unsupervised approaches* operate by searching anomalies in the features that
115 compose the transaction under evaluation. The problem in this case is that a *fraudulent* transaction might not have any anomalies in its values and, consequently, the design of effective unsupervised fraud detection approaches continues being a hard research challenge [20].

In addition, the evaluation model definition can be performed using three different
120 modes: *static*, *updating*, or *forgetting*:

- by following the *static mode*, the data under analysis are divided into several blocks of equal size and the training of the model is made by exploiting a defined number of contiguous blocks [21]. A drawback of this mode is the absence of a dynamic model of evaluation able to follow users behaviour changes.
- 125 • the *updating mode* does not work by using a unique evaluation model, since it updates the model when a new block arrives, involving in this process a defined number of the most recent and contiguous blocks [22]. The problem in this case is the impossibility to operate with small classes of data.
- 130 • the *forgetting mode* also updates the evaluation model at each new block, but it performs this operation by involving all the past *fraudulent* transactions and the *legitimate* ones present in the last two blocks [23]. However, this mode presents high computational cost.

2.3. Current Open Problems

In addition to the intrinsic issues mentioned in Section 2.2, in the following we report the most important open problems that affect the fraud detection domain, regardless of the used approach.

- **Data Scarcity:** it happens because for different reasons (commercial operators policies, privacy, legal constraints, etc.) there is not much availability of real-world datasets to use to develop and verify novel fraud detection approaches [11]. This scenario is quite understandable, given the intrinsically private nature of the involved data and it represents a big problem for the researchers, which in many cases are forced to use synthetic data [24].
- 140 • **Data Heterogeneity:** this problem is related to the difficulty to model the relationships among transaction features that are represented differently in various sets of data [25]. In other words, it is presented because each card issuer processes every day a high number of transactions, and each transaction is composed by a number of features whose values periodically might change within a single user's account or between all the accounts.
- 145 • **Model Staticity:** the classification approaches define their evaluation model based on the available data (i.e., past transactions). Considering the high level of heterogeneity of the involved information, this is a problem where the pattern that characterizes a new transaction is not present among those used to define the evaluation model [26].
- 150 • **Cold Start:** in the fraud detection scenario and in others needing the training step of an evaluation model, this problem happens when the available data are not enough [27]. In the context taken into consideration in this paper, we have a cold-start situation when a fraud detection system does not collect a sufficient number of *fraudulent* and *legitimate* transactions to perform the model training.
- 155 • **Data Imbalance:** this issue is given by the small number of *fraudulent* cases usually collected by a fraud detection system, respect to the *legitimate* cases. Considering the past transactions are used to train the evaluation models, such an occurrence leads towards a reduction of the fraud detection approaches effectiveness [2, 3, 4]. The literature presents many approaches able to face this problem, such as those proposed in [28, 29].
- 160

165 *2.4. Ensemble Methods*

Ensemble methods are largely adopted in order to perform data classification [30], since they can improve the performance achieved by a single classification method. Such methods have been much investigated in the past, as for example the work in [31], which evaluates the advantages related to the computational, statistical, and representational aspects.

170 It should be observed how the combination of more classification algorithms does not always lead to better results, because this operation can also reduce the overall classification performance. This means that the effectiveness of the resulting approach is strictly related to the strategy adopted to aggregate the single results, thus it highly depends on the definition of the global assessment model [30].

175 In scenarios such as the one considered in this paper, the ensemble methods are mainly aimed at improving the correct evaluation of a minority class label (e.g., that related to the *fraudulent* cases), since it represents an important result when the available data are strongly unbalanced [32].

180 The literature indicates the ensemble methods as one of the most effective approaches able to face the class imbalance problems. Moreover, such a scenario has been well outlined in a survey [33], where out of 527 papers taken into consideration, 218 referred to ensemble models.

185 The strategy used to combine several classification algorithms usually operates in the following two steps:

- (i) in the first step a series of different algorithms is selected based on the complementarity of their results (i.e., they get misclassifications in different places of the test set);
- (ii) in the second step their results are combined by adopting a consensus criterion, such as complete agreement, majority, absolute, correction, multi stage, weighted, confidence and ranked voting [34].

190 Approaches similar to ours are listed in the following: [35], where the authors analyzed the performance of three state-of-the-art data mining techniques in the context of a bagging ensemble classifier based on *decision tree* algorithms; [36], where the authors propose a strategy that drops a certain number of classifiers periodically and uses only a part of them for the evaluation; [37], where the authors combine the bagging and boosting techniques. However, regardless the adopted ensemble strategy, the state-of-the-art solutions do not implement any prudential criterion, such as that in our *PMC* approach, which is based on the observation that in the context taken into account (credit card fraud detection) a wrong classification of a transaction as *fraudulent* is preferable rather than a wrong classification of a transaction as *legitimate*.

200 The reader notices that some classification algorithms are ensemble in nature, such as *AdaBoost* [38] and *Gradient Boosting* [39]. Such algorithms operate by exploiting a prediction model based on an ensemble of weak prediction models [33].

205 *2.5. Performance Assessment*

Within the fraud detection scenario, especially where credit card transactions are involved, the performance assessment must follow certain criteria, due to the particular

configuration of the involved data. This is necessary because some canonical metrics (e.g. the accuracy) usually used to evaluate the classification algorithms performance might lead to unreliable results.

This happens especially when the experiments involve unbalanced data [3, 33, 40]. For example, let us assume a dataset in which the *fraudulent* cases represent the 0.01% of the entire dataset (*legitimate* and *fraudulent*), an algorithm that classifies all the samples as *legitimate* achieves the 99% of accuracy.

It should be also underlined that such an event is not rare, and it is in accordance with the real-world data.

Table 1: *Confusion Matrix*

		Algorithm classification		total
		fra	leg	
Real class	fra'	True Positive	False Negative	$ fra' $
	leg'	False Positive	True Negative	$ leg' $
total		$ fra $	$ leg $	

For the aforementioned considerations, the suitable assessment metrics should be oriented to evaluate the algorithms performance by taking into account the unbalanced configuration of data.

For this reason it is preferable to use the metrics based on the *confusion matrix* shown in Table 1 (avoiding their use in aggregate form), where *fra* stands for *fraudulent* and *leg* stands for *legitimate*.

