

---

## **Sensitive data hiding in financial anti-fraud process**

---

**Vassilios S. Verykios  
and Elias C. Stavropoulos**

School of Science and Technology,  
Hellenic Open University,  
18 Par. Aristotelous str,  
26335, Patras, Greece  
ORCID: <https://orcid.org/0000-0002-9758-0819>  
Email: [verykios@eap.gr](mailto:verykios@eap.gr)  
Email: [estavrop@eap.gr](mailto:estavrop@eap.gr)

**Vasilis Zorkadis**

Hellenic Data Protection Authority,  
Kifissias 1-3, 115 23, Athens, Greece  
ORCID: <https://orcid.org/0000-0003-4399-2495>  
Email: [zorkadis@dpa.gr](mailto:zorkadis@dpa.gr)

**George Katsikatsos  
and Evangelos Sakkopoulos**

Department of Informatics,  
University of Piraeus,  
Karaoli and Dimitriou 80, Piraeus, Greece  
ORCID: <https://orcid.org/0000-0002-6852-384X>  
Email: [gkatsik@unipi.gr](mailto:gkatsik@unipi.gr)  
Email: [sakkopul@unipi.gr](mailto:sakkopul@unipi.gr)

**Abstract:** This research presents an approach to protect personally identifiable information in compliance with the national and European institutional data protection framework in a way that still allows interoperability of information systems and applications. It is proposed to adopt privacy-preserving information hiding techniques to facilitate targeted data mining without infringing privacy restrictions. This approach is proposed as a strategic tool in the fight against financial and insurance fraud. To resolve issues related to the implementation of the protected registration interface process, the research team is turning attention to the development of algorithms and approaches based on intelligent itemset hiding. The research proposal attempts to contribute to the strategic modernisation of public authorities and financial organisations, aiming at the production of original software to provide services to them, facilitating and accelerating the work to combat fraud. The approach is analytically prevailing on previous approaches and it has experimentally shown encouraging results.

**Keywords:** data hiding; privacy preserving data mining; knowledge hiding; frequent itemset hiding; constraint-based data mining; anti-fraud; sensitive data.

**Reference** to this paper should be made as follows: Verykios, V.S., Stavropoulos, E.C., Zorkadis, V., Katsikatsos, G. and Sakkopoulos, E. (2022) 'Sensitive data hiding in financial anti-fraud process', *Int. J. Electronic Governance*, Vol. 14, Nos. 1/2, pp.7–27.

**Biographical notes:** Vassilios S. Verykios received his Diploma in 1992 from the Computer Engineering and Informatics Department at University of Patras, in Greece, and his MSc and PhD degree from Purdue University in USA, in 1997 and 1999, respectively. He is currently a Professor and the Director of the Postgraduate Programs on Information Systems, and Data Science and Machine Learning and the Founder and Director of the Big Data Analytics and Anonymization Lab, at the School of Science and Technology in the Hellenic Open University. His main research interests include data management and privacy, and his work has been internationally recognised.

Elias C. Stavropoulos holds a BSc in Mathematics and a PhD in Computer Science. Currently, he is the head officer of the IT Department and a Tutor-Counselor of HOU. He has 20 years of experience in open and distance education. He has participated in several EU and national projects as a senior researcher, developer, and course designer and creator. He has co-authored 42 scientific papers with 500+ citations in Google Scholar, and two books for Data Science using R (in Greek). His research includes educational data mining, learning analytics, association rule hiding, privacy and anonymity, and e-learning technologies and methodologies.

Vasilis Zorkadis received his Engineering Diploma from the Aristotle University of Thessaloniki, Greece and his PhD from the Faculty of Informatics, University of Karlsruhe, Germany. Since 2004, he is the Director of the Secretariats' Directorate of the Hellenic Data Protection Authority. In the past, he was visiting Professor with the Department of Informatics, University of Ioannina (1996–1999), with the Department of Education of Technology and Digital Systems, University of Piraeus (2001–2004) and with the Hellenic Open University (2000–2019), Greece. He has written two books and has published over 50 papers on the subjects of security and privacy in international journals and conference proceedings.

