

Proactivity or Retroactivity? Evaluating the Benefits in Using Proactive Transformed-domain-based Techniques in Fraud Detection Tasks

Roberto Saia and Salvatore Carta¹

*Dipartimento di Matematica e Informatica, Università di Cagliari
Via Ospedale 72, 09124 Cagliari, Italy*

Abstract

The exponential number of E-commerce transactions indicates a radical change in the way people buy and sell goods and services, a new opportunity offered by a huge global market, where they can choose their sellers or buyers on the basis of multiple criteria (e.g., economic, logistical, ethical, sustainability, etc.), instead of being forced to use the traditional brick-and-mortar criterion. If on the one hand such scenario offers enormous control to people, both at private and corporate level, allowing them to filter their needs by adopting a large range of criteria, on the other hand it has contributed to the growth of fraud cases related to the involved electronic instruments of payment, such as credit cards. The Big Data Information Security for Sustainability is a research branch aimed to face these issues in relation to the potential implications in the field of sustainability, proposing effective solutions to design safe environments in which the people can operate, exploiting the benefits offered by the new technologies. The fraud detection systems are a significant example of such solutions, although the techniques adopted by them are typically based on retroactive strategies, which are incapable of preventing fraudulent events. In this perspective, this paper aims to investigate the benefits related to the adoption of proactive fraud detection strategies, instead of the canonical retroactive ones. We evaluate two previously experimented novel proactive approaches, one based on the Fourier transform, and one based on the Wavelet transform, which are used in order to move the data (i.e., financial transactions) into a new domain, where they are analyzed and an evaluation model is defined. Such strategies allow a fraud detection system to operate proactively, since they do not need previous fraudulent examples, overcoming some important problems that reduce the effectiveness of the canonical retroactive state-of-the-art solutions. Potential benefits and limitations of the proposed approaches have been evaluated in a real-world credit card fraud detection scenario, by comparing their performance to that of one of the most used and performing retroactive state-of-the-art approaches (i.e. Random Forests).

Keywords: Business Intelligence, Sustainability, Anomaly Detection

Email address: {roberto.saia, salvatore}@unica.it (Roberto Saia and Salvatore Carta)

1. Introduction

Nowadays, the *Big Data Analytics for Sustainability (BDAS)* represents a crucial research field, since it offers us the opportunity to exploit the new technologies in a smarter way [1, 2], developing more sustainable products and processes.

In the context of the advantages in terms of sustainability offered by the E-commerce environment [3], where the great offer of goods and services allows people to choose those that respect this criterion, helped by the vast amount of information on the Internet, the fraud detection systems represent an instrument around which the interests of many economic entities revolve.

So as it happens in other fields related to technology, even in those where the sustainability represents an essential element, the potential advantages are jeopardized by who try to take advantage of the new technologies fraudulently. For the aforementioned reasons, one of the most important *BDAS* objectives is the *Big Data Information Security for Sustainability (BDISS)*. Some *BDISS* areas of great interest are, for instance, those directly related to the security of the adopted platforms (e.g., *Intrusion Detection*, *Fraud Detection*, etc.) and those indirectly related to them (e.g., *Privacy Preserving*, *Cyber Espionage*, etc.).

This paper is focused to one of these important areas, since it faces the problems related to the fraudulent use of the electronic payment instruments, nowadays an essential element for the exchange of goods and services.

Authoritative reports¹ underlined an exponential growth in the fraud losses related to the credit and debit cards, as shown in Figure 1. Several studies² have also indicated how the purchases made without authorization and the counterfeits of credit cards represent the 10-15% of total fraud cases, but the 75-80% of financial value. Only in the United States, such problem leads toward an estimated average loss per fraud case of 2 million of dollars.

The aforementioned scenario has generated an increasing in research and development investments by private and public entities, with the objective to design more and more effective methods able to face this problem.

It should be observed how the design of effective solutions represents a hard challenge due to several well-know issues, which reduce the capability of the state-of-the-art techniques used in this specific field. The most important issue consists in the fact that the *fraudulent* transactions are typically less than the *legitimate* ones, and such highly unbalanced data distribution reduces the effectiveness of the machine learning strategies [4]. In addition to this issue there is the scarcity of information that characterizes the involved financial transactions, a problem that leads toward an overlapping of the classes of expense of a user [5].

Nowadays, a fraud detection system can exploit many state-of-the-art techniques in order to evaluate a financial transaction. For instance, it can exploit: *Data Mining*

¹Nilson Report: <https://www.nilsonreport.com/>

²American Association of Fraud Examiners: <http://www.acfe.com>

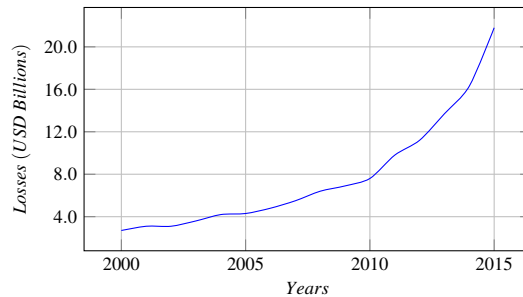


Figure 1: Annual Global Fraud Losses

techniques to generate rules from fraud patterns [6]; *Artificial Intelligence* techniques to identify data irregularities [7]; *Fuzzy Logic* techniques to perform a fuzzy analysis for a fraud detection task [8]; *Machine Learning* techniques to define ensemble methods that combine predictions from multiple models [9]; *Genetic Programming* techniques to model and detect fraud through an *Evolutionary Computation* approach [10]; *Statistical Inference* techniques that adopt a flexible Bayesian model for fraud detection [11].

However, it should be observed that regardless of the adopted technique, the principle that is commonly exploited is the detection of outliers in the transactions under analysis, a trivial approach that usually lead toward misclassifications, with all the financial consequences that derive from it (mainly money loss). The reason behind these wrong classifications is the absence of extensive evaluation criteria, since many state-of-the-art techniques are not able to manage some transaction features during the evaluation process (e.g., the non-numeric ones). For instance, *Random Forests* [12], one of the most performing approaches, is not able to manage types of data that involve a large number of categories. For the aforementioned reasons, the evaluation process performed by a fraud detection system should be able to take into account all the information about the transactions under analysis.

This paper is aimed to evaluate the benefits related to the adoption of proactive approaches, where the analysis of the transaction data is performed in a transformed-domain, instead of a canonical one, by adopting two previously experimented approaches, one based on the *Fourier transform* and one based on the *Wavelet transform*³.

In the first approach, the evaluation model is defined in terms of the spectral pattern of a transaction, by processing the information through the *Fourier transformation* [13]. In the second approach, the evaluation model is defined by following a similar criterion but by processing the information through the *Wavelet transformation* [14].

In both approaches we consider the sequence of values of each transaction feature as a *time series*, moving its representation in the *frequency-domain* by using the *Discrete Fourier Transform (DFT)* or in the *time-frequency-domain* by using the *Discrete Wavelet Transform (DWT)* process.

³These two approaches have been presented in two papers that will be published soon, but all the information needed for the purposes of this paper are here provided.

Many evaluation models adopted in the *Business Intelligence* field are defined on the basis of *time series* [15]. This happens when it is important to characterize the involved elements on the basis of the time factor [16]. The information extracted from the time series can be exploited in order to perform different tasks, such as those related to the *risk analysis* (e.g., *Credit Scoring* [17] and *Stock Forecasting* [18]) and *Information Security* (e.g., *Fraud Detection* [19] and *Intrusion Detection* [20]) ones.

In the context of the proposed approach, as *time series* we mean something slightly different to the canonical meaning given to it in literature. This because we refer to it in terms of data used as input in the *DFT* or *DWT* process, thus just in terms of sequence of values that compose a transaction (i.e., *date*, *amount*, and so on), considering them a sequence of discrete-time data taken at successive equally spaced points in time.

In other words, the relationship between *time series* and our fraud detection approach must be sought in the analysis, performed in the frequency domain, of patterns given by the feature values of a transaction. This because our goal is to define an evaluation model able to characterize the legitimate transactions, regardless of the time in which they occur, also because the *time frame* taken into account in our approaches is limited to the *feature space*.

It should be added that the involved transactions are only those related to the previous *legitimate* cases, because we do not consider the previous *fraudulent* ones.

The analysis of the information in the new domain presents some advantages. The first one is the capability for a fraud detection system to include only the previous *legitimate* transactions in the model definition process, which is a proactive approach able to face the *cold-start* issue (i.e., scarcity or absence of *fraudulent* cases during the model definition). Another advantage is related to the consideration that the new data representation reduces the data heterogeneity problem, since the new domain is less influenced by the data variations.

For exemplification reasons, from now on, we will use the term *transformed-domain* to refer to both the data domain obtained by the *Fourier transform (frequency-domain)* and the data domain obtained by the *Wavelet transform (time-frequency-domain)*.

This paper is based on a previous work presented in [19], which has been completely rewritten and extended, producing the following scientific contributions:

- (i) formalization of the procedure used to define the *time series* to use as input in the data transformation process performed by the *DFT* or *DWT* approach, introducing also an alternative method suitable for certain fraud detection contexts;
- (ii) formalization of the comparison process between the output obtained at the end of the *DFT* or *DWT* process, presenting two different modalities, one based on the *cosine similarity* measured between the entire output vectors of the *DFT* or *DWT* processes, and one based on the *punctual comparison* between each element of these output vectors;
- (iii) formalization of a generic algorithm able to classify a new transaction as *legitimate* or *fraudulent*, by exploiting the previous comparison process and the *DFT* or *DWT* process, also defining its asymptotic time complexity;

- (iv) evaluation of our proactive *DFT* and *DWT* approaches in terms of general performance and predictive power of their classification model, performed by using two real-world datasets;
- (v) evaluation of the advantage and disadvantages related to the adoption of a proactive strategies in the context of the fraud detection processes, on the basis of the experimental results.