Simple metrics as the *Sensitivity* (true positive rate) and the *Fallout* (false positive rate) give us information about the classification algorithms effectiveness in terms of *fraudulent* cases correctly classified, while metrics as the *Specificity* (true negative rate) and the *Miss Rate* (false negative rate), provide specular information on the performance related to the detection of the *legitimate* cases.

In addition to the aforementioned metrics, it is also preferable to add another one such as the *AUC* (Area Under the *ROC* Curve), since it is able to investigate the ability to discriminate between the possible destination classes (i.e., *legitimate* and *fraudulent* in our case) of the adopted evaluation model [41, 42].

3. The Proposed Approach, Notation and Problem Formulation

This section introduces the proposed approach, the formal notation which includes the formulation of the problem we address in this paper.

3.1. Proposed Approach

This paper proposes a combined approach aimed at maximizing the effectiveness of several single approaches. We employ an ensemble strategy regulated by a *Prudential*

Multiple Consensus model, which is based on a twofold criterion, probabilistic and majority based.

240 Such an idea relies on the observation that the results given by different classification algorithms are not the same in terms of correct classifications and misclassifications. It means that in many cases their results do not agree on the identification of certain *legitimate* or *fraudulent* transactions.

245 In more detail, in a preliminary study we observed that the classifications made by different algorithms are frequently in conflict, also when the different algorithms achieve good individual performances in terms of *Specificity* (i.e., legitimate cases correctly classified) and *Sensitivity* (*fraudulent* cases correctly classified). This can be exploited to increase the classification reliability, by adopting strategies that take into account the classifications made by multiple algorithms.

250 On the basis of these considerations, the proposed approach wants to exploit an aggregation strategy able to conveniently combine the correct evaluations made by the single algorithms, maximizing their effectiveness in the *fraudulent* transactions detection.

255 Instead of adopting a canonical aggregation criteria (e.g., *complete agreement*, *majority voting*, or *weighted voting*, etc.) to determine the class of destination of a new transaction and the basis of the results of the single algorithms, our approach adopts a novel *prudential criterion* which works as it follows:

- 260 (i) each single algorithm classifies a new transaction as *legitimate* only if its classification is *legitimate* and the classification probability is above the average value of the probabilities of the classifications made by all the algorithms for that transaction. When this does not happen (i.e., the algorithm classification is *fraudulent* or the classification probability is below that average value of probability) the transaction is classified as *fraudulent*;
- 265 (ii) a canonical consensus criterion based on the *majority voting* is then taken into account and the final classification of the transaction under analysis will depend on the results of all the algorithms.

3.2. Formal Notation

270 Given a set of transactions $T = \{t_1, t_2, \dots, t_N\}$ collected in the past and already classified and the subsets $T_+ = \{t_1, t_2, \dots, t_K\}$ and $T_- = \{t_1, t_2, \dots, t_J\}$, respectively related to the *legitimate* and *fraudulent* transactions in T (i.e., $T_+ \subseteq T$ and $T_- \subseteq T$), we denote as $F = \{f_1, f_2, \dots, f_M\}$ the set of features that compose each transaction $t \in T$.

In addition, we denote as $\hat{T} = \{\hat{t}_1, \hat{t}_2, \dots, \hat{t}_U\}$ a set of new transactions to classify and as $C = \{\textit{legitimate}, \textit{fraudulent}\}$ the possible classes of destination of each transaction, meaning that a transaction can belong to only one class $c \in C$.

275 Finally, we denote as $A = \{a_1, a_2, \dots, a_Z\}$ a set of classification algorithms.

Let Ψ be the classification process made by using our *PMC* model. Then we evaluate the correctness of each classification performed by *PMC* through the function *Evaluation*(\hat{e}, Ψ) that returns a Boolean value β : 1 in case of correct classification, 0

280 otherwise. In this way we can formalize our classification problem in terms of maximization of the sum of the values returned by this function, as indicated in Equation 1.

$$\max_{0 \leq \beta \leq |\hat{E}|} \beta = \sum_{u=1}^{|\hat{E}|} Evaluation(\hat{e}_u, \Psi) \quad (1)$$

4. Implementation

The architecture of our approach is shown in Figure 2.

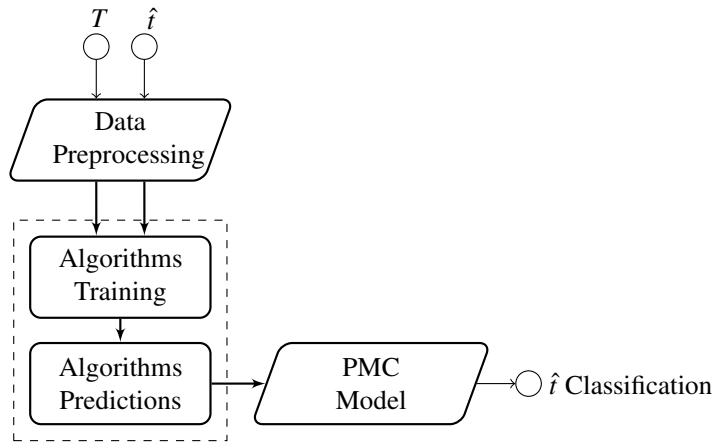


Figure 2: *PMS High-level Architecture*

285 In such an architecture, the activity made in the *Data Preprocessing* block depends on the input data and it can involve, for instance, a *class binarization* (a mapping of a multi-class learning problem to several two-class learning problems) or a *minority class oversampling* (in order to adjust the class distribution of the dataset).

4.1. Step 1: Model Definition

290 The proposed *Prudential Multiple Consensus (PMC)* model operates by combining the results of five classification algorithms (described in Section 6.3) on the basis of two criteria, one based on the *classification probability* and one based on the *majority voting*.

