

# Auditing Government Purchases with a Multicriteria Anomaly Detection Strategy

Patrícia Maia<sup>1,2</sup>, Wagner Meira Jr.<sup>1</sup>, Breno Cerqueira<sup>2</sup>, Gustavo Cruz<sup>2</sup>

<sup>1</sup> Universidade Federal de Minas Gerais, Brazil  
meira@dcc.ufmg.br

<sup>2</sup> Controladoria Geral da União, Brazil  
patricia.maia,breno.cerqueira, gustavo.cruz,@cgu.gov.br

## Abstract.

Government purchases are the usual instrument for public acquisition of goods and services. Despite extensive legislation, several control and auditing mechanisms, frauds are still diverse and commonplace at all levels of public administration, wasting public resources. Through the use of frequent patterns, temporal correlation and combined analysis of multi-criteria, this work proposes a methodology for detecting anomalies in government purchases. The methodology promotes several levels of filtering with respect to entities involved and purchases are considered as fraudulent based on diverse criteria. The applicability and effectiveness of the methodology is demonstrated through a real case study where we were able to identify a long term provider collusion.

Categories and Subject Descriptors: H.2.8 [Database Management]: Database Applications; I.2.6 [Artificial Intelligence]: Learning

Keywords: anomaly detection, government purchases, data mining

## 1. INTRODUCTION

One of the main goals of government audit is to verify the effective competition among bidders, who must provide goods and services to the Public Administration on more advantageous terms and conditions. In the auction modality, instituted by Brazilian law 10.520/2002, the bids are determined by the providers following the terms of a reverse auction, i.e., seeking the lowest price to be contracted, as an incentive to competition. However, there are diverse evidences that the competitive nature of the bid may have been frustrated by artificially limiting the competition, among other issues. Without an analysis of massive data from a large number of auctions, it is quite difficult to characterize such anti-competitive behaviors, which are usually agreed among providers prior to the auction. Some of the commonly found evidences are listed next and will be further detailed in this paper: providers that sign contracts without offering the best price; dropout providers; disqualified providers; similar profiles among disqualified or dropout providers; most frequent reasons for disqualification; overpricing; market share of each provider; and prior assignment of providers to public agencies. A big challenge here is that the majority of the evidences can hardly be noticed on an individual analyses of the auction process, as a consequence of factors such as: the team that performs the audit in one year may be different from the team in the years that follow, which favors a ‘loss of memory’ of previous work; the scope of the audit may vary from one year to another, involving different targets over the years; services’ contracts can be extended up to 60 months, which leave these contracts unaudited over a considerable period of time; individual analyses in bidding procedures do not allow to oversee the

---

Copyright©2020 Permission to copy without fee all or part of the material printed in JIDM is granted provided that the copies are not made or distributed for commercial advantage, and that notice is given that copying is by permission of the Sociedade Brasileira de Computação.

market division by the providers; the previous audits are not registered in a specific database. They only appear as *pdf* files; the audit teams are not homogeneous (standardized) in terms of techniques employed and the expertise of their auditors, although such homogeneity would be ideal.

Although some evidences may be checked through individual analysis of bidding procedures, such as providers that sign contracts without offering the best price, checking for long-term patterns may be challenging due to the aforementioned factors. Besides that, evidences such as the assignment of providers to Public Agencies in advance and the market share of each specific provider cannot be checked given a single auction. As a consequence, the use of data mining techniques emerges as an important strategy in the audit work. These techniques make the identification of pattern and collusion detection possible, since the inspection is focused on bidding procedures with the highest probability of fraud.

This paper proposes a methodology for anomaly detection that is able to find fraud and collusion among bidders. We demonstrate the applicability and effectiveness of the proposed methodology by analyzing electronic bidding for the supply of goods and/or services to Public Agencies located in the State of Minas Gerais, Brazil. Through the use of frequent-pattern mining, correlation of multivariate time series, and multi-criteria analyses of the bidding procedures, we nailed scenarios that suggest the incidence of fraud or frustration of the competitive character of the event, possibly through adjustment or agreement among participants.

This paper is organized as follows. This section introduces the subject and related works are presented in the Section 2. Section 3 presents the description of the problem. Section 4 details the methodology. In Sections 5 and 6, the experimental results are analyzed and evaluated. Finally, we present the conclusions and some future work in Section 7.

## 2. RELATED WORK

In this section, we review relevant literature about fraud detection and related techniques. We focus on approaches to tackle this problem, such as anomaly detection with supervised and unsupervised methods, temporal anomaly detection, and neural networks.

Several anomaly detection techniques have been proposed in the temporal context [Gutflaish et al. 2017] [Tian et al. 2019] [Siddiqui et al. 2018] [Song et al. 2018]. Hallac et al. [2017] for example, grouped and segmented multivariate time series to analyze different characteristics of the same object, seeking to find patterns in the data that represent distinct behaviors of a given object. Yagoubi et al. [2018] propose a *Multi-Scale Convolutional Recurrent Encoder-Decoder (MSCRED)* that is capable of detecting anomalies and diagnosing the main causes of multivariate time series. The work proposed by Siddiqui et al. [2018] incorporated analyst feedback to reduce false positives and identify the main anomalies that should be investigated. This feedback is used to adjust the anomaly ranking after each interaction with the user, ranking at the top the most interesting anomalies to the user. Approaches based on neural networks in anomaly detection may be found in [Ren et al. 2019],[Schulze et al. 2019], and [Haldar et al. 2018]. [Ren et al. 2019] propose a time-series anomaly detection service that helps customers to monitor an event stream continuously and alerts for potential incidents promptly. The algorithm is based on Spectral Residual (SR) and Convolutional Neural Network (CNN). In [Haldar et al. 2018] the authors discuss the application of neural networks to the search problem in Airbnb, which presents challenges that are similar to fraud detection. In [Schulze et al. 2019] the authors develop a contextual anomaly detector for testing cars using a neural network and compared to other approaches that use unsupervised methods, such as Isolation Forest and Autoencoder. Ramakrishnan et al. [2019] proposes an anomaly detection framework for Walmart that joins supervised (GaussianNB, Isolation Forest, and Random Forest) and unsupervised methods (Gradient Boosting Machine - GBM e Autoencoder), being able to detect items without prices or wrongly priced.