George Katsikatsos is a PhD candidate at the Department of Informatics, University of Piraeus, Greece. He is a Diploma holder on Electrical and Computer Engineering from the School of Electrical and Computer Engineering of National Technical University of Athens. He holds a masters diploma on History and Philosophy of Sciences and Technology from the joint postgraduate program of National University of Athens and National Technical University of Athens. He has also PhD studies in Philosophy at the University of Illinois at Chicago USA. He has been trained at the National Center of Administration, Greece and he has long professional experience in Special Secretariat for Financial and Economic Crime Unit, Athens, Greece.

Evangelos D. Sakkopoulos is an Assistant Professor at the Department of Informatics, University of Piraeus, Greece. He holds his Diploma, MSc and PhD from the Computer Engineering and Informatics Department, School of Engineering, University of Patras, Greece. His research interests include personalised software engineering, mobile software, secure identities software, web services, mobile ID, mDL, e-ID, e-health, and e-learning. He has more

than 90 publications in international journals and conferences at these areas. Since 1997, he has worked, designed, developed and coordinated a number of R&D national and European projects (i.e., ICT, Horizon, etc.). He has actively served as expert in mobile ID and secure documents European and international standardisation bodies.

---

## 1 Introduction

The increasing presence and spread of organised and not only fraud syndicates is of great concern to national and international institutions, as well as to private initiative actors. The impact of illegal activities extends beyond the loss of state revenue to undermining overall economic growth and breaking social cohesion. Organised economic crime currently adopts complex ways of organising and covering illegal activities at national or international level, making it more difficult for law enforcement to combat it.

The European Anti-Fraud Office (commonly known as OLAF, from the French: Office Européen de Lutte Antifraude) ([https://ec.europa.eu/anti-fraud/about-us/mission\\_en](https://ec.europa.eu/anti-fraud/about-us/mission_en)), which investigates cases of fraud against the EU budget, notes that the cost of financial fraud against the EU amounted to €3 billion for the year 2017. Similarly, fraudulent claim detection is one of the greatest challenges the insurance industry faces. Popular Electronic Commerce eShop Platform's return-freight insurance, providing return-shipping postage compensations over product return on the e-commerce platform, receives thousands of potentially fraudulent claims everyday (Liang et al., 2019).

The purpose of this research is to propose techniques to prevent, and combat forms of fraud with an emphasis on ensuring data privacy protection. In this way, it provides a critical upgrade to the operational capabilities of anti-fraud organisations and authorities. Datasets used in anti-fraud data mining usually include sensitive information or open data knowledge (Drakopoulou, 2018) which must be protected before one is being permitted to further process and to mine data.

The fight against financial crime requires complex administrative approval procedures and, finally, a special prosecutor's order for the competent authorities to access data sources due to sensitive data privacy protection requirements. Otherwise, personal data protection issues arise. In general, resource interconnection solutions and interoperability services must now take any data protection issues seriously before proceeding with data mining.

The strictest and most coherent framework is imposed by the General Data Protection Regulation 2016/679 of the EU (GDPR) ([http://ec.europa.eu/justice/policies/privacy/docs/95-46-ce/dir1995-46\\_part1\\_en.pdf](http://ec.europa.eu/justice/policies/privacy/docs/95-46-ce/dir1995-46_part1_en.pdf); DPIA, 2016; Regulation (EU) 2016/679, 2016). There are several more similar Data Protection Acts though enforced in other places around the world as the California Consumer Privacy Act (CCPA) (California Consumer Privacy Act (CCPA), 2018). Such data privacy laws insist that the protection of individuals with regard to the processing of personal data is a fundamental right and contributes to the achievement of an area of freedom, security and justice. In addition, they stipulate that the process of analysing and extracting information must comply with the principle of *data minimisation*, according to which personal data is appropriate, relevant and limited to what is necessary for the purposes for which it is processed.

A number of authorities act within the framework of anti-fraud activities in Greece, and straightforward equivalent authorities exist in most EU countries as well as further around the world. The legislature equips these public authorities with the right to access, collect, analyse and process information and data to combat financial crime at national and cross-border level. A further requirement of the legislator is the preparation of suggestions for the prioritisation of investigations for the targeted selection of ‘high profile’ cases, based on specific criteria. However, even under those mandates, the privacy preservation acts still need to be applied.