295 To get the *classification probability* we use *Logistic Function*, since it is able to measure the probability of a binary response based on more independent predictors. More formally, the probability that a new transaction $\hat{t} \in \hat{T}$ belongs to a class $c \in C$ is calculated by mapping the algorithm predictions in terms of probabilities through the sigmoid σ function⁵. Such a method is formalized in the Equation 2, where $\sigma(a_z(p))$

⁵A mathematical function characterized by a *sigmoid curve*, which maps any real value into the interval $[0,1]$.

is the probability estimate for the prediction p through the algorithm a_z , whose result is given in the range $[0, 1]$ and e denotes the base of natural log.

$$\sigma(a_z(p)) = \frac{1}{1 + e^{-p}} \quad (2)$$

Subsequently, in each classification performed by a single algorithm, a transaction is considered as *legitimate* only when its probability is above a certain value, otherwise, prudentially, the transaction is classified as *fraudulent*. The final classification is given according to the results of all the algorithms (by using the *majority voting* criterion), as shown in Equation 3, where $|A|$ is the number of classification algorithms and c is the transaction classification.

$$c = \begin{cases} \textit{legitimate}, & \textit{if } w_1 > w_2 \\ \textit{fraudulent}, & \textit{otherwise} \end{cases}$$

with

$$\mu = \frac{1}{|A|} \cdot \sum_{z=1}^{|A|} \sigma(a_z(p)) \quad (3)$$

$$w_1 = \sum_{z=1}^{|A|} 1 \textit{ if } \sigma(a_z(p)) > \mu \wedge a_z(p) = \textit{legitimate}$$

$$w_2 = \sum_{z=1}^{|A|} 1 \textit{ if } \sigma(a_z(p)) \leq \mu \vee a_z(p) = \textit{fraudulent}$$

It should be observed that *Logistic Function* represents only one of the possible approaches able to estimate the probability of a binary response given by a predictor. It means that also other approaches able to perform the same operation can be used in our model.

4.2. Step 2: Data Classification

According to the model previously formalized in Section 4.1, each new transaction $\hat{t} \in \hat{T}$ is classified by using the Algorithm 1.

The input of the Algorithm 1 is represented by the classification algorithms in the set A , the previous transactions E already classified, and a new transaction $\hat{t} \in \hat{T}$ to evaluate. The output will be the classification of the event \hat{t} as *legitimate* or *fraudulent*. At *step 4* the evaluation models related to the set A of the classification algorithms are defined, while their classifications for the transaction \hat{e} are calculated at *step 5*. The average probability value of all the performed classifications is calculated at *step 6* and saved in μ . A control aimed at checking whether the classification probability of each algorithm is above the average value in μ is performed (*step 7 to step 13*). In particular, we increase by one the value of w_1 when $p = \textit{legitimate}$ and the prediction probability is above the μ value, otherwise we increment the value of w_2 . The transaction \hat{e} is classified as *legitimate* when all predictions have been processed and $w_1 > w_2$, otherwise the transaction is classified as *fraudulent*. Such a classification is returned, and the algorithm ends, at *step 19*.

It should be noted that the functions *getProbabilityAverage()* and *getProbability()* are both based on the *Logistic Function* model formalized in Equation 2.

Algorithm 1 *Transaction classification*

Input: A =Set of algorithms, T =Past classified transactions, \hat{t} =Unevaluated transaction

Output: $result$ =Transaction \hat{t} classification

```
1: procedure CLASSIFICATION( $A, T, \hat{t}$ )
2:    $w_1 \leftarrow 0$ 
3:    $w_2 \leftarrow 0$ 
4:    $models = trainingModels(A, T)$ 
5:    $predictions = getPredictions(A, models)$ 
6:    $\mu \leftarrow getProbabilityAverage(predictions)$ 
7:   for each  $p$  in  $predictions$  do
8:     if  $getProbability(p) > \mu \wedge p == legitimate$  then
9:        $w_1 \leftarrow w_1 + 1$ 
10:    else
11:       $w_2 \leftarrow w_2 + 1$ 
12:    end if
13:  end for
14:  if  $w_1 > w_2$  then
15:     $result \leftarrow legitimate$ 
16:  else
17:     $result \leftarrow fraudulent$ 
18:  end if
19:  return  $result$ 
20: end procedure
```

5. Classification Algorithms

This section first explains the used criteria for the selection of the algorithm to use in our approach, then it describes the adopted ensemble criteria.

330 5.1. Selection Criteria

In order to implement the proposed evaluation model, we need that the classification algorithms, those in the set A of Section 3.2, not only predict the class label, but also provide the probability related to each class label. It should be observed that not all the classification algorithms provide this type of information, which represents a
335 kind of confidence level about the prediction. For this reason, during the composition of the set of algorithm A , we kept out the algorithms not providing this information and the algorithms performing a poor estimation of the class probabilities (i.e., those that, instead of a continuous probability value in $[0,1]$, returned only the 0 or 1 values).

The five algorithms that have been thus chosen for the experiments are: *Multilayer Perceptron (MLP)*, *Gaussian Naive Bayes (GNB)*, *Adaptive Boosting (ADA)*, *Gradient Boosting (GBC)*, and *Random Forests (RFC)*, and their settings are shown in Table 2.

The reason why we limited the number of algorithms to five derived by several analysis and work in literature, such as [43], which fixes to five the maximum number of algorithms to be used within an ensemble approach to obtain the best classification
345 performances.

5.2. Ensemble Criteria

We performed a set of experiments in order to try different combinations for the proposed *Prudential Multiple Consensus* model.

As a first step, we trivially used our model with single algorithms applying the
350 prudential voting defined in Section 4.1, as shown in Table 3. In order to underline the differences with respect to the native performance gained by each single algorithm, the table reports this information (*Native* columns) beside the performance gained by using the proposed approach (*Model* columns).

Afterwards, we tested our model by combining the algorithms in pairs, triples, quadruples, and finally by using all of the algorithms. The experimental results are
355 reported in Tables 4, 5, and 6.

They show that our model applied on a single algorithm reaches the best result by using *Random Forests*, which compared to the native performance indicates an improvement in terms of *fraudulent* transactions correctly detected (0.807% instead of 0.653%), slightly increasing the value of *Fallout* (0.016% instead of 0.000%) but improving that of *AUC* (0.896% instead of 0.827%).
360

Also by combining the algorithms in pairs, triples, and quadruples, we obtain the best results when *Random Forests* is involved. This is in line with other studies in literature which indicate this algorithm [44] as one of the best approaches in these kind
365 of tasks within the proposed domain [45, 4, 46].

Moreover, the results in Table 5 indicate the configuration based on four algorithms as the most promising, since by using the combination of the *MLP*, *GNB*, *GBC*, and *RFC* we get the best performances in terms of all the considered metrics.