Several fraud detection techniques have been developed in recent years, such as [Cao et al. 2019],

[Cao et al. 2017], [Zheng et al. 2018],[Nian et al. 2016]. [Cao et al. 2019] introduce the *TitAnt*, a transaction fraud detection system based on feature extraction. The system is capable of predicting an online real-time transaction fraud in milliseconds. Cao et al. [2017] presents HitFraud, an algorithm to compute meta-path based features, capturing the inter-transaction dependency. The algorithm leverages heterogeneous information networks for collective fraud detection by exploring correlated and fast evolving fraudulent behaviors. Yagoubi et al. [2018] used *LSTM-Autoencoder* to develop one-class adversarial nets (OCAN) for fraud detection using only benign users as training data. The *LSTM-Autoencoder* learns the representations of benign users from their sequences of online activities and then detects malicious users by training a discriminator of a complementary GAN model that is different from the regular one. Research in [Nian et al. 2016] proposed an unsupervised ranking method for anomaly detection on auto insurance fraud based on spectral ranking using the first non-principal eigenvector of the Laplacian matrix. In the scope of collusion and fraud detection of bidding procedures, the following works stand out [Ghedini Ralha and Sarmiento Silva 2012; Grilo Junior 2010; Fraga 2017; Balaniuk et al. 2013]. Ghedini Ralha and Sarmiento Silva [2012] analyzed auctions of the Brazilian Federal Government through mining of frequent patterns and multiagent systems. They observed the formation of cartels involving suppliers and corporate bonds among them. Grilo Junior [2010] and Fraga [2017] also used data mining techniques to detect collusion in bidding. Despite identifying collusion in tenders, their research did not explore micro-expressions, as we propose. They basically search for frequent patterns among bidders. Balaniuk et al. [2013] applied a probabilistic classifier to detect fraud in government transactions and assist audit agencies. Despite their analyses being based on multi-criteria audits, considering both suppliers and public agencies, they were used as plain features to the classifier.

This work intends to advance techniques and methodologies for detecting frauds on bidding, not only pointing collusion between providers, but also revealing how they operate. The methodology identifies the fraud and collusion through frequent patterns and the micro-expressions we propose enabled the detection of varied forms of camouflage and concealment of fraudulent behaviour.

### 3. ANOMALY IN BIDDINGS

In this section we describe bids and the problem of anomaly detection in bids, in particular patterns of anti-competitive behavior.

The Public Administration, in order to provide services to society, such as health, education, security, justice administration, and infrastructure, needs to acquire goods and contract services. In Brazil, such acquisitions and contracts have their guidelines defined by art. 37, item XXI of the Brazilian Constitution. This regulation determines that the works, services, purchases and disposals will be contracted through a public bidding process that ensures equal conditions for all competitors, and the principles of legality, impersonality, morality, publicity and efficiency must be observed. There is also infra-constitutional legislation that details the procedures to be performed by the Public Administration in acquisitions and contracts, composed of laws, decrees, ordinances and normative instructions. Therefore, the bidding procedure is the rule for acquisitions and contracts, with direct contracting (without bidding) being the exception, in specific cases defined by law.

There are two main objectives that justify carrying out the bidding procedure. The first is to meet the needs of goods and services by the Public Administration so that it can develop its activities. The second is obtaining the most advantageous conditions with regard to the amount to be paid, in order to maximize the efficiency of public expenditure, whose financial resources come from society through tax collection. For these objectives to be achieved, there must be guaranteed equal opportunity for all those interested in providing goods and services, and to participate in the bidding process. As a result, assuring effective competition among potential providers is a key requirement of the process, since real disputes among bidders tend to reduce the prices to be paid. Although auction is a modality that stimulates competition, through the bids' sharing among competitors, it does not eliminate the

possibility of an occurrence of a fraud, as we discuss next.

The bidding modalities for the acquisition of goods and services, by default Invitation, Price Taking, and Competition, as defined by the Brazilian Law number 8.666/1993 [Licitações 1993], are characterized by a sealed and single proposal submitted in an envelope, with each bidder defining its price proposal without knowing others' proposals, and the winner is the qualified proposal that requires the lowest amount to be disbursed by the Public Administration. It is also part of the process of these modalities to check the documentation required by law (qualification phase) prior to opening each envelope with the price proposal. Just the price proposal envelopes of those bidders approved in the qualification phase are opened. Auction, on the other hand, handles price secrecy differently and is described next.

Auctions have been created by Brazilian Law nº 10.520/2002 [Pregão 2002] and are characterized by an initial proposal from each participant that is also not shared with the other bidders. Later, there is a public session, during which bidders are aware of the bids submitted by other bidders, although they do not know which bidder is responsible for each bid, and each bidder is free to submit as many price proposals (bids) as it deems convenient. The rationale of auctions is that competition is stimulated, since, being aware of the best bid so far at any given time, the bidder may submit a bid with a lower price to win the bid. In the Auction, only the qualification documentation of the winning bidder is verified. Thus, there is an inversion between the qualification phases and price proposals considering other bidding modalities. The Auction can be conducted in person or through electronic media, the latter being mandatory for all entities of the Federal Public Administration, in accordance with Decree nº. 5.450/2005 [Pregão 2005] and Decree nº. 10.024/2019 [Decreto 2019]. In the Electronic Auction, bidders send their bids until the end of the session, determined randomly by the system.

Despite the automation and transparency of the Auction modality, audit still detected frauds. One type of fraud detected is the deliberate withdrawal or declassification of the lowest bidder so that the runner-up (or another runner-up) can win. Fraud occurs when the first ranked bidder bids very low and others give up competing. However, as previously agreed with one or more bidders, the winner gives up or fails to send any document to the acquirer to be purposely disqualified. Thus, the bidders ranked next are called until someone submits the required documentation and is chosen as provider. The electronic form of the auction allows automated data analysis to be carried out and facilitates the detection of anomalies such as the one described, since all auction records, including bids and withdrawals, are stored in a database. For this reason, the auction modality was chosen for the present work.

We define a context as the data associated with an audit scenario, which is characterized by a set of bids determined by the nature of the goods or services being acquired, by the participating entities, by the time interval, by the characteristics of the bids, or combinations of the aforementioned entities. A bid may be seen as a time series of events carried out by one or more entities. Both the bidding, the events, and the participating entities are characterized by attributes that hinge upon the nature of the bidding. The entities are providers and public agencies that acquire products and services. Events are actions, carried out by entities or even groups of entities during the acquisition process.

Anomalous bids differ from regular ones regarding their profiles. An example of an anomaly is the anti-competitive practice, also known as collusion, which is characterized by unusual connection among sets of entities.

The starting point of our task is a set of entities  $E$ , where each entity  $e$  is characterized by a set of attributes  $e_p$ , which include the set of actions that  $e$  may perform or participate. Each action performed by one or more entities is an event  $x$ , characterized by the participating entities  $x_e$ , at the moment in time  $x_t$  and by the attributes that define the event  $x_p$ . Given a context  $\mathcal{C}$ , the universe of events is identified by  $X_{\mathcal{C}}$ . A pattern is a logical expression of predicates and a score that quantifies its relevance. Each predicate, in turn, is a relational expression containing instances of  $x_e$ ,  $x_t$  or  $x_p$ . The