The paper is organized as follows: Section 2 introduces the background and related work of the scenario taken into account; Section 3 provides a formal notation and defines the problem taken into account; Section 4 describes our approach; Section 5 provides details on the experimental environment, on the used datasets and metrics, as well as on the adopted strategy and selected competitor, presenting the experimental results; some concluding remarks and future work are given in the last Section 6.

2. Background and Related Work

The main goal of a fraud detection system is the evaluation of the new transactions in order to classify them as *legitimate* or *fraudulent*), on the basis of an evaluation model previously defined by exploiting the information gathered by the system during the previous transactions.

Premising that the most effective state-of-the-art techniques of fraud detection operate by adopting a retroactive approach, thus they need to train their evaluation models with both the classes of transactions (i.e., *legitimate* and *fraudulent* previous cases), in this section we want to offer an overview of the today' s scenario.

This section starts by providing an overview of the *Big Data Information Security* concept, then it focuses on the most used fraud detection techniques and their important open problems. It concludes by introducing the theoretical concepts behind the proposed proactive approaches, along with the description of the competitor approach chosen to evaluate their performance.

2.1. Big Data Information Security

The concept of *Big Data Information Security (BDIS)* and its application in the context of the *sustainable technologies* are introduced in the following.

2.1.1. Overview

The main challenge of a *BDIS* process is the analysis of huge and heterogeneous data with the goal to protect them against a series of risks such as, for instance, their alteration (*integrity*) or their unauthorized use (*confidentiality*) [21, 22].

It should be observed how in this age of information the risks related to the gathering and use of data are in most of the cases tolerated in view of the great advantages that such operations offer in many fields (e.g., medical, financial, environmental, social, and so on). This kind of paradox has been discussed in literature through several studies, such as that performed in [23].

The main disadvantage of almost all the techniques used to define approaches able to face this type of risk (e.g., alteration or fraudulent use of data) need a considerable

number of examples of all possible cases to build their evaluation models (e.g., in the context of a credit card fraud detection system, they need both legitimate and fraudulent examples), precluding the adoption of proactive strategies.

2.1.2. Sustainability

As previously introduced in Section 1, the E-commerce platform allows people to have access to a huge number of goods and services, enabling them to make their own choices on the basis of different criteria. Nowadays, this has been made possible by the coexistence of two factors: a huge E-commerce platform and an equally huge source of information (i.e., Internet).

Through the Internet people are able to choose sellers and buyers not only on the basis of convenience metrics, but also by following innovative metrics such as, for example, the ethical ones.

In this context, dominated by electronic payment instruments, fraud detection systems [24, 19] play a crucial role, since they are aimed to detect the fraudulent financial transactions, allowing people to get only the benefits offered by the E-commerce infrastructure.

2.2. Fraud Detection Systems

Here are reported the most common strategies and approaches of fraud detection.

2.2.1. Strategies

The fraud detection approaches can operate by adopting *supervised* or *unsupervised* strategies [25]. By using a *supervised* strategy they exploit the previous *fraudulent* and *non-fraudulent* transactions collected by the system, and they use them to define an evaluation model able to classify a new transaction as *legitimate* or *fraudulent*. In order to perform this task they need to have a sufficient number of examples of both classes, and their recognition capability depends on the known patterns.

By using an unsupervised strategy they instead work by finding anomalies in a transaction under evaluation, in terms of substantial differences in the feature values (wrt the typical values assumed in the past). Considering that a *fraudulent* transaction can be characterized by features with values within their typical range, adopting *unsupervised* strategies [26] in a fraud detection system represents a hard challenge.

2.2.2. Approaches

The *static approach* [27] represents the most common way to operate in order to detect *fraudulent* events in a financial data stream related to a credit card activity. By following it, the data streaming is divided into equal size blocks and the evaluation model is trained by using a limited number of initial and contiguous blocks.

The *updating approach* [28] adopts instead a different modality, since at each new block the evaluation model is updated by training it with a certain number of latest and contiguous blocks.

The *forgetting approach* [29] represents another modality, where the evaluation model is updated when a new block appears, and this operation is performed by using all the previous *fraudulent* transactions, but only the *legitimate* transactions present in the last two blocks.

The models obtained through these approaches can be directly used for the evaluation process or they can be combined in order to define a biggest model of evaluation.

It should be noted that all the aforementioned approaches present some limitations: the *static approach* is not able to model the users behavior; the *updating approach* is not able to operate with small classes of data; the *forgetting approach* presents a high computational complexity. In addition, there are some common issues to overcome that reduce the effectiveness of all these approaches, as described in the following Section 2.3.

2.3. Open Problems

The most common problems related to the fraud detection tasks are reported and described in the following.

2.3.1. Data Scarcity

The scarcity of public real-world datasets [10] is the first problem that researchers working in this area have to deal with. It is related to the restrictive policies that regulate the disclosure of information in this area, which they do not allow the operators to provide information about their business activities. Such restrictive rules are related to privacy, competition, or legal reasons.

It should be added that not even a release in anonymous form of the data is usually considered acceptable by many financial operators, because even in this form the data may reveal crucial information, such as some vulnerabilities in the E-commerce infrastructure.

2.3.2. Non-adaptability

In the context of a fraud detection system, this is a problem related to the difficulty for the evaluation model to correctly classify new transactions, when they are characterized by patterns differing from those used to train the model.

This kind of problem affect both the *supervised* and *unsupervised* fraud detection approaches [30], leading toward misclassifications.

2.3.3. Data Heterogeneity

In the machine learning field, the pattern recognition is a process aimed to assign a label to a given input value. Some common applications of such process are the classification tasks, where this process is performed in order to classify each values in input into a specific class (within a finite set of classes). It can be exploited in a large number of contexts, thanks to its capability to solve a large number of real-world problems [31], although its effectiveness is usually affected by the heterogeneity of the involved data.

This is a problem described in literature as *naming problem* or *instance identification problem* and it is related to the incompatibility between similar features resulting in the same data being represented differently in different datasets [32, 33].

Given the high level of heterogeneity that characterizes the fraud scenarios (e.g., that related to the credit card transactions), an effective fraud detection system must be able to address the *data heterogeneity* issue.

2.3.4. Data Unbalance

A fraud detection task can be considered as an unbalanced data classification problem [34], because the examples used to train the evaluation model are typically composed by a large number of *legitimate* cases and a small number of *fraudulent* ones, a data configuration that reduces the effectiveness of the classification approaches [4, 35, 36].

This problem is probably worsened by a *data alteration* in the datasets publicly released by some financial operators, where in order to maintain customer trust in their services, the fraud cases have been intentional reduced, classifying part of them as legitimate.

Considering that the canonical approaches of fraud detection operate retroactively, thus they need to train their model by using both classes of examples (i.e., *legitimate* and *fraudulent*), such problem is commonly faced by preprocessing the dataset in order to obtain an artificial balance of data [37].

This kind of operation can be performed through an *over-sampling* or *under-sampling* method, where in the first case the balance is made by duplicating some of the transactions that are less in number (typically, the *fraudulent* ones), while in the second case it is made by removing some of the transactions that are in greater number (typically, the *legitimate* ones). Some studies demonstrate that the adoption of *re-sampling* methods improves the performance given by the original imbalanced data, also underlining how the *over-sampling* techniques perform better than the *under-sampling* ones [38, 39, 40].

2.3.5. Cold-start

In order to be able to operate properly, machine learning approaches need a significant amount of data to define their evaluation models. While in some contexts this is not a significant issue, in other ones such as, for example, those related to the fraud detection, it represents a big issue. It happens because the examples are characterized by a large number of *legitimate* cases and a small number of *fraudulent* ones, as described in Section 2.3.4.

This configures the so-called *cold-start problem*, i.e., the set of data used to train an evaluation model does not contain enough information about the domain taken into account, making the definition of a reliable evaluation model difficult [41]. In the context taken into account in this paper, this problem arises when the training data is not representative of all the involved classes (*legitimate* and *fraudulent*) of information [42].

2.4. Proposed Approach and Competitor

The objective of our approach can be reached by adopting two different modalities, one based on the *Fourier Transform* and one based on the *Wavelet Transform*. This section describe both the aforementioned modalities, providing also details about the state-of-the-art approach (i.e., *Random Forests*) used to evaluate the performance achieved by using them in a real-world fraud detection context.

2.4.1. Time Series Definition

A *time series* [43] usually refers to a series of values acquired by measuring the variation in the time of a specific data type (i.e., temperature, amplitude, and so on).

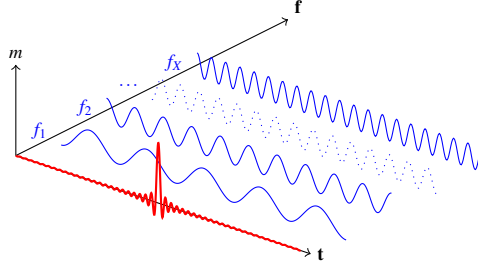


Figure 2: *Time and Frequency Domains*

In our approaches we consider as *time series* the sequence of values assumed by the transaction features in the datasets taken into account, introducing also a different modality where the *time series* are defined in terms of sequence of values assumed by each single transaction feature in the dataset domain. This last modality is suitable when we need to model the behavior of a single transaction features, instead of that of all features in the context of a transaction (i.e., how it happens in the considered credit card fraud detection task).

2.4.2. Approach 1: Fourier Transform

The idea behind this first approach is to perform the evaluation process in a *frequency-domain*, by defining the evaluation model in terms of frequency components. Such operation is performed by considering the sequence of values assumed by the transaction features as a *time series*, moving its analysis from the canonical domain to a new *transformed-domain*.