The reader notices that the above results (even in the case of the single algorithms),
370 have been calculated using our *PMC* as decision strategy.

Table 2: Algorithms Configuration

Algorithm	Parameters
MLP	activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9, beta_2=0.999, early_stopping=False, epsilon=1e-08, hidden_layer_sizes=(100,), learning_rate='constant', learning_rate_init=0.001, max_iter=200, momentum=0.9, nesterovs_momentum=True, power_t=0.5, random_state=None, shuffle=True, solver='adam', tol=0.0001, validation_fraction=0.1, verbose=False, warm_start=False
GNB	priors=None
ADA	algorithm='SAMME.R', base_estimator=None, learning_rate=1.0, n_estimators=50, random_state=None
GBC	criterion='friedman_mse', init=None, learning_rate=0.1, loss='deviance', max_depth=3, max_features=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, presort='auto', random_state=None, subsample=1.0, verbose=0, warm_start=False
RFC	bootstrap=True, class_weight=None, criterion='gini', max_depth=None, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=1, oob_score=False, random_state=None, verbose=0, warm_start=False

Table 3: Single Algorithms Performance

<i>Evaluation algorithm</i>	<i>Native Sensitivity</i>	<i>Native Fallout</i>	<i>Native AUC</i>	<i>Model Sensitivity</i>	<i>Model Fallout</i>	<i>Model AUC</i>
Multilayer Perceptron (MLP)	0.146	0.000	0.573	0.781	0.629	0.576
Gaussian Naive Bayes (GNB)	0.709	0.012	0.848	0.739	0.016	0.862
Adaptive Boosting (ADA)	0.614	0.000	0.807	0.970	0.667	0.651
Gradient Boosting (GBC)	0.506	0.000	0.753	0.699	0.082	0.808
Random Forests (RFC)	0.653	0.000	0.827	0.807	0.016	0.896

Table 4: Ensemble Algorithms Performance by Pairs and Triples

<i>Algorithms by pairs</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>	<i>Algorithms by triples</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>
MLP, GNB	0.753	0.016	0.869	MLP, GNB, ADA	0.730	0.012	0.859
MLP, ADA	1.000	0.992	0.504	GNB, ADA, GBC	0.747	0.012	0.867
MLP, GBC	0.869	0.727	0.571	ADA, GBC, RFC	0.782	0.001	0.890
MLP, RFC	0.807	0.016	0.896	MLP, ADA, RFC	0.782	0.001	0.890
GNB, ADA	1.000	0.992	0.504	MLP, GNB, RFC	0.744	0.002	0.871
ADA, GBC	1.000	0.992	0.504	MLP, GBC, RFC	0.699	0.002	0.848
GBC, RFC	0.828	0.041	0.894	GNB, ADA, RFC	0.794	0.013	0.890
GNB, GBC	0.824	0.074	0.875	MLP, ADA, GBC	0.635	0.000	0.817
GNB, RFC	0.817	0.022	0.898	MLP, GNB, GBC	0.631	0.007	0.812

Table 5: Ensemble Algorithms Performance by Quadruples

<i>Algorithms by quadruples</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>	<i>Algorithms by quadruples</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>
GNB, ADA, GBC, RFC	0.807	0.016	0.895	MLP, GNB, ADA, RFC	0.807	0.016	0.895
MLP, ADA, GBC, RFC	0.797	0.005	0.896	MLP, GNB, ADA, GBC	0.768	0.013	0.878
MLP, GNB, GBC, RFC	0.800	0.005	0.897				

Table 6: Ensemble All Algorithms Performance

<i>Algorithms</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>
MLP, GNB, ADA, GBC, RFC	0.769	0.002	0.884

6. Experimental Environment

This section provides details on the experimental environment, on the dataset and the performed strategy, and on the metrics used to evaluate the classification performance.

375 6.1. Technological Environment

The development environment used to implement the approach presented in this paper is based on the *Python* language: the *scikit-learn*⁶ libraries have been used to implement the state-of-the-art algorithms. In order to ensure the reproducibility of the experiments we have carried out, the seed of the pseudo-random number generator
380 used by the *scikit-learn* classification algorithms has been set to *1*.

6.2. DataSet

The real-world dataset⁷ used for the experiment contains a series of transactions related to European cardholders and executed in two days of *2013*. As shown in Table 7, such a dataset presents a high degree of data imbalance [47], since only 492 out
385 of 284,807 transactions are classified as *fraudulent* (i.e., the 0.0017%). All the information in the dataset have been anonymized, except those related to the *time* and the *amount*, which contain, respectively, the seconds elapsed since the first transaction in the dataset and the amount of the underlying transaction.

Table 7: Dataset Details

Transactions	Legitimate	Fraudulent	Features	Classes
$ T $	$ T_+ $	$ T_- $	$ F $	$ C $
284,807	284,315	492	30	2

⁶<http://scikit-learn.org>

⁷<https://www.kaggle.com/mlg-ulb/creditcardfraud>

6.3. Strategy

390 In order to respect the transaction chronology, instead of a canonical *k-fold cross-validation* criterion we used the *TimeSeriesSplit* *scikit-learn* function to perform a *time series cross-validation* criterion. Such a function allows us splitting our dataset in a series of training and test sets, respecting the transactions chronology. For the experiments we used the *TimeSeriesSplit* function with *n_splits=10*.

395 The data imbalance problem, previously described in Section 2.3, has not been faced during the experiments. As suggested in [48], we have preferred to evaluate the effectiveness of our approach without any kind of data preprocessing (e.g., *under-sampling* or *over-sampling* balancing process) because in some cases the undersampling can potentially remove important samples whereas the oversampling can lead to
400 overfitting and increase the computational load when the dataset is already fairly large.

The existence of a statistical significance between the obtained results has been verified by using the independent-samples *two-tailed Student's t-tests* ($p < 0.05$).

6.4. Metrics

405 According to the considerations made in Section 2.5, the performance of the involved algorithms has been evaluated by using three metrics: the *Sensitivity*, the *Fallout*, and the *AUC* (i.e., Area Under the *ROC Curve*). As we mentioned before, such metrics have been chosen because they provide information about the performance in terms of *fraudulent* transactions correctly classified (*Sensitivity* and *Fallout*), a crucial indicator in the context taken into account, and in terms of effectiveness of the adopted
410 evaluation model (*AUC*).