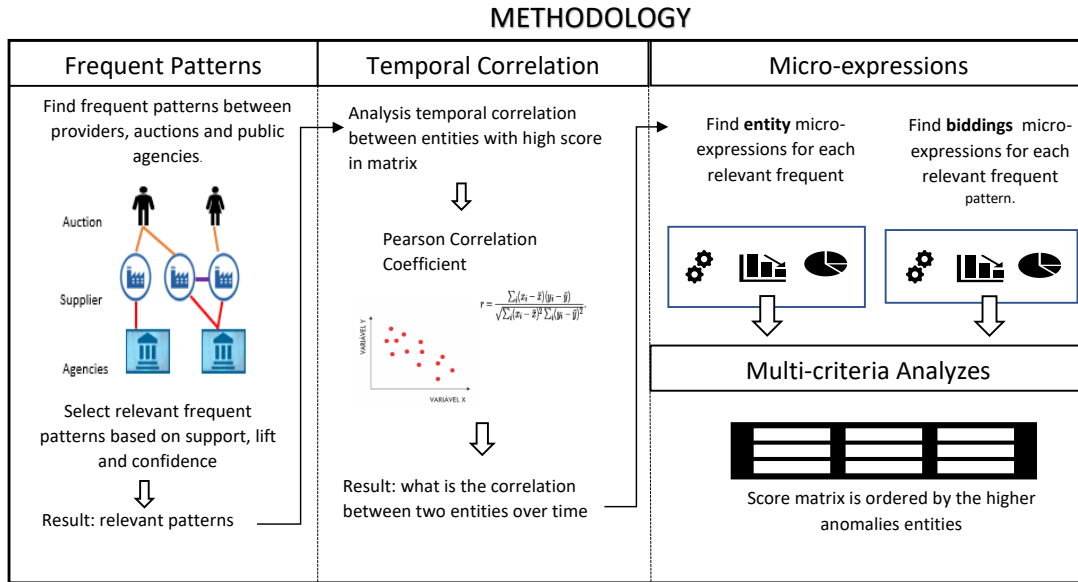


Fig. 1. Our methodology and its four serialized steps. The frequent patterns are used to detect anomalous biddings. Temporal Correlation looks for correlation among entities across time. Micro-expressions materialized audit trails and Multi-criteria Analysis is based on a score matrix of entities.

score may be a combination of predicate scores or a score expression, among other functions. Given an events' set and the domain expressions associated with it, our target is to determine the events' set that maximizes the score of the various possible subsets and satisfied expressions. A brute force strategy to solve the problem is to list all possible events sets as well as the universe of expressions associated with them, to calculate the score for each set and expression, and to rank them according to their score, which is clearly not feasible even for scenarios with few events, since the complexity of the strategy is  $O(2^{|E|} \times X_c^2 \times 2^{|X^c|})$ .

In the next section, we introduce a greedy strategy for determining the desired patterns and their scores.

#### 4. METHODOLOGY

Our methodology is based on four serialized steps. The goal of the first step is to detect frequent patterns. The next step is to identify temporal correlations with significant Pearson correlation coefficient. The third step consists of instantiating micro-expressions, divided into entity and bid micro-expressions, that materialize audit trails. Last, we perform a multi-criteria analysis, based on the score matrix of anomalies entities. These steps are detailed next and depicted in Figure 1.

##### 4.1 Frequent Patterns

In the first step we search for entities that often participate together in the same bidding. Entities  $e$  may be public agencies that organized/placed bidding or providers that bid. According to law 10.520/02 [Pregão 2002], the auction is the main bidding modality for acquisition of goods and services.

It is performed as a public session, through a bidding dispute, ranking the bidders according to their bid value and the qualification phase takes place after the end of the public session, when the least expensive proposal is chosen. Each auction consists of a series of events held by a public agency, until the auctioneer closes the process on due time.

Frequent patterns are used here to detect anomalies in these bids, considering frequent co-occurrence of entities as an anomaly, given a previously determined minimum support (frequency). Notice that, just the co-occurrence of entities in a bidding session should not be considered an anomalous practice. However, such recurring joint participation and other features (formalized here as micro-expressions) may characterize an anomaly known as collusion. Collusion can be formed between providers or between providers and public agencies / auctioneers. Each bidding is considered as a transaction and providers and public agencies as items of that transaction. Thus, frequent sets of providers or providers and public agencies that stood out the most will be selected as anomalous, that is, those that appear repeatedly participating in the same tenders together.

This stage is crucial for collusion identification, because if the entities were identified separately, we probably discover just individual frauds. Here we use anomaly detection to differentiate those providers that stand out. From these frequent sets found, we deepen the analysis of their behavior to ascertain the actual occurrence of a collusion and detail how it is carried out in the steps that follow. The determination of anomalous competing entities is carried out using frequent pattern mining algorithms, among them *Apriori* [Agrawal and Srikant 1994] and *Eclat* [Zaki et al. 1997]. The algorithm input is a file containing the entities under analysis, such as providers, for example. The resulting frequent sets are organized according to the usual interestingness metrics, being *lift* [Zaki and Meira Jr. 2020] the most relevant to our scenario. These sets are the input for the next stage.

## 4.2 Temporal Correlation

After the individual analysis of the entities that compose a frequent pattern, the set as a whole is scrutinized, taking into account the temporal dimension. In other words, we look for the correlation between two entities over time and the influence of each participating entity on the identification process of collusion and/or cartel formation.

A relevant frequent pattern found in the previous step identifies a group of providers with frequent co-participation in an auction. However, the actual participation of providers within the group may vary, some of them may have a higher correlation than others, or even a greater correlation in a given period. This type of analysis aims not only to confirm the pattern previously found, but to decompose how this pattern has occurred across time, and to determine the most interconnected entities within the anomalous pattern. Stronger collusion over a long period or even in a short, but recent, period, should be prioritized by the audit in order to identify the main fraudsters within the collusion, or even to prevent deviations still in progress.

The temporal correlation is assessed through the Pearson Correlation Coefficient. A strong correlation between the frequent set associated with the score defined in the previous step may identify not only how the anomaly elapsed, but also to what extent this behavior persists over time.

We used the input file containing the bidding dates and the number of bids in which each of the anomalous provider participated. This share is calculated between these providers using Pearson's coefficient. Together, given eight providers considered anomalous w.r.t. frequent patterns, for example, their pairwise correlations may not be the same, demonstrating that four of them are more interconnected, or even that part of these providers present a strong correlation (Pearson coefficient between 0.7 and 0.9) <sup>1</sup> over a certain period, but this correlation may decline later. This can be explained,

<sup>1</sup>Pearson's coefficient above 0.9 presents very strong correlation and presents strong correlation when between 0.7 and 0.9. Between 0.5 and 0.7 presents moderate correlation and Pearson's coefficient below 0.5 presents no correlation.

in some cases, by the fact that the two providers with a very strong correlation are dominating and dividing the market between them, and, therefore, present greater correlation, but other providers, despite being part of the collusion, play a less relevant role, like participating just to create false competition.

### 4.3 Micro-expressions

After defining the sets of interest entities, based on co-occurrences and temporal correlations, we instantiate micro-expressions for all sets. These micro-expressions materialize audit trails, which can be determined based on past experience or legislation. Micro-expressions vary over time among different groups, according to the context or entity types. Therefore, for a better division among the different groups or entities, it is necessary to identify whether the context refers to a work, a service provision, purchases of which types of objects or any other type of labor in a given market. The greater the number of satisfied micro-expressions (if binary) or their magnitude for the selected sets of entities (if continuous), the more attention has to be paid to the detected potential anomaly. Micro-expressions can be associated with entities or bids as described here after.