The result is a spectral pattern composed by the frequency components, as shown in Figure 2, where t , f , and m , respectively stand for *time*, *frequency*, and *magnitude*.

We made this by recurring to the *Discrete Fourier Transform (DFT)*, whose formalization is shown in Equation 1, where i denotes the imaginary unit.

$$F_n \stackrel{\text{def}}{=} \sum_{k=0}^{N-1} f_k \cdot e^{-2\pi i k n / N}, \quad n \in \mathbb{Z} \quad (1)$$

The result is a set of sinusoidal functions, each of them related to a specific frequency component. We can return to the original time domain by using the *inverse Fourier transform* reported in Equation 2.

$$f_k = \frac{1}{N} \sum_{n=0}^{N-1} F_n \cdot e^{2\pi i k n / N}, \quad n \in \mathbb{Z} \quad (2)$$

A periodic wave is characterized by a frequency f and a wavelength λ (i.e., the distance in the medium between the beginning and end of a cycle $\lambda = \frac{w}{f_0}$, where w stands for the wave velocity), which are defined by the repeating pattern, the non-periodic waves that we take into account during the *Discrete Fourier Transform* process do not have a frequency and a wavelength. Their fundamental period τ is the period

where the wave values were taken and sr denotes their number over this time (i.e., the acquisition frequency).

Assuming that the time interval between the acquisitions is constant, on the basis of the previous definitions applied in the context of this paper, the considered non-periodic wave is given by the sequence of values assumed by each distinct feature $v \in V$ that characterize the transactions in the set T_+ (i.e., the past *legitimate* transactions), and this sequence of values represents the *time series* taken into account in the *DFT* process.

Their fundamental period τ starts with the value assumed by the feature in the oldest transaction of the set T_+ and it ends with the value assumed by the feature in the newest transaction, thus we have that $sr = |T_+|$; the sample interval si is instead given by the fundamental period τ divided by the number of acquisition, i.e., $si = \frac{\tau}{|T_+|}$.

The *transformed-domain* representation, obtained by the *DFT*, process gives us information about the magnitude and phase of the signal at each frequency. Denoting as x the output of the process, it represents a series of complex numbers, where x_r is the real part and x_i is the imaginary one (i.e., we have that $x = (x_r + ix_i)$).

Premising that the magnitude can be calculated by using $|x| = \sqrt{(x_r^2 + x_i^2)}$ and that the phase can be calculated by using $\phi(x) = \arctan\left(\frac{x_i}{x_r}\right)$, in the context of the presented approach we will only take into account the frequency magnitude.

In the context of the presented approach we use the *Fast Fourier Transform (FFT)* algorithm in order to perform the Fourier transformations, since it allows us to rapidly compute the *DFT* by factorizing the input matrix into a product of sparse (mostly zero) factors. This is a largely used algorithm because it is able to reduce the computational complexity of the process from $O(n^2)$ to $O(n \log n)$ (where n denotes the data size).

2.4.3. Approach 2: Wavelet Transform

The idea behind this second approach is to move the evaluation process from the canonical domain to a new *transformed-domain* by exploiting the *Discrete Wavelet Transformation (DWT)* [44, 45]. In more detail, we use the *DWT* process in a *time series* data mining context.

The evaluation of the transactions in the new domain offered by the *DWT* leads toward interesting advantages. Such process transforms a *time series* by exploiting a set of functions named *wavelets* [45], and in literature it is usually performed in order to reduce the data size (e.g., in the image compression tasks) or to reduce the data noise (e.g., in the filtering tasks). The wavelets are mathematical functions that allow us to decompose the original data into different frequencies at different scales, then they move the data representation from the time domain (sequence of transaction feature values) to a new domain where the data is represented both in terms of frequency and time.

The so-called *time-scale multiresolution* offered by the *DWT* represents an important aspect of this process, since it allows us to observe the original *time series* from different points of view, each containing interesting information on the original data. As frequency we mean the number of occurrences of a value in a *time series* over a unit of time and as scale we mean the time interval that characterize the *time series*. The capability in the new domain to observe the data by using multiple scales (multiple res-

olution levels), allows our approach to define a more stable and representative model of the transactions, wrt the canonical approaches at the state of the art.

The process of transformation operated by the *DWT* is different from that carried out by similar approaches, such as the Fourier transforms [13], characterized by a constant resolution for all the frequencies, because it analyzes the data at multiple resolution for different frequencies. Formally, a *Continuous Wavelet Transform (CWT)* is defined as shown in Equation 3, where $\Psi(t)$ represents a continuous function in both the time and frequency domain (called *mother wavelet*) and the $*$ denoting the complex conjugate.

$$X_w(a,b) = \frac{1}{|a|^{1/2}} \int_{-\infty}^{\infty} x(t) \Psi^* \left(\frac{t-b}{a} \right) dt \quad (3)$$

Given the impossibility to analyze the data by using all wavelet coefficients, usually it is sufficient to consider a discrete subset of the upper half-plane to be able to reconstruct the data from the corresponding wavelet coefficients. The considered discrete subset of the half-plane are all the points $(a^m, na^m b)$, where $m, n \in \mathbb{Z}$, and this allows us to define the so-called *child wavelets* as shown in Equation 4.

$$\Psi_{m,n}(t) = \frac{1}{\sqrt{a^m}} \Psi \left(\frac{t-nb}{a^m} \right) \quad (4)$$

The use of small scales (i.e., that corresponds to large frequencies, since the scale is given by the formula $\frac{1}{\text{frequency}}$) compresses the data, giving us an overview of the involved information, while large scales (i.e., low frequencies) expand the data, offering a detailed analysis of the information. On the basis of the characteristics of the wavelet transformation, although it is possible to use many basis functions as *mother wavelet* (e.g., *Daubechies*, *Meyer*, *Symlets*, *Coiflets*, etc), for the scope of our approach we decided to use one of the simplest and oldest formalization of wavelet, the *Haar wavelet* [46]. We perform this choice because the *Haar wavelet* has the capability of measuring the contrast directly from the responses of low and high frequency subbands. This mother wavelet is shown in Equation 5.

$$\Psi(t) = \begin{cases} 1, & 0 \leq t < \frac{1}{2} \\ -1, & \frac{1}{2} \leq t < 1 \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

For exemplification purposes, considering a *time series* $TS = \{ts_1, ts_2, \dots, ts_N\}$, for instance $TS = \{8, 5, 6, 7, 5, 4, 6, 5\}$ (then with $|TS| = N = 8$), the transformation operated by using the pyramid algorithm of *Haar wavelet* in order to obtain a representation of data based on the average, gives the values reported in Equation 6 as result.

$$\Psi(TS) = \{6.5, 6.5, 4.5, 5.5\} \quad (6)$$

The result is obtained by following the criterion shown in the following Equation 7.

$$\frac{ts_2 + ts_1}{2}, ts_2 = \frac{ts_4 + ts_3}{2}, ts_3 = \frac{ts_6 + ts_5}{2}, ts_4 = \frac{ts_8 + ts_7}{2} \quad (7)$$

We can apply the *Haar wavelet* function on the *time series* multiple times, reducing the result length according to the sequence $\frac{N}{2}, \frac{N}{4}, \frac{N}{8}$, and so on. Such process reduces the level of detail and increases the overview on the data. The *Haar wavelet* function assumes that the length of the input is 2^n , with $n > 1$. When this is not possible, other solutions can be used to overcome this problem, e.g., the *Ancient Egyptian Decomposition* process [47].

2.4.4. Competitor

Taking into account that the most effective fraud detection approaches at the state of the art need to train their model by using both the *fraudulent* and *legitimate* previous cases, in this paper we do not compare our approach to many of them, limiting the comparison to only one of the most used and effective ones, being *Random Forests* [12]. The *Random Forests* approach represents one of the most common and powerful state-of-the-art techniques for data analysis, because in most of the cases it outperforms the other ones [48, 35, 49].

It consists in an ensemble learning method for classification and regression based on the construction of a number of randomized decision trees during the training phase. The conclusion are inferred by averaging the obtained results and this technique can be used to solve a wide range of prediction problems, with the advantage that it does not need any complex configuration, because it only requires the adjustment of two parameters: the number of trees and the number of attributes used to grow each tree.

Our aim is to prove that through our approach is possible to define effective evaluation models built by using only a class of transactions (i.e., the *legitimate* one), granting some advantages.

3. Preliminaries

This section provides the formal notation adopted in this paper and some basic assumptions, as well as the formal definition of the faced problem.

3.1. Formal Notation

The formal notation adopted in this paper is reported in Table 1. It should be observed that a transaction can only belong to one class $c \in C$.

Table 1: Formal Notation

Notation	Description	Note
$T = \{t_1, t_2, \dots, t_N\}$	Set of classified transactions	
$T_+ = \{t_1, t_2, \dots, t_K\}$	Subset of legitimate transactions	$T_+ \subseteq T$
$T_- = \{t_1, t_2, \dots, t_J\}$	Subset of fraudulent transactions	$T_- \subseteq T$
$V = \{v_1, v_2, \dots, v_M\}$	Set of transaction features	
$\hat{T} = \{\hat{t}_1, \hat{t}_2, \dots, \hat{t}_U\}$	Set of unclassified transactions	
$C = \{\textit{legitimate}, \textit{fraudulent}\}$	Set of possible classifications	
$F = \{f_1, f_2, \dots, f_X\}$	Output of DFT or DWT process	

3.2. Problem Definition

Denoting as Ξ the process of comparison between the *DFT* (or *DWT*) output of the *time series* in the set T_+ (i.e., the sequence of feature values in the previous *legitimate* transactions) and the *DFT* (or *DWT*) output of the *time series* related to the unevaluated transactions in the set \hat{T} (processed one at a time), the objective of our approach is the classification of each transaction $\hat{t} \in \hat{T}$ as *legitimate* or *fraudulent*.