In order to evaluate the algorithm performance also in terms of correct and incorrect classification of the *legitimate* transactions, we took into account two more metrics, which provide specular information with respect to the *Sensitivity* and *Fallout*: the *Specificity* and the *Miss Rate*.

415 The formulation of all the aforementioned metrics is presented below:

6.4.1. Sensitivity

The *Sensitivity* is calculated as reported in Equation 4, where \hat{T} is the set of new transactions to classify, *TP* is the number of transactions correctly classified as *fraudulent*, and *FN* is the number of *legitimate* transactions erroneously classified as *fraudulent*.
420

$$Sensitivity(\hat{T}) = \frac{TP}{(TP + FN)} \quad (4)$$

6.4.2. Fallout

The *Fallout* is calculated as reported in Equation 5, where \hat{T} is the set of new transactions to classify, *FP* is the number of *fraudulent* transactions erroneously classified as *legitimate*, and *TN* is the number of transactions correctly classified as *legitimate*.

$$Fallout(\hat{T}) = \frac{FP}{(FP + TN)} \quad (5)$$

425 **6.4.3. AUC**

The *AUC* is calculated as reported in Equation 6, where given the subsets of the past *legitimate* T_+ and the past *fraudulent* transactions L_- , Ψ indicates all the possible comparisons between these subsets (i.e., T_+ and T_-). Its result will be given by averaging all the comparisons and will lie within the interval $[0, 1]$, where I denotes the best performance.

$$\Psi(i_+, i_-) = \begin{cases} 1, & \text{if } i_+ > i_- \\ 0.5, & \text{if } i_+ = i_- \\ 0, & \text{if } i_+ < i_- \end{cases} \quad AUC = \frac{1}{|T_+||L_-|} \sum_1^{|T_+|} \sum_1^{|L_-|} \Psi(i_+, i_-) \quad (6)$$

435 **6.4.4. Specificity**

The *Specificity* is calculated as reported in Equation 7, where \hat{T} is the set of new transactions to classify, TN is the number of transactions correctly classified as *legitimate*, and FP is the number of *fraudulent* transactions erroneously classified as *legitimate*.

$$Specificity(\hat{T}) = \frac{TN}{(TN + FP)} \quad (7)$$

440 **6.4.5. Miss Rate**

The *Miss Rate* is calculated as reported in Equation 8, where \hat{T} is the set of new transactions to classify, FN is the number of *legitimate* transactions erroneously classified as *fraudulent*, and TP is the number of transactions correctly classified as *fraudulent*.

$$Miss\ Rate(\hat{T}) = \frac{FN}{(FN + TP)} \quad (8)$$

7. Results

This section reports the results of the performed experiments by comparing our solution to single and multiple algorithms approaches. Discussions on the results are also highlighted.

445 **7.1. Single Algorithm**

Figure 3 summarizes the comparison between our solution based on the *PMC* model and each single algorithm. Results indicate *Sensitivity*, *Fallout*, *Specificity*, *Miss Rate*, and *AUC* values. Please note that the *Sensitivity*, *Fallout*, and *AUC* values for the *PMC* model are the same shown in Table 5 using the best combination found (*MLP*, *GNB*, *GBC*, *RFC*). The reader notices that, we have not included the *Adaptive Boosting* (*ADA*) algorithm as indicated by a preliminary study discussed in Section 5.

450 As reported in Section 6.3, all the experiments have been performed according to a *time series cross-validation* criterion and after a thorough analysis of their results shown in Figure 3, we can made the following observations:

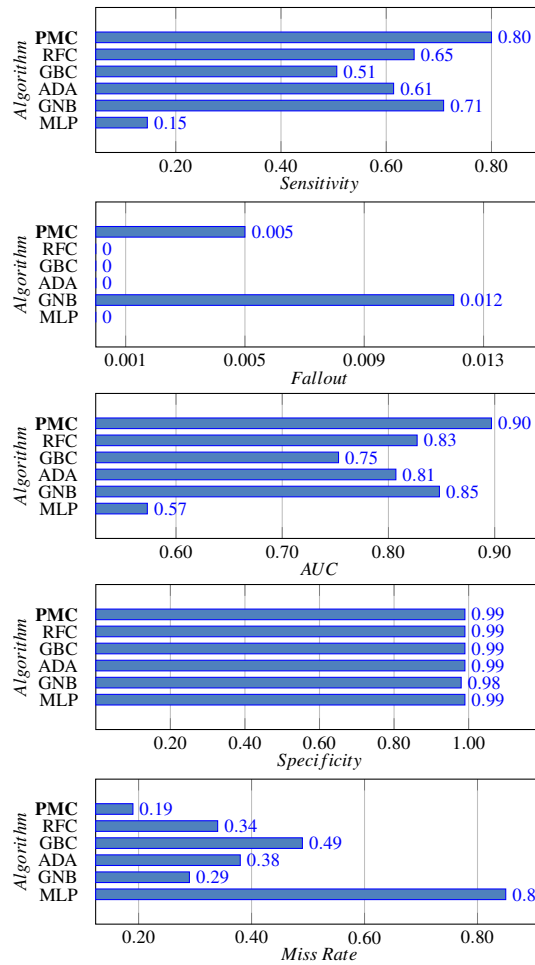


Figure 3: *Specificity, Miss Rate, AUC, Sensitivity, and Fallout*

- 455 • As *fraudulent* transactions in the dataset were 492, our approach correctly detected 394 of them (*Sensitivity* = 0.800), with only 2 misclassifications (*Fallout* = 0.005);
- compared to the best single approach (i.e., *GNB*), which correctly detected 349 *fraudulent* transactions, with 6 misclassifications, this means that our approach
460 had a gain of 9.1%;
- our solution was able to outperform the other algorithms in terms of *AUC* metric;
- the improvement of our approach is further confirmed in terms of *Specificity* and *Miss Rate*, as they prove that the obtained gain in terms of *Sensitivity* does not depend on a mere increase of *fraudulent* classifications made by our evaluation
465 model.

7.2. Multiple Algorithms

Here, we report the results of the last sets of experiments, which were aimed to compare the performance of the proposed approach, based on the *PMC* model, with that of the canonical state-of-the-art models used to manage the ensemble strategy of
470 classification, and to a state-of-the-art solution that operates by adopting a *Bagging* method based on the *Decision Tree* algorithm [35].