#### 4.3.1 Entity Micro-expressions.

The entity micro-expressions are related to the set of biddings by an entity over the analyzed period. In this work, we employed the following micro-expressions:

- $\mathcal{E}_{noprim}$ : The total percentage of signed contracts in which the chosen entity was not the winner or first ranked. The bidding consists of a sequence of bids submitted by the providers and, at the end of the bidding period, the provider with the best price is ranked first. In this case, the winning provider is expected to resume the hiring process. Otherwise, the next provider in the ranking is summoned to do so. The objective of the bidding is to purchase for the best price, but when the original winner is not chosen, there is a financial loss, at least.
- $\mathcal{E}_{motdesc}$ : Distribution of the most frequent reasons for disqualification. Considering repeated withdrawals or disqualifications, the minutes of the sessions were investigated in order to discover the main reasons for disqualification. This verification is done manually because that information is not in the database of the *SIASG* system, but in a document file, available at the *Comprasnet* platform.
- $\mathcal{E}_{market}$ : Market share of a provider. It analyzes the distribution of the companies that participate in the context being analyzed. More specifically,  $C$  is the context or audit scenario and the distribution shows all ‘provider’ type entities that are part of context  $C$ , in particular  $v$ , which is the contract value per item. The percentage of each existing provider in the context will be represented by the total value of all the approved contracts to this provider divided by the total value of all approved contracts of the providers from this scenario. The *SIASG* database contains several kinds of biddings, including goods purchases and service provision. However, there is no market concept or perspective in the database that supports the assessment of the market share. It contains only the description of the bid item or the field object of the bid. These fields are solely textual and have no mandatory parameter. To find the market in a context, entities that participated in any bidding, in which at least one of the correlated providers participated or whose keywords characterize the market, were selected. This analysis does not consider the entire market in a context. There may be bids in which some of the related providers did not participate in and which do not have the keywords highlighted in the description of their item. We assume that it covers most of the market.
- $\mathcal{E}_{agreement}$ : Sets of entities frequently associated with the public agency. The participation of providers in bids from given public agencies possibly indicates a market division among providers. This micro-expression highlights a presumed arrangement among correlated providers. The public agencies were grouped by providers, characterized by the number of contracts. The public agencies are analyzed in descending order of cartel occurrence likelihood.

- $\mathcal{E}_{discount}$ : The percentage of total discounts in the context. This number is based on the difference between the summation of the estimated value and the final agreed value on signing the contract.
- $\mathcal{E}_{PositionAverage}$ : It is the average of all ranks in which the provider was classified. Zero is assigned to each occurrence an entity wins as the first ranked or the rank average otherwise.
- $\mathcal{E}_{successRate}$ : The success rate represents the number of successful biddings divided by all biddings in which each entity participated.
- $\mathcal{E}_{materiality}$ : Categorization of contracts according to value ranges.
- $\mathcal{E}_{Amendment}$ : The percentage of total amendments in the contracts.

4.3.2 *Bidding Micro-expressions*. The bidding micro-expressions are instantiated considering the set of biddings in relation to the sets of competing entities that are temporally correlated. That is, those biddings where all entities of a given context participate in the set. In this work we considered the following bidding micro-expressions:

- $\mathcal{L}_{Dropout}$ : Percentage of dropout providers. Dropout providers are those who offer the first bid above the estimated value, as defined by the bidding call, and are, therefore, immediately disqualified.
- $\mathcal{L}_{disqualified}$ : Percentage of disqualified providers. Disqualified providers are those that win the bidding but, for some reason, do not resume the hiring, leaving room to the next provider in the ranking order.
- $\mathcal{L}_{perfildd}$ : The ratio between the dropout percentage or provider abstention in the selected bids within the context. The providers that resigned most frequently are examined w.r.t how they perform on the market. Some of them do not hold a history of successfully approved contracts when they sign a contract. Or, in other cases, they do not have any contract signed with the government, participating in the biddings and never winning. This fact may indicate a possible agreement among providers to disguise a competition, and actually represent a collusion. This is quantified by ratio between the number of dropouts/disqualified and the number of signed contracts.
- $\mathcal{L}_{Overpricing}$ : Percentage of overpricing per entity. The total percentage of overpricing is calculated by the difference between the total sum of the values offered by the first ranked (winning bidder) and the total sum of the values of the signed contracts, in auctions where the winning bidder is not hired.
- $\mathcal{L}_{Waiver}$ : The proportion of bid waivers. The proportion of bid waivers within the context is the ratio between the number of waivers and the total number of contracts.
- $\mathcal{L}_{1Participants}$ : Individual competition. The percentage of biddings in which a single provider participates or a non-competitive bidding occurs.
- $\mathcal{L}_{FewParticipants}$ : Weak competition. The percentage of biddings in which only few providers participate; no more than 5 per bidding.

#### 4.4 Multi-criteria Analyses

The result of each micro-expression will compose the entity score index of the frequent pattern set. Micro-expressions are not applicable to all biddings and some of them need to satisfy specific requirements to justify their instantiation. The combination of all, or part of them, provides more strength towards the identification of collusion. Each occurrence of the micro-expressions will add the respective score value to the entity score index. The micro-expression overpricing, for example, identifies whether a provider overpriced or not the service or good in its bid. If so, the percentage of overprice will be converted into a scale ranging from 0 to 1, where 1 represents maximum values of overprice and 0 minimum values. In other words, a score of 0.8 presented by a provider in this micro-expression means the practice of overpricing in 80% of the bids he or she submitted. The sum of the scores of all satisfied micro-expressions of the provider will compose the final score, determining the risk index of this provider. Providers with a score of zero or close to zero indicate low or no risk. The higher



this index is, the greater the risk associated with that provider and, thus, the greater the likelihood of fraud. Therefore, the micro-expressions, analyzed individually, serve to describe the evidence of fraud of a given provider, but, together with the frequent patterns and correlation, they detail how the collusion works. We can consider a set of providers regarded as anomalous and the indexes of micro-expressions  $\mathcal{E}_{agreement}$  and  $\mathcal{E}_{noprim}$  being close to 1. These anomalous providers are probably part of a collusion where the slice of the public agencies in the context occurs, materialized through the successive disqualifications of the first ranked provider.

The frequent patterns indicate the group of providers that operate in the collusion. The correlation identifies which ones stand out in this anomalous set, as well as the strength of their interconnection through time. The micro-expressions give transparency to the individual performance of the providers, demonstrating how the participants in the collusion really slice the market. The multi-criteria analysis defines the score of collusion entities, determining who should be prioritized for sake of auditing. The combination of these stages supports the creation of an automated machine learning model for fraud detection.

## 5. CASE STUDY

Our methodology was applied in a real case study, using data from biddings issued by federal agencies in the state of Minas Gerais, from January 2013 to October 2019, stored at the *SIASG* (database with information from Comprasnet<sup>2</sup>).

In particular, we considered 259.345 bidding items and 7.652 providers. These items refer to the bids made for contracting services, in the period analyzed. Each bid in our dataset consists of the following information: bid item, public agency, bid amount, date of purchase, auctioneers, final contract value, estimated price by the public agency, reasons for dropout, description of the object and all providers present in the auction. The result of applying our methodology to this context detected five providers with frequent significant participation, that is, always acting together in the same events in the same public agencies, indicating a possible collusion with restricted market share among them. Next we detail the findings at each stage of the methodology that led us to this conclusion.