Defining a function $EVAL(\hat{t}, \Xi)$ that performs this operation based on our approach, returning a boolean value β (0 =*misclassification*, 1 =*correct classification*) for each classification, we can formalize our objective function (Equation 8) in terms of maximization of the results sum.

$$\max_{0 \leq \beta \leq |\hat{T}|} \beta = \sum_{u=1}^{|\hat{T}|} EVAL(\hat{t}_u, \Xi) \quad (8)$$

4. Proposed Approach

The proposed approach was implemented by performing the steps listed below and detailed in the Sections 4.1, 4.2, and 4.3.

1. **Data Definition:** definition of the *time series* to use as input in the *DFT* or *DWT* process, in terms of sequence of values assumed by the transaction features;
2. **Data Processing:** conversion of the *time series* obtained in the previous step into a *transformed-domain* by using the *DFT* or *DWT* process;
3. **Data Evaluation:** formalization of the algorithm able to classify a new transaction as *legitimate* or *fraudulent* on the basis of a comparison process made in the *transformed-domain*.

4.1. Data Definition

As previously introduced in Section 2.4.1, a *time series* is a sequence of data points stored by following the time order and, in most of the cases, it is a sequence of discrete-time data measured at successive equally spaced points in time.

In the context of our approach, we considered as *time series* (ts) the sequence of values $v \in V$ assumed by the features of the transactions in T_+ and \hat{T} , as shown in Equation 9 and Equation 10.

$$T_+ = \begin{bmatrix} v_{1,1} & v_{1,2} & \dots & v_{1,M} \\ v_{2,1} & v_{2,2} & \dots & v_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ v_{K,1} & v_{K,2} & \dots & v_{K,M} \end{bmatrix} \quad \hat{T} = \begin{bmatrix} v_{1,1} & v_{1,2} & \dots & v_{1,M} \\ v_{2,1} & v_{2,2} & \dots & v_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ v_{U,1} & v_{U,2} & \dots & v_{U,M} \end{bmatrix} \quad (9)$$

$$\begin{aligned} ts(T_+) &= (v_{1,1}, v_{1,2}, \dots, v_{1,M}), (v_{2,1}, v_{2,2}, \dots, v_{2,M}), \dots, (v_{K,1}, v_{K,2}, \dots, v_{K,M}) \\ ts(\hat{T}) &= (v_{1,1}, v_{1,2}, \dots, v_{1,M}), (v_{2,1}, v_{2,2}, \dots, v_{2,M}), \dots, (v_{U,1}, v_{U,2}, \dots, v_{U,M}) \end{aligned} \quad (10)$$

The *time series* related to an item $\hat{t} \in \hat{T}$ will be compared to the *time series* related to all the items $t_+ \in T_+$, by following the criteria explained in the next steps.

An alternative method defines the *time series* in terms of sequence of values assumed by each transaction feature $v \in V$ in the set T_+ and \hat{T} , as shown in Equation 11. Such different modality is suitable when the aim of the evaluation model is to detect atypical values in a single feature, rather than in the whole set of features (i.e., as it occurs in the context taken into account in this paper).

$$\begin{aligned} ts(T_+) &= (v_{1,1}, v_{2,1}, \dots, v_{K,1}), (v_{1,2}, v_{2,2}, \dots, v_{K,2}), \dots, (v_{1,M}, v_{2,M}, \dots, v_{K,M}) \\ ts(\hat{T}) &= (v_{1,1}, v_{2,1}, \dots, v_{U,1}), (v_{1,2}, v_{2,2}, \dots, v_{U,2}), \dots, (v_{1,M}, v_{2,M}, \dots, v_{U,M}) \end{aligned} \quad (11)$$

4.2. Data Processing

The *time series* defined in the previous step are here processed in order to move their representation to the *transformed-domain*, by using the *DFT* or *DWT* process.

In a preliminary study we compared some patterns in the time domain (i.e., the *time series*) to their representation in the *transformed-domain*. Without going deep into the merits of the formal characteristics of *Fourier* and *Wavelet* transformations, but by limiting our analysis to the context taken into account, we underlined the properties described below:

4.2.1. Exploited Fourier Properties

1. **Phase invariance:** the first property, shown in Figure 3, demonstrates that there are not variations in the spectral pattern in case of a value translation⁴. More formally, it is one of the *phase properties* of the Fourier transform [50], i.e., a shift of a *time series* in the time domain leaves the magnitude unchanged in the *transformed-domain* [50]. It means that the representation in the *transformed-domain* allows us to detect a specific pattern, regardless of the position of the values assumed by the transaction features that originate it;
2. **Amplitude correlation:** the second property, shown in Figure 4, instead proves the existence of a direct correlation between the values assumed by the features in the time domain and the corresponding magnitudes assumed by the spectral components in the *transformed-domain*. More formally, it is the *homogeneity property* of the Fourier transform [50], i.e., when the amplitude is altered in one domain, it is altered by the same entity in the other domain⁵. This ensures that the proposed approach is able to evaluate the differences in terms of feature values, i.e., it is able to differentiate identical spectral patterns on the basis of the values assumed by their transaction features;
3. **Additivity quality:** another interesting property, shown in Figure 5, allows us to define patterns able to represent particular user behaviors, simply by adding the *time series* related to the involved transactions. More formally, it represents the *additivity property* of the Fourier transform [50], i.e., to the addition in the time domain corresponds an addition in the frequency domain. It means that we can

⁴A translation in time domain corresponds to a change in phase in the frequency domain.

⁵Scaling in one domain corresponds to scaling in the other domain

merge two patterns in the time domain, without losing information in the spectral pattern representation.

By using the *Fourier* approach, in this step we move the *time series* of the transactions to the *transformed-domain* by a *DFT* process performed through the *FFT* algorithm introduced in Section 2.4.2 and Section 2.4.3.. Basically, we extract the spectral pattern of each transaction by processing the related *time series* defined in the previous step.

4.2.2. Exploited Wavelet Properties

1. **Dimensionality reduction:** the *DWT* process represents an effective method for the *time series* data reduction, since the orthonormal transformation operated reduces the dimensionality of a *time series*, providing a compact representations of data, which however preserves the original information in its coefficients. By exploiting this property a fraud detection system can reduce the computational complexity of the involved processes;
2. **Multiresolution analysis:** applied on the *time series* context, the *DWT* allows us to define separate *time series* on the basis of the original one, distributing the information in these new representations of data in terms of the wavelet coefficient. The most important aspect of such transformations is that the *DWT* process performs an orthonormal transformation, preserving the original information, allowing us to restore the original data representation. A fraud detection system can exploit this mechanism in order to detect rapid changes in the data under analysis, observing the *data series* under two different points of view (i.e., types of wavelet coefficient), an approximated and a detailed one. The approximate point of view provides an overview on the data, while the detailed point of view provides information useful to evaluate data changes.

By using the *Wavelet* approach, in this step we transform the original *time series* given by the sequence of values assumed by the transaction features (as explained in Section 4.1) by performing the *Haar wavelet* process described in Section 2.4.3. The approximation coefficients of $\frac{N}{2}$ level is preferred to a detailed one in order to define a more stable model (i.e., less influenced by the data heterogeneity) for the evaluation of the new transactions.

4.3. Data Evaluation

The process of evaluation of a new transaction is performed by comparing the *DFT* or *DWT* outputs of the previous *legitimate* transactions to those of the transactions to evaluate.

For each transaction $\hat{t} \in \hat{T}$ we compare its *transformed-domain* representation $F(\hat{t})$ (i.e., the series of values $f \in F$) to the *transformed-domain* representation $F(t_+)$ of each *legitimate* previous transaction $t_+ \in T_+$.

The comparison process can be done in the *transformed-domain* (i.e., *DFT* or *DWT* outputs vectors) by using one of the two different methods described in the following:

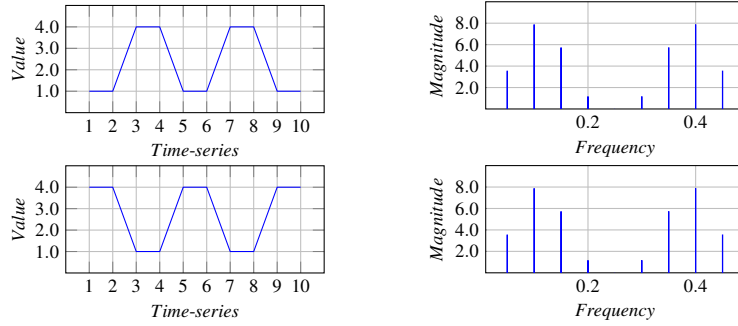


Figure 3: *Fourier : Phase Invariance*

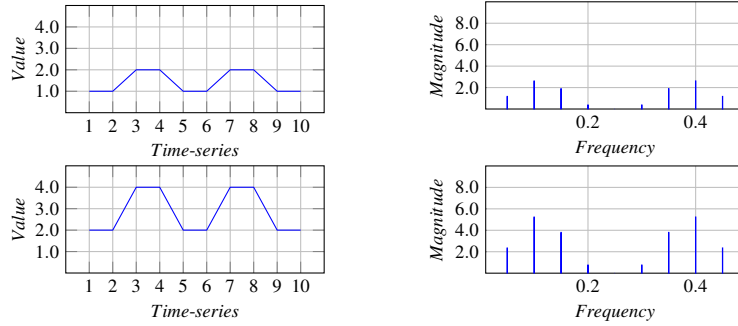


Figure 4: *Fourier : Amplitude Correlation*

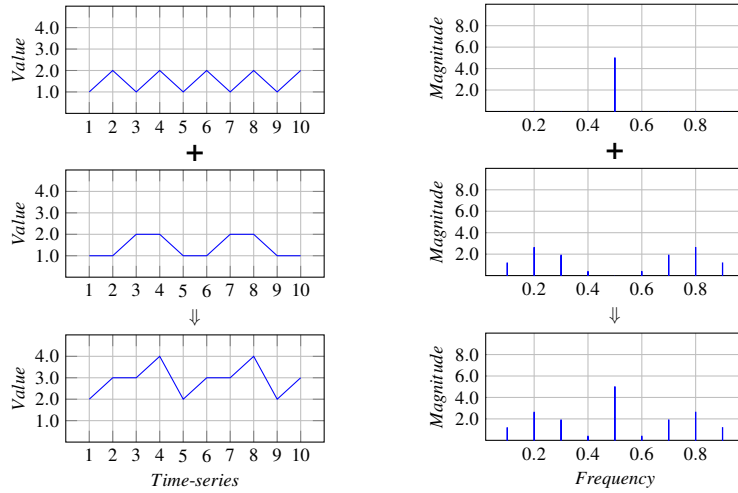


Figure 5: *Fourier : Additivity Quality*

1. the first method is based on the *cosine similarity*, a well-known metric described in Section 5.3.1. This first method is suitable when we need to evaluate the similarity between transactions in a global manner, thus by jointly evaluating the behavior of all the elements that compose the output vectors of the *DFT* or *DWT* process;
2. the second method is based on the *punctual comparison* between the values assumed by each element of the output vectors, with regard to the minimum or maximum value assumed by the element in the dataset (the result will be a boolean value, then 0 or 1). The similarity is evaluated with respect to a threshold, e.g., a transaction is considered similar to another one, when the sum of the comparison results of all the elements is above the $\frac{|F|}{2}$ value. Such method is suitable when we need to evaluate the similarity on the basis of the behavior of each single feature (e.g., in some *Intrusion Detection Systems* [51]).

For the aforementioned considerations, in our approaches we adopt the first method (i.e., *cosine similarity*), as shown in Equation 12, where Δ represents the similarity value in terms of *cosine similarity*, α is a threshold value experimentally defined in Section 5.4.2, and c is the resulting classification.

We repeat this process by using the transaction to evaluate \hat{t} and all the transactions $t_+ \in T_+$, obtaining the final classification of the transaction by averaging over these comparisons.

$$\Delta = \cos(F(t), F(\hat{t})), \text{ with } c = \begin{cases} \Delta \geq \alpha, & \text{legitimate} \\ \Delta < \alpha, & \text{fraudulent} \end{cases} \quad (12)$$

4.3.1. Algorithm

The final classification of a new transaction \hat{t} , which takes into account all the comparisons (Equation 12) between the transaction \hat{t} and all the transactions in T_+ , is performed by using the Algorithm 1.

This process takes as input the set T_+ of past *legitimate* transactions, a transaction \hat{t} to evaluate, and the threshold value α to use in the spectral pattern comparison process (i.e., in the context of the *cosine similarity* evaluation). It returns as output a boolean value that indicates the \hat{t} classification (i.e., *true=legitimate* or *false=fraudulent*).

From *step 1* to *step 16* we process the unevaluated transaction \hat{t} , by starting with the definition of the *time series* related to the transaction \hat{t} (*step 2*), moving it in the *transformed-domain* (*step 3*).

In the *steps* from 4 to 8, we compare in the *transformed-domain* the transaction \hat{t} to that of each transaction $t_+ \in T_+$ (obtained at the *steps 5* and *6*), adding the result (i.e., the *cosine similarity* value) to the variable *cos* (*step 7*).

The average of the final value of the variable *cos* (*step 9*) is compared to the threshold value α (*steps* from 10 to 14), and the final classification of the transaction \hat{t} , returned by the algorithm at the *step 15*, depends on the result of this operation.

4.3.2. Complexity

Here we calculate the cost in time needed to perform a classification of a single transaction \hat{t} , since this type of information allows us to evaluate the performance of

Algorithm 1 Transaction evaluation

Input: T_+ =Legitimate previous transactions, \hat{t} =Unevaluated transaction, α =Threshold value

Output: β =Classification of the transaction \hat{t}

```
1: procedure TRANSACTIONEVALUATION( $T_+$ ,  $\hat{t}$ )
2:    $ts1 \leftarrow getTimeseries(\hat{t})$ 
3:    $sp1 \leftarrow getTransformedDomain(ts1)$ 
4:   for each  $t_+$  in  $T_+$  do
5:      $ts2 \leftarrow getTimeseries(t_+)$ 
6:      $sp2 \leftarrow getTransformedDomain(ts2)$ 
7:      $cos \leftarrow cos + getCosineSimilarity(sp1, sp2)$ 
8:   end for
9:    $avg \leftarrow \frac{cos}{|T_+|}$ 
10:  if  $avg > \alpha$  then
11:     $\beta \leftarrow true$ 
12:  else
13:     $\beta \leftarrow false$ 
14:  end if
15:  return  $\beta$ 
16: end procedure
```

the proposed approach in the context of a *real-time system* [52], a scenario where the *response-time* represents a primary aspect.

We perform this operation by analyzing the theoretical complexity of the classification Algorithm 1, previously formalized. Denoted as N the dimension of the set T_+ (i.e., $N = K = |T_+|$), the asymptotic time complexity of a single evaluation, in terms of *Big O notation*, can be determined on the basis of the following observations:

- (i) as shown in Figure 6, the Algorithm 1 presents two nested loops given by the outer loop that starts at *step 4* (*L1* loop), which executes N times the inner loop *L2*, plus other operations (*getTimeseries* and *getTransformedDomain*), respectively with complexity $O(n)$ and $O(n \log n)$;
- (ii) it represents the worst case, since the *Discrete Wavelet Transform* takes only $O(n)$ in certain cases, as compared to $O(n \log n)$ takes by the *Fast Fourier Transform* algorithm (i.e., that we use in order to perform the *Discrete Fourier Transform*);
- (iii) the inner loop *L2* performs the comparison of the spectral patterns by recurring to the *cosine similarity* metric, which is characterized by a $O(N^2)$ complexity;
- (iv) the other involved operations of *comparisons* and *assignments* are characterized by a $O(1)$ complexity.

The aforementioned considerations allow us to determine that the asymptotic time complexity of the proposed algorithm is $O(N^2)$, a complexity that can be effectively reduced by parallelizing the process over several machines, e.g., by exploiting large scale distributed computing models.

5. Experiments

This section reports information about the experimental environment, the used datasets and metrics, the adopted strategy, as well as the results of the performed experiments.

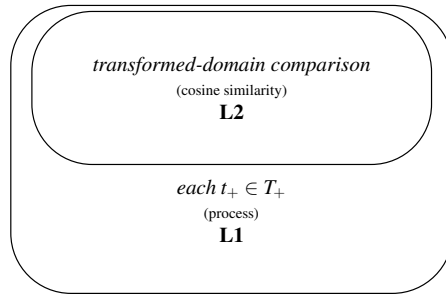


Figure 6: Algorithm Nested Loops

5.1. Environment

The proposed approach was developed in Java, where we use the *JTransforms*⁶ library to operate the Fourier transformations, and the *JWave*⁷ library to operate the Wavelet transformations..

The state-of-the-art approach and the metrics used to evaluate its results were implemented in R^8 , by using *randomForest*, *DMwR*, and *ROCR* packages.

It should be further added that we verified the existence of a statistical difference between the results, by using the independent-samples *two-tailed Student's t-tests* ($p < 0.05$).

5.2. DataSets

The two real-world datasets used in the experiments (i.e., *European Transactions*⁹ and *German Credit*¹⁰) represent two benchmarks in this research field. We chose two datasets with different levels of data imbalance, whose characteristics are described in the following.

5.2.1. European Transactions (ET)

This dataset contains the transactions carried out in two days of September 2013, for a total of 492 frauds out of 284,807 transactions. It should be observed how this represents an highly unbalanced dataset [53], considering that the *fraudulent* cases are only the 0.0017% of all the transactions.

For confidentiality reasons, all fields of the dataset have been anonymized, except the *time* (that we do not take into account in the Fourier transformation process) and *amount* features that report, respectively, the number of seconds elapsed between the first transaction in the dataset and the current transaction, and the amount of the credit card transaction. As usual, the last field contains the transaction classification ($0=legitimate$ and $1=fraudulent$).

⁶<https://sourceforge.net/projects/jtransforms/>

⁷<https://github.com/cscheiblich/JWave/>

⁸<https://www.r-project.org/>

⁹<https://www.kaggle.com/dalpozz/creditcardfraud/>

¹⁰<ftp://ftp.ics.uci.edu/pub/machine-learning-databases/statlog/>

5.2.2. German Credit (GC)

This dataset is composed by 1,000 transactions and 300 of them are frauds. Also in this case it represents an unbalanced dataset, since the *fraudulent* cases are the 30.00% of all the transactions.

The dataset is released with all the features modified for confidentiality reasons, and we used the version with all numeric features (i.e., without categorical variables). Each transaction is composed by 20 fields, plus a classification field (1=*legitimate* and 2=*fraudulent*).

5.3. Metrics

This section introduces the metrics used in the context of this paper.

5.3.1. Cosine Similarity

The *cosine similarity* (*Cosim*) between two non-zero vectors \vec{v}_1 and \vec{v}_2 is calculated in terms of cosine angle between them, as shown in the Equation (13).