7.2.1. Model Comparison

This set of experiments is aimed at comparing our strategy based on the *PMC* model with the strategies such as the *complete agreement*, the *majority voting* (i.e., classification based on the majority of the classifications made by the algorithms), and the
475 *weighted voting* (i.e., classification based on the weight of the classifications made by the algorithms, in terms of class probability) between the algorithms. The experiments have been conducted by using the four chosen algorithms presented in Section 5 (i.e., the same algorithms used in our approach) and all the algorithms. The results reported
480 in Table 8 and Table 9 indicate that our *PMC* approach outperforms other ensemble approaches which use different decision strategies. Please note that the *Sensitivity*, *Fallout*, *AUC*, *Specificity*, and *Miss Rate* values for the *PMC* model have been found using the best combination found of Table 5 (MLP, GNB, GBC, RFC).

Table 8: Ensemble Strategies Comparison (Four Algorithms)

<i>Strategy</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>	<i>Specificity</i>	<i>Miss Rate</i>
Complete agreement	0.08	0.000	0.54	0.99	0.91
Majority voting	0.68	0.000	0.84	0.99	0.31
Weighted voting	0.55	0.000	0.77	0.99	0.44
PMC	0.80	0.005	0.89	0.99	0.19

7.2.2. Algorithm Comparison

485 The last set of experiments has been performed in order to evaluate our approach with a state-of-the-art approach employing a *Bagging* method based on the *Decision*

Table 9: Ensemble Strategies Comparison (All Algorithms)

<i>Strategy</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>	<i>Specificity</i>	<i>Miss Rate</i>
Complete agreement	0.06	0.000	0.53	1.00	0.93
Majority voting	0.63	0.000	0.81	0.99	0.36
Weighted voting	0.63	0.000	0.81	0.99	0.36
PMC	0.80	0.005	0.89	0.99	0.19

490 *Tree* algorithm [35] (denoted as *BDT*). According to the experimental criteria formalized in [35], our dataset has been divided into four parts by following the same percentage criterion (i.e., $P1 = 21.27\%$, $P2 = 27.19\%$, $P3 = 39.02\%$, $P4 = 12.52\%$), as shown in Table 10. The results of the experiments are reported in Table 11, which contains both the performances in terms of *Sensitivity*, *Fallout*, and *AUC* calculated on each part of the dataset (Table 10), and their average value calculated on all the dataset parts.

Table 10: Dataset Composition

Dataset part	Legitimate	Fraudulent	Total	Fraud rate
<i>P1</i>	60,415	163	60,578	0.0026%
<i>P2</i>	77,338	101	77,439	0.0013%
<i>P3</i>	110,942	189	111,131	0.0017%
<i>P4</i>	35,620	39	35,659	0.0010%
<i>Total</i>	284,315	492	284,807	

Table 11: Ensemble Algorithm Comparison (Bagging and Decision Tree)

<i>Approach</i>	<i>Dataset</i>	<i>Sensitivity</i>	<i>Fallout</i>	<i>AUC</i>
BDT	<i>P1</i>	0.75	0.010	0.81
BDT	<i>P2</i>	0.70	0.010	0.80
BDT	<i>P3</i>	0.65	0.019	0.79
BDT	<i>P4</i>	0.71	0.019	0.78
BDT	<i>Average</i>	0.70	0.014	0.79
PMC	<i>P1</i>	0.85	0.005	0.88
PMC	<i>P2</i>	0.80	0.005	0.89
PMC	<i>P3</i>	0.83	0.004	0.87
PMC	<i>P4</i>	0.77	0.006	0.80
PMC	<i>Average</i>	0.81	0.005	0.86

7.3. Discussion

495 The results highlighted in Section 7 indicate that the proposed approach based on our *PMC* model is able to improve the performance of a fraud detection system in terms of number of *fraudulent* transactions correctly classified.

Such an achievement is related to the value of *Sensitivity* (i.e., 0.800) and *Fallout* (i.e., 0.005) that, respectively, indicate its capability to correctly classify 9.1% of

500 *fraudulent* transactions more than the best competitor algorithm (GNB, which has a *Sensitivity* value of 0.71).

The results in terms of *AUC* metric underline the effectiveness of our evaluation model, proving its ability to classify new transactions as *legitimate* or *fraudulent*.

505 The evaluation in terms of *Specificity* and *Miss Rate* confirms the above results, showing that the increase of correctly classified *fraudulent* transactions implies a more robust and precise model.

The experiments aimed at comparing our method with other combined approaches show the effectiveness of our evaluation model (compared to the state-of-the-art based on the *complete agreement*, *majority voting*, and *weighted voting* criteria) in two different configurations, i.e. by using the four algorithms selected in Section 5 and by using all the algorithms.

510 The last set of experiments, where the performance of our approach has been compared to that of a performing state-of-the-art approach that uses a *Bagging* method based on the *Decision Tree* algorithm, also show that our *PMC* model outperforms its competitor, since it obtains best average performances in terms of *Sensitivity*, *Fallout*, and *AUC*.

515 Summarizing, we have proved that in real-world scenarios, characterized by a high degree of data imbalance, the proposed *PMC* model can significantly improve a fraud detection system, reducing the losses related to the misclassification of the fraudulent events.

520 The reader notices that the rationale of the *PMC* method, and the reason why it works well in the proposed domain, is because legitimate transactions are much higher in number and usually share a similar pattern easy to recognize. During the classification, several algorithms are thus able to assess with higher precision whether a transaction is legitimate. On the other hand, when a sample is fraudulent, most of the algorithms return a lower probability (confidence value) on their classification (either legitimate or fraudulent) and that is likely to be fraudulent. We have modeled this behaviour in our proposed *PMC* algorithm and this is why we obtain such high performances.

8. Conclusions and Future Work

530 In our era of big data, data intelligence and data security are very important research topics, and present constant challenges for academia and industry. The rapid evolution of the E-commerce platforms is an example of the increasing number of financial transactions made by electronic instruments of payment such as credit cards. Malicious people try to steal sensitive information from these transactions creating huge risks for the entire ecosystem. This is why, *Fraud Detection Systems*, especially those oriented to discover credit card frauds, are becoming more and more important.