### 5.1 Frequent Patterns

In order to determine the frequent competitors, the *Apriori* algorithm was executed on 259.345 transactions (all bidding items from the context), using 0.001 for support and 0.6 for confidence. As mentioned, the input data contain purchases from public agencies that carried out biddings for outsourced services in the period from 2013 to 2019, as mentioned. Since we are looking for collusion patterns that are still operative or can be punished, frauds prior to 2013 may have been prescribed, and therefore it may be impossible to charge anyone.

In this stage, the bid items, the providers participating in each bid, as well as the responsible public agency and the auctioneer were used as identifier. To select the anomalous sets, *lift* values were taken into consideration. The major anomalies detected were between providers or between provider and public agency. The auctioneers in charge of these auctions were also analyzed but, at first, no significant relationships were found between auctioneers and providers, what may be explained by the fact that the biddings were held by different entities across the considered period. Support values were low due to the dispersion of bids in the database, since it comprises bids from all public agencies in the state of Minas Gerais. Therefore, high support values would probably rule out relevant sets, considering that providers usually vary from one context to another.

<sup>2</sup>Comprasnet is a website that provides information about Federal Government purchases - <https://www.comprasgovernamentais.gov.br/>.

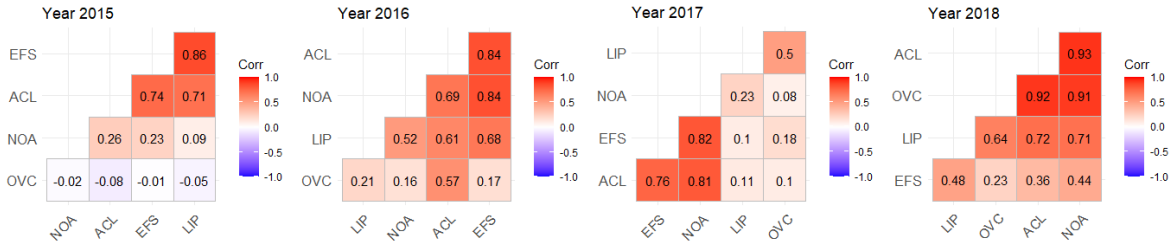


Fig. 2. The variation of the temporal correlation over the years. Not all five providers showed the same correlation across the verified period.

Among the anomalous sets identified in this stage for contracting services, some attracted attention and were further explored in depth. The case study exposed here stood out, identifying five providers strongly correlated w.r.t. participation, indicating probable collusion and a significant final score indicative of fraud, as detailed next.

## 5.2 Temporal Correlation

We generated the temporal correlation per year, based on the purchase date of the item biddings and the bid frequency of each entity per item on the same dates. The temporal correlation was calculated using the multivariate Pearson Correlation Coefficient. If the entity competed for 10 items in the same bidding process, it will have 10 participations on that date. We adopted such analysis granularity in order to not discard information from entities that compete for a different number of items in the same bidding process.

In this step, we analyzed the correlation between the five providers identified by the frequent-pattern stage. However, not all five showed the same correlation across the verified period. As we can see in graph 2, this correlation varied over the years. ACL has a strong correlation (above 0.7) with EFS from 2015 to 2017. However, in 2018 this correlation decreases to 0.36 (considered a weak correlation). NOA and ACL, despite no significant correlation in 2015 (0.26), present a strong correlation from 2016 to 2017, above 0.8, and very strong correlation in 2018, with Pearson's coefficient of 0.93. LIP just present a strong correlation with ALC in 2015 and 2018.

Thus, through the temporal correlation we were able to verify, for those five anomalous providers verified in the previous stage, that not all showed the same correlation throughout the analyzed period. This relationship may indicate a dominance of some providers in the collusion or even an alternation among them across time.

## 5.3 Entity Micro-expressions

Next, we instantiate the micro-expressions related to the most relevant provider – provider and provider - public agency relationships. The information present in the bids of the five providers highlighted as anomalous are analyzed here using micro-expressions. Thus, the bidding information was used considering the providers, the bids' items in which they participated, contract values, providers rank, dropout provides, public agencies, auctioneers, date of purchase and estimated contract values. The results of the micro-expressions considered relevant for this anomalous set are shown below. Here we considered the period from 2008 to 2019 because we are trying to assess how the entities behave over a longer time period.

$\mathcal{E}_{noprim}$ : The result of the micro-expression  $\mathcal{E}_{noprim}$  is shown in Table I. The fields  $\#l$  e  $\#i$  are the number of biddings where entity  $e$  participated and the number of items that  $e$  competed,

Table I. Approved Values and Quantities X Approved in Higher Ranks 2008 to 2019 - Brazil

Ent.	# $l$	# $i$	# $i_{prim}$	# $l_{noprim}$	% $l_{noprim}$	\$Homolog	\$Homolog#1 <sup>o</sup>
ACL	10015	32135	1548	1244	80%	R\$500.960.788,76	R\$442.193.399,55
NOA	1071	3219	825	546	66%	R\$249.792.951,15	R\$179.298.059,90
OVC	793	4310	420	245	58%	R\$235.020.112,88	R\$145.002.096,90
EFS	846	2525	143	108	76%	R\$89.848.858,41	R\$45.137.363,41
LIP	746	2214	284	155	55%	R\$39.310.272,56	R\$26.191.665,02

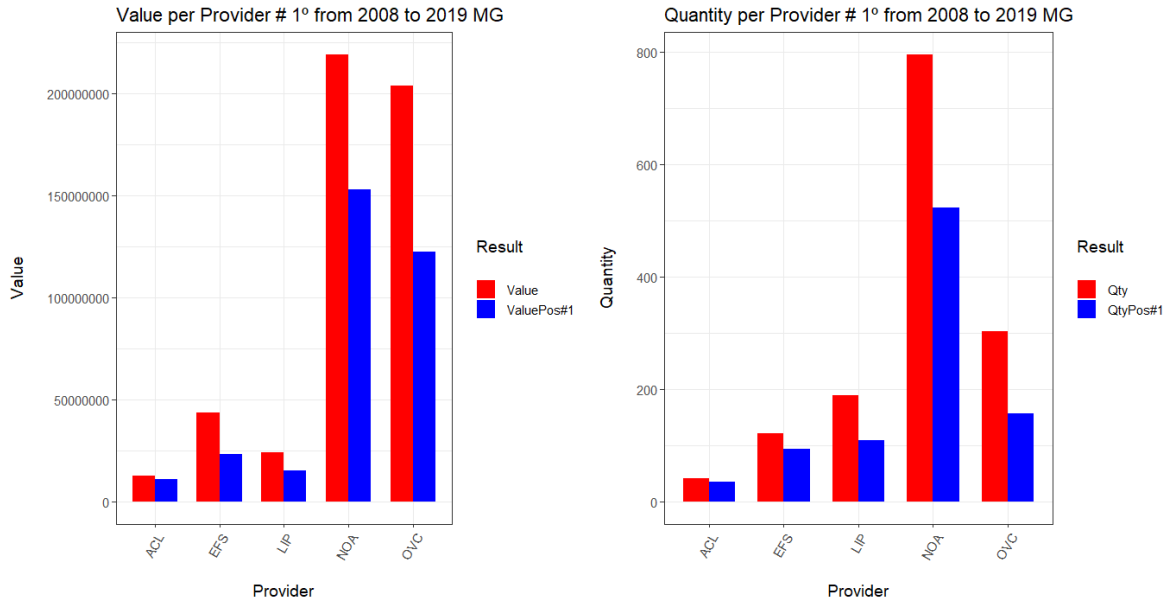


Fig. 3. Contracts signed in higher auction ranks - by value and by quantity. The red bars indicate the number of contracts signed by the provider in any rank and the blue bars indicate the total of contracts that it signed without being the first ranked.