It represents a widespread measure that allows us to evaluate the similarity between two transaction patterns by comparing the vectors given by the values of their components in the *transformed-domain*.

$$\text{Cosim}(\vec{v}_1, \vec{v}_2) = \cos(\vec{v}_1, \vec{v}_2) = \frac{\vec{v}_1 \cdot \vec{v}_2}{\|\vec{v}_1\| \cdot \|\vec{v}_2\|} \quad (13)$$

5.3.2. F-score

The *F-score* is considered an effective performance measures for unbalanced datasets [53, 54]. It represents the weighted average of the *Precision* and *Recall* metrics and it is a largely used metric in the statistical analysis of binary classification, returning a value in a range $[0, 1]$, where 0 is the worst value and 1 the best one.

More formally, given two sets $T^{(P)}$ and $T^{(R)}$, where $T^{(P)}$ denotes the set of performed classifications of transactions, and $T^{(R)}$ the set that contains the actual classifications of them, this metric is defined as shown in Equation 14.

$$\begin{aligned} F\text{-score}(T^{(P)}, T^{(R)}) &= 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \\ \text{with} \\ \text{Precision}(T^{(P)}, T^{(R)}) &= \frac{|T^{(R)} \cap T^{(P)}|}{|T^{(P)}|} \\ \text{Recall}(T^{(P)}, T^{(R)}) &= \frac{|T^{(R)} \cap T^{(P)}|}{|T^{(R)}|} \end{aligned} \quad (14)$$

5.3.3. AUC

The *Area Under the Receiver Operating Characteristic* curve (*AUC*) is a performance measure used to evaluate the effectiveness of a classification model [55]. Its result is in a range $[0, 1]$, where 1 indicates the best performance.

More formally, according to the notation of Section 3.1, given the subset of previous *legitimate* transactions T_+ and the subset of previous *fraudulent* ones T_- , the formalization of the *AUC* metric is reported in the Equation 15, where Θ indicates all

possible comparisons between the transactions of the two subsets T_+ and T_- . It should be noted that the result is obtained by averaging over these comparisons.

$$\Theta(t_+, t_-) = \begin{cases} 1, & \text{if } t_+ > t_- \\ 0.5, & \text{if } t_+ = t_- \\ 0, & \text{if } t_+ < t_- \end{cases} \quad AUC = \frac{1}{|T_+||T_-|} \sum_1^{|T_+|} \sum_1^{|T_-|} \Theta(t_+, t_-) \quad (15)$$

5.4. Strategy

This section provides information about the strategy adopted during the execution of the experiments.

5.4.1. Cross-validation

In order to reduce the impact of data dependency, improving the reliability of the obtained results, all the experiments have been performed by using the *k-fold cross-validation* criterion, with $k=10$.

Each dataset is divided in k subsets, and each k subset is used as the test set, while the other $k-1$ subsets are used as the training set. The final result is given by the average of all results.

5.4.2. Threshold Tuning

Before starting the experiments we carried out a study aimed to identify the best value of the threshold parameter α to use in the evaluation process, according to the Equation 12.

In order to maintain a proactive approach, we perform this operation by using only the *legitimate* transactions in the dataset, calculating the average value of the *cosine similarity* related to all pairs of different transactions $t_+ \in T_+$, according to the Algorithm 2.

Through this process we want to obtain the average value of *cosine similarity* measured between the *transformed-domain* representation of all pairs of previous *legitimate* transactions, so that we can use it to identify the different ones (i.e., the potential *fraudulent* transactions).

It takes as input the set T_+ of past *legitimate* transactions and returns the threshold value α . The two nested loops that start at *step 2* and at *step 3* select only the different transaction pairs (*step 4*) in the set T_+ . For these pairs (i.e., t'_+ and t''_+), we calculate their *time series* and we move them in the *transformed-domain* (*steps* from 6 to 9), then we sum the *cosine similarity* between them (*step 10*). The average of all *cosine similarity* evaluations is calculated at *step 14* and it is returned by the algorithm at the *step 15*.

The evaluation was stopped when the value of α did not present significant variations. In both datasets, the results indicate $\alpha = 0.90$ as the optimal threshold to use in the *DFT* approach and $\alpha = 0.91$ as the optimal threshold to use in the *DWT* approach.

5.5. Competitor

As introduced in Section 5.1, we compare our approach to *Random Forests*.

Algorithm 2 *Threshold tuning*

Input: T_+ =Legitimate previous transactions

Output: α =Threshold value

```
1: procedure GETALPHA( $T_+$ )
2:   for each  $t_+^i$  in  $T_+$  do
3:     for each  $t_+^j$  in  $T_+$  do
4:       if  $t_+^i \neq t_+^j$  then
5:         evaluations  $\leftarrow$  evaluations + 1
6:         ts1  $\leftarrow$  getTimeseries( $t_+^i$ )
7:         sp1  $\leftarrow$  getTransformedDomain(ts1)
8:         ts2  $\leftarrow$  getTimeseries( $t_+^j$ )
9:         sp2  $\leftarrow$  getTransformedDomain(ts2)
10:        cos  $\leftarrow$  cos + getCosineSimilarity(sp1, sp2)
11:      end if
12:    end for
13:  end for
14:   $\alpha \leftarrow \frac{\text{cos}}{\text{evaluations}}$ 
15:  return  $\alpha$ 
16: end procedure
```

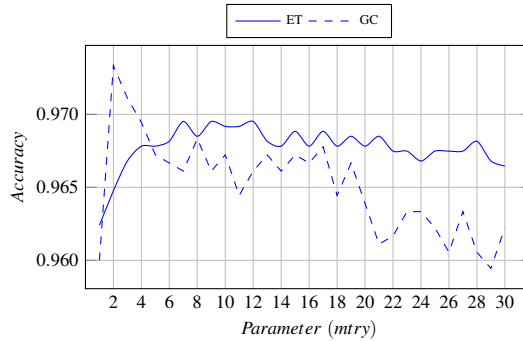


Figure 7: *Random Forests Tuning*

5.5.1. Description

It is implemented in *R* language, by using the *randomForest* and *DMwR* packages.

The *DMwR* package allows *Random Forests* to manage the class imbalance problem through the *Synthetic Minority Over-sampling Technique (SMOTE)* [56]. It represents a very popular sampling technique able to create new synthetic data by randomly interpolating pairs of nearest neighbors.

The combined use of *Random Forests* and *SMOTE* allows us to verify the performance of our approach compared to one of the best solutions for fraud detection at the state of the art.

For reasons of reproducibility of the *RF* experiments, the *R* function *set.seed()* has been used in the code to fix the seed of the random number generator. The *RF* parameters have been experimentally tuned by searching those that maximize the performance.

5.5.2. Tuning

In order to maximize the performance of the *RF* approach, we need to detect the optimal value of the *mtry* (number of variables randomly sampled as candidates at each

split) and the *ntree* (number of trees to grow) parameters.

We performed this operation by exploiting the *caret R* package, which provides an excellent tuning functionality. It supports only those algorithm parameters that have a crucial role in the tuning process, such as the *mtry* one. We proceeded by using the so-called *grid search* approach, where each axis of the grid represents an algorithm parameter and the values in the grid represent specific parameters combinations.

In more detail, we exploited the *SMOTE* technique functionalities (implemented through the *DMwR* package), in order to preprocess the original unbalanced dataset, obtaining as result a balanced dataset to use for the tuning process. This has been done in order to replicate the fraud detection operative context during the parameter tuning.

About the *SMOTE* configuration, we set to 200 the *perc.over* parameter¹¹, and we set to 150 the *perc.under* parameter¹²

The obtained datasets was used in order to tune the *mtry* parameter, obtaining the results shown in Figure 7, which indicates *mtry* = 12 as optimal value for the *ET* dataset, since it is the value that leads toward the maximum *Accuracy* (i.e., 0.969%), and *mtry* = 2 as optimal value for the *GC* dataset, since it is the value that leads toward the maximum *Accuracy* (i.e., 0.973%).

5.6. Results

The observations that arise by examining the experimental results are summarized and discussed in this section.

5.6.1. Overview

- (i) the first set of experiments was focused on the evaluation of the proposed approach in terms of *F-score*. The results, shown in Figure 8, indicate that both the *DFT* and *DWT* performance are similar to that of *Random Forests (RFS)*, in the context of the two datasets taken into account. This happens despite our approaches do not use any previous *fraudulent* transaction to train their models, adopting a pure proactive strategy. It means that they are able to operate without training their models with both classes of transactions (*legitimate* and *fraudulent*).
- (ii) the second set of experiments was aimed to evaluate the performance of the *DFT* and *DWT* approaches in terms of *AUC*. As described in Section 5.3.3, this metric evaluates the predictive power of a classification model, and the results in Figure 9 show how our model achieves performance close to that of *RFS*, in the context of both datasets, although they get better performance with the *ET* dataset. This is because they define their models on the basis of legitimate cases, therefore a greater number of these cases allows them to achieve better performance.

¹¹A number that drives the *over-sampling*, i.e., how many extra cases from the minority class we want to create.

¹²A number that drives the *under-sampling*, i.e., how many extra cases from the majority class are selected for each case generated from the minority class.