The Strategic Plan for the Fight against Corruption (2018–2021), which defines corruption as a multidisciplinary phenomenon, provides for coordination and synergy between the actors involved in the form of specific commitments and actions. Among other things, the need to develop operational actions related to the interconnection of information databases is emphasised in order to create reliable data files for the prosecution of crime, the establishment of a case monitoring procedure and the imposition of sanctions.

In order to provide a general solution for privacy preservation to anti-fraud authorities while being able to access sensitive data coming from many more data sources this paper specifically focuses on the frequent itemset hiding problem, a specific subdomain of knowledge hiding, where the goal is to sanitise a database from a set of sensitive frequent itemsets, in such a way that

- a those sensitive itemsets cannot be mined from the sanitised database
- b the quality of the sanitised data is maximised
- c the non-sensitive itemsets remain as close as possible to those that are mined from the original database.

A data curator should look into the results of the Apriori algorithm and should decide upon the sensitivity of the induced frequent itemsets, based on certain confidentiality rules and security regulations. Sanitising data may include record removal that affects the support of each itemsets in the resulting dataset. Maximising the quality of sanitised data is about the process of itemset handling (removal or even additions) in order to main support frequencies relationship throughout the database as intact as possible.

A major bottleneck in the hiding process is the tradeoff between the hiding of sensitive knowledge and the utility of the sanitised data. The hiding of the sensitive itemsets will have on the ideal positive border of the non-sensitive frequent itemsets alone.

The drawback of the border idea is its high complexity, since the Apriori algorithm needs to run first, in order to produce the frequent itemsets along with their negative border. Secondly, the sensitive patterns have to be selected and finally the initial border has to be revised, before the hiding algorithm can start running.

We adopt a level-wise algorithm for solving the constraint-based model formulation, by directly computing the revised border in one minimal phase, given that the sensitive itemsets are known beforehand (which is usually the case), without going back and forth traversing the itemset lattice. Additionally, an application of the proposed hiding model is presented. More specifically, we build a new exact approach for the frequent itemset hiding problem, by formulating a linear program based on the extension of the original database.

In summary, this research paper explores how sensitive data from different sources are hidden efficiently, in a way that allows mining of information coming from multiple systems and services, in accordance with the national and European institutional framework for data protection. The proposed approach is shown to be highly efficient after analysis and experimental evaluation.

The paper is organised as follows: Section 2 discusses privacy issues in mining and processing big data. Section 3 presents a short background on authorities that manage anti-fraud activities in Greece as well as data protection provisions of the law. Section 4 presents data hiding related work and a short introduction to the proposed approach. Section 5 dives into constraint based mining proposed approach specifics. Section 6 presents the experimental evaluation of the constraint-based mining algorithm. Section 7 describes the conceptual framework that we are proposing highlighting the overall steps needed for implementation. Section 8 discusses the conceptual framework for sensitive data hiding. Finally, Section 9 concludes the paper.

## **2 Privacy issues in big data**

Research still shows a data privacy preservation lag adding up to the anti-fraud process complexity and even hinder it from going further. In research community, there are reports that even the exchange of ideas in fraud detection and specifically in payment card fraud detection is severely limited due to security and privacy concerns (Ryman-Tubb et al., 2018; Sahin et al., 2013). Application of blockchain technologies though disruptive they need to get wider governmental acceptance and be regulated (Lacerda et al., 2021) in order to be used in anti-fraud cases. There are open information systems designed to enable transparency of public expenses and discourage corruption in the public sector (Gritzalis et al., 2019), however privacy risks of such approaches are identified even by supporters of the ‘right to know’ principle. On the governmental stakeholders side, prosecuting authorities are limited to e.g., in the participation and use of the existing networks of mutual administrative assistance and exchange of information as simple secure communication channels. Privacy awareness may be low and it needs to be increased (Sideri et al., 2019). However, there is a greater obstacle in the process. The lack of modern data analysis and mining tools for combining and utilising efficiently and with privacy different databases slows down the project and reduces their effectiveness especially as e-governance project implementers may, on top of that, struggle with the inter- and intra-organisational collaborations (Pandey and Suri, 2020). In addition, the inability to utilise information does not allow the audit authorities to take the next step, i.e., to implement the legislator’s requirement for a risk analysis system that will allow the prioritisation of audits and the selection of ‘high profile’ cases based on criteria and based on documentation with the help of analytical tools.