respectively. The field  $\#i_{prim}$  indicates the number of items for which  $e$  won classified in first place. The field  $\#l_{noprim}$  is the number of items in which  $e$  signed a contract but was not ranked first.  $\%l_{noprim}$  is, among all biddings where the provider was hired, the percentage of biddings where it was not ranked first. Finally,  $\$Homolog$  and  $\$Homolog\#1^o$  are the approved cumulative value in contracts that the provider celebrated and the approved cumulative value in the contracts that the provider celebrated but was not ranked first, respectively. As we can see, among all auctions where one of the five entities won, on average, they were not ranked first in more than 50% of them. In some cases, this ratio reached 80%. It means that another provider won the bidding but was disqualified or dropped out, and one of the five providers of the collusion signed the contract. Figure 3 shows the percentage of approved contracts when the provider did not win the event ranked first. These ratios are represented as quantities and contract values. The red bars indicate the number of contracts signed by the provider in any rank and the blue bars indicate the total of contracts that it signed without being the first ranked. As shown in table I and figure 3, the five providers, although often do not submit the best proposal in terms of value and do not win the auction, at the end sign the contracts with public agencies.

$\mathcal{E}_{motdesc}$ : Concerning the micro-expression  $\mathcal{E}_{motdesc}$ , the most frequent reasons for disqualification were entities that do not present the necessary documentation; entities disqualified by the auctioneer; unaffordable price; cost spreadsheet problems; entities disqualified by their own request. The reasons for withdrawal presented would be considered normal when referring to one bid or even a small number of bids. However, when analyzing this micro-expression together with the previous one ( $\mathcal{E}_{noprim}$ ), it is

Table II. Overprice in Contracts Signed by the Five Entities Analyzed

Entity	Overprice
NOA	R\$ 26.960.195,00
LIP	R\$ 4.112.675,00
ACL	R\$ 3.577.980,00
EFS	R\$ 2.099.826,00
OVC	R\$ 304.336,00

clear that these five providers exhibited high rates of signing contracts without being the first ranked, followed by disqualifications from the first ranked for reasons that indicate, at least, a non-compliance with the entire bidding process. Therefore, the two micro-expressions analyzed together indicate the likely agreement between the winning providers and the disqualified providers.

$\mathcal{E}_{market}$ : For market micro-expression  $\mathcal{E}_{market}$ , we found that entities NOA, ACL and OVC present the highest contract values between 2008 and 2019, being this value up to four times the value of other similar entities from the same market. That is, they seem to dominate the market compared to other providers in the same activity sector. Providers ACL, LIP and EFS are within the first 13 largest providers in the same market w.r.t. signed contracts.

$\mathcal{E}_{agreement}$ : Another analysis of micro-expressions that provided fraud evidence was the participation of providers in the public agencies selected according to the micro-expression  $\mathcal{E}_{agreement}$ . The collusion entities seem to be providing services and goods to the same public agencies for the last ten years, indicating a possible market share between them, according to Figure 4. This graph shows the main public agencies that contracted the same providers. The colored spheres represent the five providers and the size of each sphere indicates the monetary value of signed contracts. The X axis represents the years in which the providers had the contract approved<sup>3</sup> and the Y axis represents the main public agencies where the five providers operate. This chart does not display other entities that are not part of the group of five analyzed and that submitted bids to auctions from any of these public agencies because the objective was to show the existence or not of a prior division among the five entities for being hired by public agencies. When analyzing in detail the public agencies separately and comparing the contracts of the five providers to contracts of other providers with public agencies, we noticed that, in some cases, such as at *U109*, the provider NOA dominates the outsourcing service contracts<sup>4</sup>. Considering the contracts from public agency *A101*, the provider OVC signed large contracts, providers NOA and EFS signed some small contracts and the rest of the contracts is dispersed among other outsourcing providers.

#### 5.4 Bidding Micro-expression

In this step we evaluated all biddings where the five providers participated and identified possible distortions.

$\mathcal{L}_{dd}$  e  $\mathcal{L}_{perfildd}$ : Regarding the micro-expressions  $\mathcal{L}_{dd}$  and  $\mathcal{L}_{perfildd}$ , when analyzing providers' dropouts, it is clear that some of them did not participate in bidding events where the cartel entities did not participate either. In addition, several of these entities did not win any bids whenever the cartel entities participated, but won when the cartel entities did not participate. In some cases, these providers have participated in more than 100 auctions and have not won or signed any contracts.

$\mathcal{L}_{overprice}$ : In the micro-expression  $\mathcal{L}_{overprice}$ , all five entities of the analyzed cartel practiced overpricing. The entity NOA presented more than 20 million Reais in overpricing according the defined rules. This comparison is showed in Table II.

<sup>3</sup>The year was extracted based on the reference date field of the database purchase *SIASG*.

<sup>4</sup>Here we selected the bids to the public agencies where the five providers dominate. In addition, bids were filtered by the words *cleaning, maintenance, maintenance and service* in the description field of the bidding item.