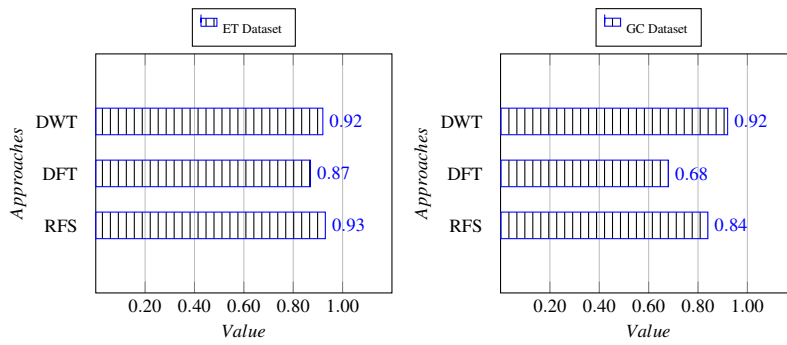


Figure 8: *F-score Performance*

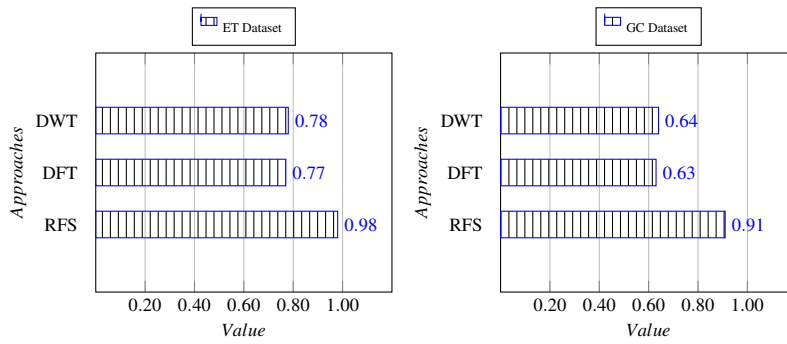


Figure 9: *AUC Performance*

5.6.2. Discussion

In the light of the obtained results, we can observe how the modelization of the transactions made in a *transformed-domain* is able to face the problems related to the *non-adaptability* and the *heterogeneity issues* described in Section 2.3, thanks to the stability offered by the new data representation.

We can also observe how the proactive strategy followed by our approach is able to reduce/overcome the *data imbalance* and *cold-start* issues, which are also described in Section 2.3, since only a class of transactions has been used.

The most important aspect of such proactivity is represented by the fact that it allows a real-world fraud detection system to operate even in the absence of previous *fraudulent* cases, with all the obvious advantages that derive from it.

With regard to the retroactive and proactive aspect of the fraud detection techniques, it should be evaluated on the basis of the operative scenario. In the scenario taken into account it is in fact reasonable to tolerate that a proactive approach gets worse performance than that of a retroactive one.

This statement is based on the consideration that through a proactive approach a fraud detection system can operate without the need to collect a number of *fraudulent* transactions to use for the model training (used by the retroactive approaches), reducing the economic losses.

It should also be added how the proposed proactive approaches can be considered in the context of hybrid techniques, since their correct classifications can be used in order to improve the effectiveness of the canonical retroactive state-of-the-art approaches, reducing the data unbalance issue. This means that a combined approach, which adopts retroactive and proactive techniques, can be used to design a very effective fraud detection system, where the capabilities of the single approaches are optimized.

They indicate that the differences between the competitor and our best performing approach (i.e., *DWT*) are really minimal, despite it adopts a pure proactive strategy. This minimum difference in performance must be further reduced in the light of the fact that the misclassifications made by our approaches do not necessarily lead toward loss of money, as they are related to both *false positive* and *false negative* cases.

Summarizing, through the adoption of proactive approaches, such as those proposed in this paper, we contrast the issues discussed in Section 2.3, as their processes do not involve *fraudulent* examples (facing the *data scarcity* and *data unbalance* issues), adopting a *transform-domain-based* model able to well characterize a specific class of transactions (i.e., the *legitimate* one) that results less influenced by the data variation (facing the *non-adaptability* and *data heterogeneity* issues), presenting the positive side effect of solving the *cold-start* issue.

6. Conclusions and Future Work

Today, the sustainability represent an imperative paradigm to preserve the planet's resources. In this context, the new technologies allow us, in a more or less direct way, to adopt this paradigm in many day-to-day choices.

The research that stand behind the *Big Data Analytics for Sustainability* is a representative example of such scenario, since it aims to offer solutions able to allow the people to exploit the new technologies in a smarter and securely way.

Table 2: Performance

Dataset	Approach	F-score	AUC
ET	DWT	- 0.01	- 0.20
ET	DFT	- 0.06	- 0.21
GC	DWT	+ 0.08	- 0.27
GC	DFT	- 0.16	- 0.28

The exponential growth in the number of sellers and buyers who work through the E-commerce platform offers people the opportunity to make their own choices also by following non-traditional paradigms such as, for instance, that of sustainability. However, this scenario is jeopardized by the risks related to the fraudulent use of electronic payment instruments, which represent the most common payment means in such environment.

For the aforementioned reason, the research that revolves around the *Big Data Information Security*, in this case that aimed to define effective fraud detection systems, assumes an increasingly central role, involving large investments by public and private entities. The risk scenario under consideration is mainly given by the combination of two factors: the exponential growth in the use of the E-commerce environment and the exponential growth in the use of credit cards by people.

Considering that, as a result of these two factors, even money losses have reported an exponential trend in recent years, through this paper we wanted to investigate the benefits given by the adoption of proactive strategies of fraud detection.

More than wanting to replace the existing retroactive state-of-the-art approaches, by recurring to the *Fourier Transform* or the *Wavelet Transform*, we introduced a novel proactive strategy, which performs the data analysis and the definition of the evaluation model in a new *transformed-domain*. Such proactivity allowed us to face some well-known issues that affect the canonical retroactive state-of-the-art approaches, the most important of which are the *data unbalance* and the *cold-start* ones.

The obtained results can be considered interesting, since it is necessary to consider that the state-of-the-art competitor taken into account (i.e., *Random Forests*), in addition to using both classes of transactions to train its model also preprocesses the dataset by using an effective balancing technique (i.e., *SMOTE*). The minimal differences in performance with regard to the retroactive state-of-the-art competitor approach, clearly indicate the capability of our proactive approaches to improve the fraud detection tasks, by operating stand-alone or by working in the context of a hybrid approach.

The use of the proposed proactive strategies can be considered a valuable contribution in several *BDAS* research fields, such as that of the *Big Data Information Security for Sustainability* previously mentioned, or that of the *Computational intelligence and algorithms for Sustainability*, since they allow us to improve the state-of-the-art solutions, providing them the capability to define an evaluation model on the basis of a single class of data.

For the aforementioned considerations, a possible future work could be focused on the definition of a new fraud detection approach that combines the characteristics of the canonical non-proactive state-of-the-art approaches with those of our proactive approaches, in order to define a hybrid strategy that maximizes the performance of both the approaches.

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* (2013).

References

- [1] V. Chang, Towards data analysis for weather cloud computing, *Knowl.-Based Syst.* 127 (2017) 29–45. doi:10.1016/j.knosys.2017.03.003. URL <https://doi.org/10.1016/j.knosys.2017.03.003>
- [2] M. V. M. Cano, F. Terroso-Saenz, A. González-Vidal, M. Valdés-Vela, A. F. Skarmeta, M. A. Zamora, V. Chang, Applicability of big data techniques to smart cities deployments, *IEEE Trans. Industrial Informatics* 13 (2) (2017) 800–809. doi:10.1109/TII.2016.2605581. URL <https://doi.org/10.1109/TII.2016.2605581>
- [3] T. Ferreira, I. Pedrosa, J. Bernardino, Business intelligence for e-commerce: Survey and research directions, in: Á. Rocha, A. M. R. Correia, H. Adeli, L. P. Reis, S. Costanzo (Eds.), *Recent Advances in Information Systems and Technologies - Volume 1 [WorldCIST'17, Porto Santo Island, Madeira, Portugal, April 11-13, 2017]*., Vol. 569 of *Advances in Intelligent Systems and Computing*, Springer, 2017, pp. 215–225. doi:10.1007/978-3-319-56535-4_22. URL https://doi.org/10.1007/978-3-319-56535-4_22
- [4] N. Japkowicz, S. Stephen, The class imbalance problem: A systematic study, *Intell. Data Anal.* 6 (5) (2002) 429–449.
- [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.
- [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.

- [7] 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.
- [8] M. J. Lenard, P. Alam, Application of fuzzy logic fraud detection, in: M. Khosrow-Pour (Ed.), Encyclopedia of Information Science and Technology (5 Volumes), Idea Group, 2005, pp. 135–139.
- [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] 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.
- [11] B. Hooi, N. Shah, A. Beutel, S. Günnemann, L. Akoglu, M. Kumar, D. Makhija, C. Faloutsos, BIRDNEST: bayesian inference for ratings-fraud detection, in: S. C. Venkatasubramanian, W. M. Jr. (Eds.), Proceedings of the 2016 SIAM International Conference on Data Mining, Miami, Florida, USA, May 5-7, 2016, SIAM, 2016, pp. 495–503. doi:10.1137/1.9781611974348.56.
URL <https://doi.org/10.1137/1.9781611974348.56>
- [12] L. Breiman, Random forests, *Machine Learning* 45 (1) (2001) 5–32. doi:10.1023/A:1010933404324.
- [13] P. Duhamel, M. Vetterli, Fast fourier transforms: a tutorial review and a state of the art, *Signal processing* 19 (4) (1990) 259–299.
- [14] P. Chaovalit, A. Gangopadhyay, G. Karabatis, Z. Chen, Discrete wavelet transform-based time series analysis and mining, *ACM Comput. Surv.* 43 (2) (2011) 6:1–6:37. doi:10.1145/1883612.1883613.
URL <http://doi.acm.org/10.1145/1883612.1883613>
- [15] V. Chang, The business intelligence as a service in the cloud, *Future Generation Comp. Syst.* 37 (2014) 512–534. doi:10.1016/j.future.2013.12.028.
URL <https://doi.org/10.1016/j.future.2013.12.028>
- [16] E. J. Keogh, A decade of progress in indexing and mining large time series databases, in: U. Dayal, K. Whang, D. B. Lomet, G. Alonso, G. M. Lohman, M. L. Kersten, S. K. Cha, Y. Kim (Eds.), Proceedings of the 32nd International Conference on Very Large Data Bases, Seoul, Korea, September 12-15, 2006, ACM, 2006, p. 1268.
URL <http://dl.acm.org/citation.cfm?id=1164262>