It is obvious that the effective fight against financial crime presupposes the selection of records from many sources into a single database, often distributed in different information systems, which belong to either public or private organisations (banking institutions, tax and customs authorities, telecommunications groups, telecommunications groups, energy, hospitals, insurance companies, etc.). An integrated data management approach is required, utilising interoperable tools and interconnection tools for many different data sources such as data warehousing.

The problem, in principle, is created due to the combination of sensitive data such as financial transactions and presence audit trails, that are coming from different and/or multiple distributed databases. It has been argued that the use of big data analysis technologies can have severe implications even for group privacy, including (political) targeting of particular groups (Mavriki and Karyda, 2019). As a result, multi-sourced data combinations bring a number of useful information but also *sensitive data fields* together endangering further data privacy in mining operations. This resulted in the rapid development of a specialised research area privacy preserving data mining (Clifton and Marks, 1996; O’Leary, 1991). The goal of this area is twofold. First, it includes the modification of the original database to exclude sensitive raw data as personally identifiable information (PII) (<https://www.gsa.gov/reference/gsa-privacy-program/rules-and-policies-protecting-pii-privacy-act>). Second, it protects sensitive knowledge that can be mined from a database.

The main objective of our proposal is to present intelligent techniques that will allow direct access to a dataset with sensitive information as described above. The general approach is to properly and efficiently hide sensitive information. In this way, the dataset will allow the controller to deal with and detect fraud without disclosing PII. If there is a fraud detected, then more steps may be taken for additional information disclosure.

Encryption is one approach that has been employed in order to protect privacy in fraud detection process. However, it may not scale well, as it increases up to an order of magnitude of five (5) the computation and communication overhead of processing such encrypted data as study shows (Canillas et al., 2018). Alternative approaches have also been proposed based on rule based approaches. Taxonomies are used to generalise concrete values appearing in fraud detection rules to higher level concepts which conform to some privacy/utility requirements set by the owner (Deutch et al., 2018). The approach however needs continuous fine-tuning and tweaking of the initial rules manually.

Within the data mining domain, more techniques have been introduced to deal with issues related to the privacy of the input data known as input privacy techniques (Evfimievski et al., 2004; Rizvi and Haritsa, 2002), as well as with the privacy of the induced knowledge in a data mining setting, known as output privacy techniques (Clifton, 1999; Kantarcioglu et al., 2004; Bu et al., 2007). The input privacy techniques are specialised in approaching the problem of how to guarantee the privacy of the input data during its publication by ensuring the maximum utility of the data for data mining purposes. The latter approaches take into consideration issues that touch upon the privacy of the induced patterns, and they aim at protecting the disclosure of sensitive patterns from the data in such a way that other non-sensitive patterns can be routinely produced from the so-called sanitised data (data from which sensitive knowledge has been removed). Because of the specific methodology that is used by these approaches to protect the sensitive patterns, they are collectively known as knowledge hiding techniques (Bonchi and Ferrari, 2011).

### **3 Data hiding related work**

A lot of techniques have been introduced to deal with privacy issues of the input data known as input privacy techniques (Evfimievski et al., 2004), as well as with the privacy

of the induced knowledge in a data mining setting, known as output privacy techniques (Kantarcioglu and Clifton, 2004). The former techniques aim at ensuring the maximum utility of the data for data mining purposes. Developed techniques in this category include various randomisation, perturbation and anonymisation algorithms. The latter approaches take into consideration the privacy of the induced patterns, and they aim at protecting the disclosure of sensitive patterns from the data in such a way that other non-sensitive patterns can be routinely produced from sanitised data.

In this paper we work on another research direction that is known as privacy preserving data mining. This kind of mining has as its goal to sanitise the data so that its privacy is protected against knowledge originating from the data itself. In such a scenario, adversarial rules that correlate values of sensitive attributes, are pinpointed and the data are minimally modified so that these rules are blocked out.

In particular, this paper focuses on the frequent itemset hiding problem, a specific subdo- main of knowledge hiding, where the goal is to sanitise a database from a set of sensitive frequent itemsets, in such a way that

- a those sensitive itemsets cannot be mined from the sanitised database
- b the quality of the sanitised data is maximised
- c the non-sensitive itemsets remain as close as possible to those that are mined from the original database.