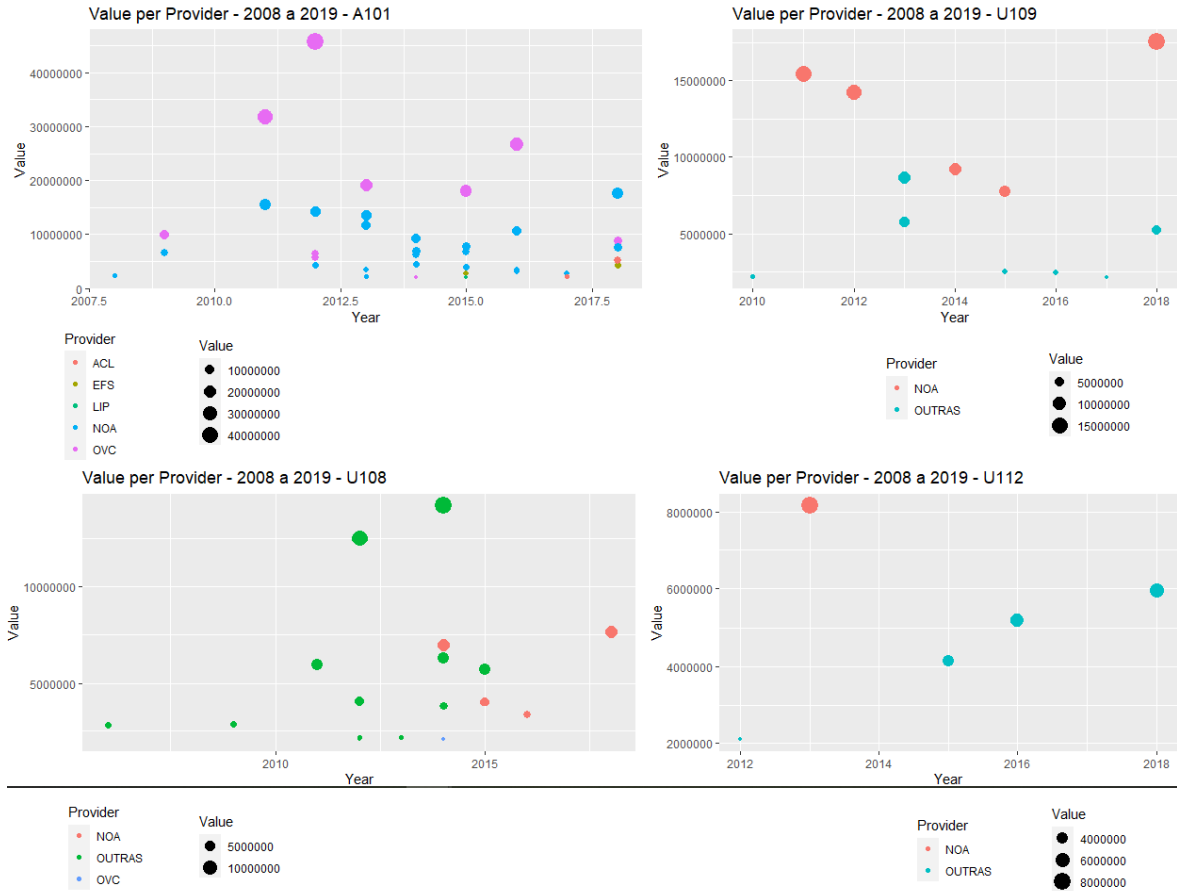


Fig. 4. Providers x public agencies select by agencies. The main public agencies that contracted the same providers.

### 5.5 Multi-criteria Analyzes

The instantiation of all aforementioned micro-expressions is shown in Table III, which presents their score for each of the five providers identified by our methodology with strong correlation. The values were scaled accordingly to ease the understanding. The sum of all scores results in the entity final score. Since it is possible to verify, all of them are associated with scores well above zero, and some even had an index close to, or even above, 0.5 in the micro-expression  $\mathcal{E}_{noprim}$ . Their materiality is also high, especially for NOA and OVC.

Table III. Providers Frequent Patterns- Service Outsourcing

Micro-expressions - Service Outsourcing								
Provider	$\mathcal{E}_{Suc.}$	$\mathcal{E}_{noprim}$	$\mathcal{E}_{Disp.}$	$\mathcal{E}_{Desc}$	$\mathcal{E}_{Amend.}$	$\mathcal{E}_{AveRate.}$	$\mathcal{E}_{Material.}$	FinalScore
NOA	0,2	0,66	0,2	0,167	0,5	0,5	0,45	2,67
LIP	0,1	0,57	0,05	0,66	0,5	0,5	0,25	2,63
ACL	0,1	0,88	0,01	0,25	0,5	0,5	0,25	2,49
EFS	0,1	0,78	0,01	0,14	0,5	0,5	0,3	2,33
OVC	0,1	0,52	0,038	0,03	0,2	0,2	0,4	1,48

Through micro-expressions it is possible to understand how the collusion operates. Market dominance and contract slicing, probably signed through an agreement among the five providers and using other providers (dropouts) to establish this procedure. The dropout providers serve to disguise

the competition and to raise the amounts paid by the public administration. This is because these providers won the bid for a price and, when they gave up, the public administration signs a contract with the next provider at a higher price. In addition, the five providers appear to act in some bids only to mask a false competition, thus presenting a normal success rate compared to the other providers, but this rate is much higher than normal when considering the public agencies in which each one dominates.

Therefore, through frequent patterns, entity micro-expressions and temporal correlation, the methodology detected five providers with frequent relevant participation, that is, always acting together in the same events, extending their contracts in the same public agencies, indicating a possible collusion between outsourcing providers. Also, through the methodology, it was possible to confirm that this collusion, a priori, is an agreement between the providers themselves, without any interference from the public agencies, since no evidence was found indicating correlations between providers and auctioneers. In summary, the application of the proposed methodology to a real case allowed the detection of a collusion in an efficient and assertive manner. Further details on the analyses were omitted due to data confidentiality.

## 6. EVALUATION

The results generated by the methodology were compared to actual audits performed under the mentioned scenarios. It was shown that some audits and inspections identified providers and public agencies that failed in obeying the rules and laws and were properly identified by the methodology. Next, we present some results of the methodology application confronted with the audits and inspections carried out.

### 6.1 Frequent Patterns

Providers were identified through the methodology as frequent participants, with possible market sharing and division of the public agencies with which they interact. They were subject to inspection and some of them are liable for prosecution in other states for putting together a cartel and frustrating the competitive nature of the event, getting to the point of even having their assets frozen.

### 6.2 Temporal Correlation

During the audit, some sets of providers were identified over the years, but representing weak evidences. On the other hand, the temporal correlation of the methodology was able to detail how this participation occurred over the years and which providers were more correlated in a given period. Such analysis is particularly effective for non symmetric correlations throughout the period, that is, those that have varied among providers.

### 6.3 Micro-expressions

Regarding the micro-expressions, some of them were detected in both, audit and methodology. Amongst them,  $\mathcal{E}_{overprice}$  and  $\mathcal{E}_{ScoreDeclassification}$  were the ones that stood out the most.

$\mathcal{E}_{overprice}$  Within the methodology, a potential overprice was indicated in the contracts of the providers analyzed, since their hirings were at a value above the market practice, generating losses for the public administration. During inspections performed in some agencies, it was proven that providers were actually charging prices above market value.

$\mathcal{E}_{ScoreDeclassification}$  The providers that are part of the potential collusion indicated by the methodology simulate and put on a false competition where other providers enter the contest just to disguise the process, then withdraw or are disqualified, as pointed out by the micro-expressions of entities.

Audits and inspections carried out by the CGU, in the public agencies where these providers operate, pointed out that, in processes where contracts had been formalized, the highest ranked providers gave up on the event or did not present the necessary documents to complete the process.

Another micro-expressions, such as  $\mathcal{E}_{agreement}$ ,  $\mathcal{E}_{motdesc}$ ,  $\mathcal{L}_{perfield}$ , and  $\mathcal{L}_{dd}$ , provided insights that were not detected in previous audits.