540 The approach proposed in this paper, based on a novel *Prudential Multiple Consensus* model, addresses this problem and its risks associated with the aim of identify fraudulent transactions with higher precision than several state-of-the-art classification approaches. Our ensemble approach is able to reduce some well known problems that

affect this kind of classification tasks, first of all the issue related to the data imbalance, improving the classification performance in terms of number of frauds correctly detected.

545 All the performed experiments have been conducted by involving a real-world dataset characterized by a high degree of data imbalance, and the performances of our approach have been compared to those of several state-of-the-art solutions, both single and combined strategies, proving its effectiveness in terms of *Sensitivity* and *AUC*.

550 A future work would be to evaluate the proposed approach in other scenarios also characterized by a high degree of data imbalance, as well as the experimentation of new aggregation strategies based on the Artificial Neural Network.

Acknowledgments

This research is partially funded by: *Regione Sardegna* under project *Next generation Open Mobile Apps Development (NOMAD)*, *Pacchetti Integrati di Agevolazione (PIA) - Industria Artigianato e Servizi - Annualità 2013*; Italian Ministry of Education, University and Research - Program Smart Cities and Communities and Social Innovation project ILEARNTV (D.D. n.1937 del 05.06.2014, CUP F74G14000200008 F19G14000910008); Sardinia Regional Government (Convenzione triennale tra la Fondazione di Sardegna e gli Atenei Sardi Regione Sardegna L.R. 7/2007 annualità 2016 DGR 28/21 del 17.05.2016, CUP: F72F16003030002).

560 **References**

- [1] A. Abdallah, M. A. Maarof, A. Zainal, Fraud detection system: A survey, *J. Network and Computer Applications* 68 (2016) 90–113.
- [2] N. Japkowicz, S. Stephen, The class imbalance problem: A systematic study, *Intell. Data Anal.* 6 (5) (2002) 429–449.
- 565 [3] H. He, E. A. Garcia, Learning from imbalanced data, *IEEE Trans. Knowl. Data Eng.* 21 (9) (2009) 1263–1284. doi:10.1109/TKDE.2008.239.
- [4] I. Brown, C. Mues, An experimental comparison of classification algorithms for imbalanced credit scoring data sets, *Expert Syst. Appl.* 39 (3) (2012) 3446–3453. doi:10.1016/j.eswa.2011.09.033.
- 570 [5] R. C. Holte, L. Acker, B. W. Porter, Concept learning and the problem of small disjuncts, in: N. S. Sridharan (Ed.), *Proceedings of the 11th International Joint Conference on Artificial Intelligence*. Detroit, MI, USA, August 1989, Morgan Kaufmann, 1989, pp. 813–818.
- 575 [6] M. Lek, B. Anandarajah, N. Cerpa, R. Jamieson, Data mining prototype for detecting e-commerce fraud, in: S. Smithson, J. Gricar, M. Podlogar, S. Avgerinou (Eds.), *Proceedings of the 9th European Conference on Information Systems, Global Co-operation in the New Millennium, ECIS 2001, Bled, Slovenia, June 27-29, 2001, 2001*, pp. 160–165.

- 580 [7] N. Carneiro, G. Figueira, M. Costa, A data mining based system for credit-card fraud detection in e-tail, *Decision Support Systems* 95 (2017) 91–101.
- [8] A. J. Hoffman, R. E. Tessendorf, Artificial intelligence based fraud agent to identify supply chain irregularities, in: M. H. Hamza (Ed.), *IASTED International Conference on Artificial Intelligence and Applications*, part of the 23rd Multi-Conference on Applied Informatics, Innsbruck, Austria, February 14-16, 2005, IASTED/ACTA Press, 2005, pp. 743–750.
- 585 [9] D. G. Whiting, J. V. Hansen, J. B. McDonald, C. C. Albrecht, W. S. Albrecht, Machine learning methods for detecting patterns of management fraud, *Computational Intelligence* 28 (4) (2012) 505–527.
- [10] I. Nolan, Transaction fraud detection using random forest classifier and logistic regression, *Neural Networks & Machine Learning* 1 (1) (2017) 2–2.
- 590 [11] C. Assis, A. M. Pereira, M. de Arruda Pereira, E. G. Carrano, Using genetic programming to detect fraud in electronic transactions, in: C. V. S. Prazeres, P. N. M. Sampaio, A. Santanchè, C. A. S. Santos, R. Goularte (Eds.), *A Comprehensive Survey of Data Mining-based Fraud Detection Research*, Vol. abs/1009.6119, 2010, pp. 337–340.
- 595 [12] A. Mead, T. Lewis, S. Prasanth, S. Adams, P. Alonzi, P. Beling, Detecting fraud in adversarial environments: A reinforcement learning approach, in: *Systems and Information Engineering Design Symposium (SIEDS)*, 2018, IEEE, 2018, pp. 118–122.
- 600 [13] W. Wang, D. Lu, X. Zhou, B. Zhang, J. Mu, Statistical wavelet-based anomaly detection in big data with compressive sensing, *EURASIP J. Wireless Comm. and Networking* 2013 (2013) 269.
- [14] R. Saia, A discrete wavelet transform approach to fraud detection, in: *NSS*, Vol. 10394 of *Lecture Notes in Computer Science*, Springer, 2017, pp. 464–474.
- 605 [15] R. Saia, S. Carta, Evaluating credit card transactions in the frequency domain for a proactive fraud detection approach, in: *SECRYPT*, SciTePress, 2017, pp. 335–342.
- [16] R. Saia, S. Carta, A frequency-domain-based pattern mining for credit card fraud detection, in: *IoTBDS*, SciTePress, 2017, pp. 386–391.
- 610 [17] S. Y. Huang, C. Lin, A. Chiu, D. C. Yen, Fraud detection using fraud triangle risk factors, *Information Systems Frontiers* 19 (6) (2017) 1343–1356.
- [18] R. Saia, Unbalanced data classification in fraud detection by introducing a multi-dimensional space analysis, in: *IoTBDS*, SciTePress, 2018, pp. 29–40.
- [19] R. J. Bolton, D. J. Hand, Statistical fraud detection: A review, *Statistical Science* (2002) 235–249.
- 615