The problem of hiding sensitive knowledge from the data mining process has been the field of a lot of active research since it was first introduced by Atallah et al. (1999) Following the aforementioned proof that the problem of hiding sensitive frequent patterns can have an NP-hard optimal solution, several more works were proposed to improve on their heuristic sanitisation algorithm, namely the works of Dasseni et al. (2001) and Kantarcioglu and Clifton (2004).

Verykios et al. (2004) extended the initial work (Dasseni et al., 2001) by proposing algorithms for hiding not only frequent itemsets but also association rules, and by evaluating these approaches according to different performance and data quality metrics. Verykios et al. (2007) propose a different approach by turning either 1's or 0's to question marks (implying unknown values), so that the hiding is achieved without falsifying the data.

Moustakides and Verykios (2008) to build an algorithm that implements the *maxmin* criterion and is similar in accuracy but much more efficient than the border-based algorithm. Gkoulalas-Divanis and Verykios (2009) applied the border-based principle too, in order to develop two linear programming techniques for optimally solving the hiding problem. The first technique (Lacerda et al., 2021), the so called inline approach, was introducing binary variables into the original database while the second one (Gkoulalas-Divanis and Verykios, 2009), known as hybrid, was extending the original database with synthetically generated transactions. In both approaches the goal was to fix the contents of specific items in the database, or in its extension thereof, so that to control the support of sensitive and non-sensitive itemsets.

The work proposed in this paper relies also on the constraint-based data mining area apart from the frequent itemset hiding area which was previously presented.

## 4 Data hiding introduction

Following the short related work we shall present an introduction to the approach for hiding (Verykios et al., 2019). Table 1 describes the notation used in the paper onwards.

**Table 1** Notation used in the paper

<i>Symbol</i>	<i>Description</i>
$D'$	The sanitised database obtained from in $D$
$F_D^\sigma$	Set of $\sigma$ -frequent itemsets in $D$
$S_D^\sigma$	Set of sensitive itemsets in $D$
$Bd^-(F_D^\sigma)$	negative border of $F_D^\sigma$
$Bd^+(F_D^\sigma)$	positive border of $F_D^\sigma$

Let  $I = \{i_1, i_2, i_3, \dots, i_n\}$  be a set of distinct literals called items. An itemset  $X$  is a nonempty subset of  $I$ , and a  $k$ -itemset is an itemset of length  $k$  (i.e.,  $|X| = k$ ). A transaction  $T$  over  $I$  is a 2-tuple  $T = \langle \text{tid}; t \rangle$  where  $\text{tid}$  is the identifier of transaction  $T$  and  $t$  is a set of items such that  $t \subseteq I$ . We say that a transaction  $T = \langle \text{tid}; t \rangle$  supports an itemset  $X$  iff  $X \subseteq t$ . A transaction database  $D$  is a collection of transactions. The support count of an itemset  $X$  in database  $D$ , denoted by  $\text{supc}_D(X)$ , is the cardinality of the set of transactions supporting  $X$ . Equivalently, we define the support of an itemset  $X$ , denoted by  $\text{sup}_D(X)$ , as the fraction of the support count of transactions supporting  $X$ , over the total count of transactions in the database  $D$ .

Given a user-specified support threshold  $\sigma$ , we call an itemset  $X$   $\sigma$ -frequent or, simply frequent, in  $D$  iff  $\text{sup}_D(X) \geq \sigma$ . Given a support threshold  $\sigma$ , and a database  $D$ , let  $F_D^\sigma$  be the collection of  $\sigma$ -frequent itemsets in  $D$ , where  $F_D^\sigma \subseteq P(I)$  and  $P(I)$  is the powerset of  $I$ . The positive border of the collection  $F_D^\sigma$  denoted as  $Bd^+(F_D^\sigma)$  is given by  $Bd^+(F_D^\sigma) = \{X \in P(I) \mid X \subset Y \text{ implies } Y \notin F_D^\sigma\}$ , while the negative border is given by  $Bd^-(F_D^\sigma) = \{X \in P(I) \mid Y \subset X \text{ implies } Y \notin F_D^\sigma\}$ .