## 7. CONCLUSIONS

This paper presents an audit methodology for collusion detection through the use of frequent pattern mining, temporal correlation and joint multi-criteria analysis. The results demonstrate that our proposal is able to identify not only possible collusion, but also to explain how they work. We present a case study where we identified five providers that present high correlation and suspicious behavior. In, at least 50% of the contracts that they signed, they were not the winners ranking first. Regarding their provision relationship to public agencies, they seem to follow a pattern, extending their contracts always with the same public agencies. Beyond that, the number of dropout providers that participated in the biddings and never won, whenever the same five providers participate together is expressive. Last, the most common reasons for disqualification do not vary much and the five providers signed a huge amount of contracts and extensions, considering the local market in the state of Minas Gerais.

The results generated by the methodology were compared to actual audits and inspections performed by CGU office in Minas Gerais. It was shown that some audits and inspections revealed that those five anomalous providers presented the same irregularities and the collusion pointed out by the methodology. It is worth mentioning that the methodology was not based on the audit findings, nor the audits and inspections were based on the methodology. Both were developed separately, without influencing each other. However, the results of audits and inspections are in line with the results of the methodology, demonstrating its effectiveness. The methodology is now being used to direct new work, thereby reducing the scope yet to be audited.

For future work, we intend to automate the derivation of micro-expressions and their process of combined analyses. We intend to compare the micro-expressions generated until now with the main micro-expressions from the previous audit that pointed out fraud detection in government purchases. These and other improvements will also be analyzed in other government acquisition scenarios.

## Acknowledgments

This work was partially supported by FAPEMIG, CNPq and CAPES, by the projects InWeb, MASWeb, EUBra-BIGSEA, INCT-Cyber, ATMOSPHERE and by Brazilian CGU.

## REFERENCES

- AGRAWAL, R. AND SRIKANT, R. Fast algorithms for mining association rules in large databases. In *Proceedings of the 20th International Conference on Very Large Data Bases*. VLDB '94. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 487–499, 1994.
- BALANIUK, R., BESSIERE, P., MAZER, E., AND COBBE, P. Corruption risk analysis using semi-supervised naïve bayes classifiers. *International Journal of Reasoning-based Intelligent Systems* vol. 5, pp. 237 – 245, 2013.
- CAO, B., MAO, M., VIIDU, S., AND YU, P. S. Hitfraud: A broad learning approach for collective fraud detection in heterogeneous information networks. *CoRR* vol. abs/1709.04129, 2017.
- CAO, S., YANG, X., CHEN, C., ZHOU, J., LI, X., AND QI, Y. A. Titant: Online real-time transaction fraud detection in ant financial. *CoRR* vol. abs/1906.07407, 2019.
- DECRETO. Decreto n.10.024/2019 - pregão, 2019. Acessado: 09 mar. 2020.
- FRAGA, A. *Deteção de Casos Suspeitos de Fraude em Licitações Realizadas no Município da Paraíba*. Ph.D. thesis, Universidade Federal da Paraíba, Brasil, 2017.
- GHEDINI RALHA, C. AND SARMENTO SILVA, C. V. A multi-agent data mining system for cartel detection in brazilian government procurement. *Expert Syst. Appl.* 39 (14): 11642–11656, Oct., 2012.

- GRILO JUNIOR, T. *Aplicação de Técnicas de Data Mining para Auxiliar o Processo de Fiscalização*. Ph.D. thesis, Universidade Federal da Paraíba, 2010.
- GUTFLAISH, E., KONTOROVICH, A., SABATO, S., BILLER, O., AND SOFER, O. Temporal anomaly detection: calibrating the surprise. *CoRR* vol. abs/1705.10085, pp. 1705.10085, 2017.
- HALDAR, M., ABDOOL, M., RAMANATHAN, P., XU, T., YANG, S., DUAN, H., ZHANG, Q., BARROW-WILLIAMS, N., TURNBULL, B. C., COLLINS, B. M., AND LEGRAND, T. Applying deep learning to airbnb search. *CoRR* vol. abs/1810.09591, 2018.
- HALLAC, D., VARE, S., BOYD, S. P., AND LESKOVEC, J. Toeplitz inverse covariance-based clustering of multivariate time series data. *CoRR* vol. abs/1706.03161, pp. 1706.03161, 2017.
- LICITAÇÕES, L. lei n.8666/93 - licitações e contratos, 1993. Acessado: 09 jun. 2019.
- NIAN, K., ZHANG, H., TAYAL, A., COLEMAN, T., AND LI, Y. Auto insurance fraud detection using unsupervised spectral ranking for anomaly. *The Journal of Finance and Data Science* 2 (1): 58 – 75, 2016.
- PREGÃO, D. Decreto n.5450/2005 - pregão, 2005. Acessado: 09 mar. 2020.
- PREGÃO, L. lei n.10.520/2002 - pregão, 2002. Acessado: 09 jun. 2019.
- RAMAKRISHNAN, J., SHAABANI, E., LI, C., AND SUSTIK, M. A. Anomaly detection for an e-commerce pricing system, 2019.
- REN, H., XU, B., WANG, Y., YI, C., HUANG, C., KOU, X., XING, T., YANG, M., TONG, J., AND ZHANG, Q. Time-series anomaly detection service at microsoft. *CoRR* vol. abs/1906.03821, 2019.
- SCHULZE, J.-P., MROWCA, A., REN, E., LOELIGER, H.-A., AND BÖTTINGER, K. Context by proxy: Identifying contextual anomalies using an output proxy. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. KDD '19. ACM, New York, NY, USA, pp. 2059–2068, 2019.
- SIDDIQUI, M. A., FERN, A., DIETTERICH, T. G., WRIGHT, R., THERIAULT, A., AND ARCHER, D. W. Feedback-guided anomaly discovery via online optimization. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. KDD '18. ACM, New York, NY, USA, pp. 2200–2209, 2018.
- SONG, D., XIA, N., CHENG, W., CHEN, H., AND TAO, D. Deep  $r$ -th root of rank supervised joint binary embedding for multivariate time series retrieval. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. KDD '18. ACM, New York, NY, USA, pp. 2229–2238, 2018.
- TIAN, K., ZHOU, S., FAN, J., AND GUAN, J. Learning competitive and discriminative reconstructions for anomaly detection. *CoRR* vol. abs/1903.07058, pp. 1903.07058, 2019.
- YAGOUBI, D. E., AKBARINIA, R., KOLEV, B., LEVCHENKO, O., MASSEGLIA, F., VALDURIEZ, P., AND SHASHA, D. Parcorr: efficient parallel methods to identify similar time series pairs across sliding windows. *Data Mining and Knowledge Discovery* vol. 32, 08, 2018.
- ZAKI, M. J. AND MEIRA JR., W. *Data Mining and Machine Learning: Fundamental Concepts and Algorithms*. Cambridge University Press, 2020.
- ZAKI, M. J., PARTHASARATHY, S., AND LI, W. A localized algorithm for parallel association mining. In *Proceedings of the Ninth Annual ACM Symposium on Parallel Algorithms and Architectures*. SPAA '97. Association for Computing Machinery, New York, NY, USA, pp. 321–330, 1997.
- ZHENG, P., YUAN, S., WU, X., LI, J., AND LU, A. One-class adversarial nets for fraud detection, 2018.