- [17] E. Baidoo, J. L. Priestley, An analysis of accuracy using logistic regression and time series.
- [18] K. R. Lai, C. Fan, W. Huang, P. Chang, Evolving and clustering fuzzy decision tree for financial time series data forecasting, *Expert Syst. Appl.* 36 (2) (2009) 3761–3773. doi:10.1016/j.eswa.2008.02.025. URL <https://doi.org/10.1016/j.eswa.2008.02.025>
- [19] R. Saia, S. Carta, A frequency-domain-based pattern mining for credit card fraud detection, in: M. Ramachandran, V. M. Muñoz, V. Kantere, G. Wills, R. J. Walters, V. Chang (Eds.), *Proceedings of the 2nd International Conference on Internet of Things, Big Data and Security, IoTBDS 2017, Porto, Portugal, April 24-26, 2017*, SciTePress, 2017, pp. 386–391. doi:10.5220/0006361403860391. URL <https://doi.org/10.5220/0006361403860391>
- [20] D. Zheng, F. Li, T. Zhao, Self-adaptive statistical process control for anomaly detection in time series, *Expert Syst. Appl.* 57 (2016) 324–336. doi:10.1016/j.eswa.2016.03.029. URL <https://doi.org/10.1016/j.eswa.2016.03.029>
- [21] K. A. Salleh, L. Janczewski, Technological, organizational and environmental security and privacy issues of big data, *Procedia Computer Science* 100 (2016) 19 – 28, international Conference on ENTERprise Information Systems/International Conference on Project MANagement/International Conference on Health and Social Care Information Systems and Technologies, CENTERIS/ProjMAN / HCist 2016. doi:http://dx.doi.org/10.1016/j.procs.2016.09.119. URL <http://www.sciencedirect.com/science/article/pii/S1877050916322864>
- [22] N. G. Miloslavskaya, A. Makhmudova, Survey of big data information security, in: M. Younas, I. Awan, J. E. Haddad (Eds.), *4th IEEE International Conference on Future Internet of Things and Cloud Workshops, FiCloud Workshops 2016, Vienna, Austria, August 22-24, 2016*, IEEE Computer Society, 2016, pp. 133–138. doi:10.1109/W-FiCloud.2016.38. URL <https://doi.org/10.1109/W-FiCloud.2016.38>
- [23] S. Kokolakis, Privacy attitudes and privacy behaviour: A review of current research on the privacy paradox phenomenon, *Computers & Security* 64 (2017) 122–134. doi:10.1016/j.cose.2015.07.002. URL <https://doi.org/10.1016/j.cose.2015.07.002>
- [24] R. Saia, L. Boratto, S. Carta, Multiple behavioral models: A divide and conquer strategy to fraud detection in financial time series, in: A. L. N. Fred, J. L. G. Dietz, D. Aveiro, K. Liu, J. Filipe (Eds.), *KDIR 2015 - Proceedings of the International Conference on Knowledge Discovery and Information Retrieval, part of the 7th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management (IC3K 2015), Volume 1, Lisbon, Portugal, November 12-14, 2015*, SciTePress, 2015, pp. 496–503. doi:10.5220/0005637104960503. URL <https://doi.org/10.5220/0005637104960503>

- [25] R. J. Bolton, D. J. Hand, Statistical fraud detection: A review, *Statistical Science* (2002) 235–249.
- [26] 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.
- [27] 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.
- [28] 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.
- [29] 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.
- [30] S. Sorounejad, 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.
URL <http://arxiv.org/abs/1611.06439>
- [31] G. Garibotto, P. Murrieri, A. Capra, S. D. Muro, U. Petillo, F. Flammini, M. Esposito, C. Pragliola, G. D. Leo, R. Lengu, N. Mazzino, A. Paolillo, M. D’Urso, R. Vertucci, F. Narducci, S. Ricciardi, A. Casanova, G. Fenu, M. D. Mizio, M. Savastano, M. D. Capua, A. Ferone, White paper on industrial applications of computer vision and pattern recognition, in: *ICIAP (2)*, Vol. 8157 of *Lecture Notes in Computer Science*, Springer, 2013, pp. 721–730.
- [32] A. Chatterjee, A. Segev, Data manipulation in heterogeneous databases, *ACM SIGMOD Record* 20 (4) (1991) 64–68.
- [33] D. Che, M. S. Safran, Z. Peng, From big data to big data mining: Challenges, issues, and opportunities, in: B. Hong, X. Meng, L. Chen, W. Winiwarer, W. Song (Eds.), *Database Systems for Advanced Applications - 18th International Conference, DASFAA 2013, International Workshops: BDMA, SNSM, SeCoP, Wuhan, China, April 22-25, 2013. Proceedings*, Vol. 7827 of *Lecture Notes in Computer Science*, Springer, 2013, pp. 1–15. doi:10.1007/978-3-642-40270-8_1.
URL https://doi.org/10.1007/978-3-642-40270-8_1
- [34] A. Ghosh, R. K. De, S. K. Pal (Eds.), *Pattern Recognition and Machine Intelligence*, Second International Conference, PReMI 2007, Kolkata, India, Dec Vol. 4815 of *Lecture Notes in Computer Science*, Springer, 2007. doi:10.1007/978-3-540-77046-6.
URL <https://doi.org/10.1007/978-3-540-77046-6>

- [35] 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.
- [36] 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.
- [37] V. Vinciotti, D. J. Hand, Scorecard construction with unbalanced class sizes, *Journal of Iranian Statistical Society* 2 (2) (2003) 189–205.
- [38] A. I. Marqués, V. García, J. S. Sánchez, On the suitability of resampling techniques for the class imbalance problem in credit scoring, *JORS* 64 (7) (2013) 1060–1070. doi:10.1057/jors.2012.120. URL <http://dx.doi.org/10.1057/jors.2012.120>
- [39] S. F. Crone, S. Finlay, Instance sampling in credit scoring: An empirical study of sample size and balancing, *International Journal of Forecasting* 28 (1) (2012) 224–238.
- [40] O. Loyola-González, J. F. Martínez Trinidad, J. A. Carrasco-Ochoa, M. García-Borroto, Study of the impact of resampling methods for contrast pattern based classifiers in imbalanced databases, *Neurocomputing* 175 (2016) 935–947. doi:10.1016/j.neucom.2015.04.120. URL <https://doi.org/10.1016/j.neucom.2015.04.120>
- [41] P. Donmez, J. G. Carbonell, P. N. Bennett, Dual strategy active learning, in: *ECML*, Vol. 4701 of *Lecture Notes in Computer Science*, Springer, 2007, pp. 116–127.
- [42] 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>
- [43] R. Agrawal, C. Faloutsos, A. N. Swami, Efficient similarity search in sequence databases, in: D. B. Lomet (Ed.), *Foundations of Data Organization and Algorithms*, 4th International Conference, FODO'93, Chicago, Illinois, USA, October 13-15, 1993, Proceedings, Vol. 730 of *Lecture Notes in Computer Science*, Springer, 1993, pp. 69–84. doi:10.1007/3-540-57301-1_5. URL http://dx.doi.org/10.1007/3-540-57301-1_5
- [44] M. R. Chernick, Wavelet methods for time series analysis, *Technometrics* 43 (4) (2001) 491. doi:10.1198/tech.2001.s49. URL <http://dx.doi.org/10.1198/tech.2001.s49>
- [45] D. B. Percival, A. T. Walden, *Wavelet methods for time series analysis*, Vol. 4, Cambridge university press, 2006.

- [46] S. Mallat, A theory for multiresolution signal decomposition: The wavelet representation, *IEEE Trans. Pattern Anal. Mach. Intell.* 11 (7) (1989) 674–693. doi:10.1109/34.192463. URL <http://dx.doi.org/10.1109/34.192463>
- [47] P. M. Higgins, *Professor Higgins’s Problem Collection*, Oxford University Press, 2017.
- [48] 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.
- [49] 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>
- [50] S. W. Smith, et al., *The scientist and engineer’s guide to digital signal processing*.
- [51] A. L. Buczak, E. Guven, A survey of data mining and machine learning methods for cyber security intrusion detection, *IEEE Communications Surveys and Tutorials* 18 (2) (2016) 1153–1176. doi:10.1109/COMST.2015.2494502. URL <https://doi.org/10.1109/COMST.2015.2494502>
- [52] J. T. S. Quah, M. Sriganesh, Real-time credit card fraud detection using computational intelligence, *Expert Syst. Appl.* 35 (4) (2008) 1721–1732. doi:10.1016/j.eswa.2007.08.093. URL <http://dx.doi.org/10.1016/j.eswa.2007.08.093>
- [53] 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. URL <http://dx.doi.org/10.1109/SSCI.2015.33>
- [54] V. Chang, M. Ramachandran, Towards achieving data security with the cloud computing adoption framework, *IEEE Trans. Services Computing* 9 (1) (2016) 138–151. doi:10.1109/TSC.2015.2491281. URL <https://doi.org/10.1109/TSC.2015.2491281>
- [55] D. Faraggi, B. Reiser, Estimation of the area under the roc curve, *Statistics in medicine* 21 (20) (2002) 3093–3106.
- [56] K. W. Bowyer, N. V. Chawla, L. O. Hall, W. P. Kegelmeyer, SMOTE: synthetic minority over-sampling technique, *CoRR abs/1106.1813*. URL <http://arxiv.org/abs/1106.1813>