Let  $S_D^\sigma \subseteq F_D^\sigma$  be the set of sensitive (frequent) itemsets that need to be hidden. Note here that  $F_D^\sigma$  and  $S_D^\sigma$  determine the ideal set  $F_D^\sigma$  of non-sensitive frequent itemsets based on the Apriori property. For example, if  $ab$  and  $abc$  belong to  $S_D^\sigma$  then it suffices that  $ab$  is hidden, since based on the antimonotonicity property of the Apriori,  $abc$  will also be hidden in the process. In order to ensure the minimum impact on the quality of the original database, the set  $Bd^-(S_D^\sigma)$  of the minimal, with respect of the above property, itemsets of  $S_D^\sigma$  should be transferred to the ideal negative border. Our goal then is to transform  $D$  to  $D'$  by selectively removing some items from the transactions of  $D$  in such a way that we minimise

- $|F_D^\sigma - \tilde{F}_D^\sigma|$  representing the number of the sensitive item-sets that are not hidden
- $|\tilde{F}_D^\sigma - F_D^\sigma|$  representing the number of hidden non-sensitive itemsets.

We give a formal definition of the problem in the sequel.

**Definition 1** (Frequent itemset hiding problem): Given a transaction database  $D$  over a set of items  $I = \{i_1, i_2, i_3, \dots, i_n\}$  a support threshold  $\sigma$ , and a set of sensitive frequent itemsets  $S_D^\sigma$  transform  $D$  to  $D'$  such that:

- 1  $\sup_D(X) < \sigma$  for every  $X \in S_D^\sigma$
- 2  $|F_D^\sigma - \tilde{F}_D^\sigma|$  is minimised
- 3  $|\tilde{F}_D^\sigma - F_D^\sigma|$  is minimised.

Due to the large number of frequent itemsets in  $\tilde{F}_D^\sigma$  it is computationally inefficient to keep track of all these itemsets during the selective removal of items for the hiding of the itemsets in  $S_D^\sigma$ . We can instead focus on the positive border of  $\tilde{F}_D^\sigma$ , which consists of the set of maximal non-sensitive frequent itemsets in  $\tilde{F}_D^\sigma$ . At the same time we make sure that we hide the minimal set of sensitive itemsets from  $S_D^\sigma$ . This minimal set comprises the negative border  $Bd^-(S_D^\sigma)$  of the sensitive itemsets in  $S_D^\sigma$ .

**Definition 2:** Given a transaction database  $D$  over a set of items  $I$ , a support threshold  $\sigma$ , and a set of sensitive frequent itemsets  $S_D^\sigma$  of  $D$ , the negative border  $Bd^-(S_D^\sigma)$  of  $S_D^\sigma$  is the set of minimal itemsets in  $S_D^\sigma$  with respect to set inclusion.

## 5 Data hiding approach using constraint based mining

At first, our approach relies on modelling the set of sensitive itemsets  $S_D^\sigma$  presented in the previous discussion, as a Boolean formula defined over a set of variables that correspond to the items of  $D$ . In this way we can take advantage of the simplification of this formula by using rules from the Boolean Algebra. In addition, we can easily get rid of the supersets of sensitive itemsets, that need to be neglected upon the computation of the borders.

Given a set of sensitive frequent itemsets  $S_D^\sigma$  of a transaction database  $D$  over a set of items  $I$ , we construct a Boolean formula  $B_D^\sigma$  over a set of variables  $V$  as follows: every item in  $I$  is mapped to a variable in  $V$  and every itemset of  $S_D^\sigma$  is mapped to a positive term, i.e., a conjunction of positive variables corresponding to the items of the itemset. The disjunction of all these terms comprises a Boolean formula  $B_D^\sigma$  in disjunctive normal form (DNF) with positive terms.

Continuing with our example database, the set of sensitive frequent itemsets  $S_D^\sigma = \{ac; bd; abc; acd\}$   $B_D^\sigma = ac + bd + abc + acd$ . By the absorption law, we can remove redundancies and simplify the original formula as  $B_D^\sigma = ac + bd$  (the rest of the terms are removed since they are subsumed by the term  $ac$ ). This irredundant positive DNF Boolean formula corresponds to the set of sensitive frequent itemsets  $S_D^\sigma$

**Definition 3:** The irredundant positive DNF Boolean formula  $B_D^\sigma = C_1 + C_2 + \dots + C_m$  where the  $i$ th term is of the form  $C_i = a_{i_1} \dots a_{i_2} \dots a_{|C_i|}$  is the Boolean formula obtained by the set of sensitive frequent itemsets  $S_D^\sigma$  when every item of  $D$  is mapped to a variable of  $B_D^\sigma$  and every itemset of  $S_D^\sigma$  is mapped to a positive term of  $B_D^\sigma$  after removing redundancies.