- [20] C. Phua, V. C. S. Lee, K. Smith-Miles, R. W. Gayler, A comprehensive survey of data mining-based fraud detection research, CoRR abs/1009.6119.
- [21] A. D. Pozzolo, O. Caelen, Y. L. Borgne, S. Waterschoot, G. Bontempi, Learned lessons in credit card fraud detection from a practitioner perspective, Expert Syst. Appl. 41 (10) (2014) 4915–4928. doi:10.1016/j.eswa.2014.02.026.
- [22] H. Wang, W. Fan, P. S. Yu, J. Han, Mining concept-drifting data streams using ensemble classifiers, in: L. Getoor, T. E. Senator, P. M. Domingos, C. Faloutsos (Eds.), Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 24 - 27, 2003, ACM, 2003, pp. 226–235. doi:10.1145/956750.956778.
- [23] J. Gao, W. Fan, J. Han, P. S. Yu, A general framework for mining concept-drifting data streams with skewed distributions, in: Proceedings of the Seventh SIAM International Conference on Data Mining, April 26-28, 2007, Minneapolis, Minnesota, USA, SIAM, 2007, pp. 3–14. doi:10.1137/1.9781611972771.1.
- [24] E. L. Barse, H. Kvarnström, E. Jonsson, Synthesizing test data for fraud detection systems, in: ACSAC, IEEE Computer Society, 2003, pp. 384–394.
- [25] A. Chatterjee, A. Segev, Data manipulation in heterogeneous databases, ACM SIGMOD Record 20 (4) (1991) 64–68.
- [26] S. Sorournejad, Z. Zojaji, R. E. Atani, A. H. Monadjemi, A survey of credit card fraud detection techniques: Data and technique oriented perspective, CoRR abs/1611.06439.
- [27] J. Attenberg, F. J. Provost, Inactive learning?: difficulties employing active learning in practice, SIGKDD Explorations 12 (2) (2010) 36–41. doi:10.1145/1964897.1964906. URL <http://doi.acm.org/10.1145/1964897.1964906>
- [28] M. Zareapoor, J. Yang, A novel strategy for mining highly imbalanced data in credit card transactions, Intelligent Automation & Soft Computing (2017) 1–7.
- [29] V. Vinciotti, D. J. Hand, Scorecard construction with unbalanced class sizes, Journal of Iranian Statistical Society 2 (2) (2003) 189–205.
- [30] H. M. Gomes, J. P. Barddal, F. Enembreck, A. Bifet, A survey on ensemble learning for data stream classification, ACM Comput. Surv. 50 (2) (2017) 23:1–23:36.
- [31] T. G. Dietterich, Ensemble methods in machine learning, in: Multiple Classifier Systems, Vol. 1857 of Lecture Notes in Computer Science, Springer, 2000, pp. 1–15.
- [32] S. Akila, U. S. Reddy, Risk based bagged ensemble (rbe) for credit card fraud detection, in: Inventive Computing and Informatics (ICICI), International Conference on, IEEE, 2017, pp. 670–674.

- [33] H. Guo, Y. Li, J. Shang, G. Mingyun, H. Yuanyue, G. Bing, Learning from class-imbalanced data: Review of methods and applications, *Expert Syst. Appl.* 73 (2017) 220–239.
- 655 [34] M. Faghani, M. J. Nordin, S. Shojaeipour, Optimization of the performance face recognition using adaboost-based, in: R. Chen (Ed.), *Intelligent Computing and Information Science*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2011, pp. 359–365.
- [35] M. Zareapoor, P. Shamsolmoali, Application of credit card fraud detection: Based on bagging ensemble classifier, *Procedia Computer Science* 48 (2015) 679–685.
- 660 [36] D. Wu, Y. Liu, G. Gao, Z. Mao, W. Ma, T. He, An adaptive ensemble classifier for concept drifting stream, in: *Computational Intelligence and Data Mining, 2009. CIDM'09. IEEE Symposium on, IEEE, 2009*, pp. 69–75.
- [37] Y. Bian, M. Cheng, C. Yang, Y. Yuan, Q. Li, J. L. Zhao, L. Liang, Financial fraud detection: a new ensemble learning approach for imbalanced data., in: *PACIS, 2016*, p. 315.
- 665 [38] Y. Freund, R. E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, *J. Comput. Syst. Sci.* 55 (1) (1997) 119–139.
- [39] A. Natekin, A. Knoll, Gradient boosting machines, a tutorial, *Front. Neurorobot.* 2013.
- 670 [40] A. O. Adewumi, A. A. Akinyelu, A survey of machine-learning and nature-inspired based credit card fraud detection techniques, *International Journal of System Assurance Engineering and Management* 8 (2) (2017) 937–953.
- [41] D. M. Powers, Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation.
- 675 [42] D. Faraggi, B. Reiser, Estimation of the area under the roc curve, *Statistics in medicine* 21 (20) (2002) 3093–3106.
- [43] A. T. Sergio, T. P. F. de Lima, T. B. Ludermir, Dynamic selection of forecast combiners, *Neurocomputing* 218 (2016) 37–50.
- 680 [44] L. Breiman, Random forests, *Machine Learning* 45 (1) (2001) 5–32. doi:10.1023/A:1010933404324.
- [45] S. Lessmann, B. Baesens, H. Seow, L. C. Thomas, Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research, *European Journal of Operational Research* 247 (1) (2015) 124–136. doi:10.1016/j.ejor.2015.05.030.
- 685 [46] S. Bhattacharyya, S. Jha, K. K. Tharakunnel, J. C. Westland, Data mining for credit card fraud: A comparative study, *Decision Support Systems* 50 (3) (2011) 602–613. doi:10.1016/j.dss.2010.08.008. URL <http://dx.doi.org/10.1016/j.dss.2010.08.008>

- ⁶⁹⁰ [47] A. D. Pozzolo, O. Caelen, R. A. Johnson, G. Bontempi, Calibrating probability with undersampling for unbalanced classification, in: IEEE Symposium Series on Computational Intelligence, SSCI 2015, Cape Town, South Africa, December 7-10, 2015, IEEE, 2015, pp. 159–166. doi:10.1109/SSCI.2015.33.
- ⁶⁹⁵ URL <http://dx.doi.org/10.1109/SSCI.2015.33>
- [48] N. V. Chawla, N. Japkowicz, A. Kotcz, Special issue on learning from imbalanced data sets, ACM Sigkdd Explorations Newsletter 6 (1) (2004) 1–6.