At this point, the only requirement regarding  $S_D^\sigma$  is the fact that the association among the items in the itemsets of  $S_D^\sigma$  should be protected. No prior assumptions on support or on other related metrics are made. This is a legitimate knowledge in the hands of a data curator that tries to block out knowledge from the data. Also notice that, by removing redundancies, every term  $C_j$  of  $B_D^\sigma$  is minimal (it is not subsumed by no other term of  $B_D^\sigma$ ), and maps a minimal sensitive itemset, say  $S_j$  in  $Bd^-(S_D^\sigma)$ .

**Lemma 1:** *The DNF Boolean formula  $B_D^\sigma$  corresponds to  $Bd^-(S_D^\sigma)$  the border-based theory (Verykios et al., 2019).*

**Theorem 1:** *If  $X$  is a non-sensitive frequent itemset of the ideal sanitised database, it corresponds to a negated pattern (a truth assignment with zero values) that satisfies  $\overline{B_D^\sigma}$  (Verykios et al., 2019).*

The following proposition states a trivial but important property of the Boolean formula  $\overline{B_D^\sigma}$  that is related to a constraint defined in the next section.

**Lemma 2:** *The Boolean formula  $\overline{B_D^\sigma}$  is equivalent to an irredundant negative DNF Boolean formula  $B$  (Verykios et al., 2019).*

The maximum number of terms in  $B = \prod_{i=1}^M |C_i|$  where  $C_i \in B_D^\sigma$

A Constraint-Based Theory for Mining of Borders presentation follows (Verykios et al., 2019). The next problem to deal with is how we can use the  $\overline{B_D^\sigma}$  in order to efficiently compute the required, by the hiding algorithm, borders  $Bd^+(\tilde{F}_D^\sigma)$  and  $Bd^-(S_D^\sigma)$

We address this problem by considering the computation of these borders as a constraint-based mining problem. It is well known that constraint-based mining allows the unearthing of interesting knowledge

- a by reducing the number of extracted patterns to only those of interest
- b by pushing constraints inside the mining algorithm in order to achieve better performance.

A constraint on itemsets is a function  $C: 2^T \rightarrow \{\text{true}, \text{false}\}$  We say that an itemset  $X$  satisfies a constraint  $C$  iff  $C(X) = \text{true}$ . We define the theory of a constraint  $C$ , as the set of itemsets that satisfies the constraint  $C$ , and we denote this theory by  $\text{Th}(C) = \{X \in 2^I \mid C(X) = \text{true}\}$ . Fortunately, as we will show soon enough, the constraints that we have to deal with in our hiding problem formulation are antimonotone constraints.

**Definition 4:** For every pair of itemsets  $X$  and  $Y$ , a constraint  $C$  is antimonotone if  $Y \subseteq X : C(X) \Rightarrow C(Y)$  (Verykios et al., 2019).

The first constraint in our problem formulation is the support constraint  $C_{sup}$  which is satisfied by itemsets having support greater than, or equal, to the support threshold  $\sigma$ . The support constraint is well known to be an antimonotone constraint. The second constraint, which is related to itemsets that are non-sensitive denoted as  $C_{sen}^-$ , should also hold for all the induced itemsets from the sanitised database. This constraint specifies that an interesting itemset for our hiding problem can be a proper subset of any sensitive itemset, but not a sensitive itemset or a superset of a sensitive itemset. The next proposition associates  $B$  to  $C_{sen}^-$ .

**Proposition 1:** The constraint  $C_{sen}^-$  holds for an itemset  $X$  if  $X$  satisfies  $B$ .

The antimonotonicity property of the  $C_{sen}^-$  constraint is shown analytically in Verykios et al. (2019).

As a result, we can easily push those constraints into the frequent itemset mining algorithm since it is well known that any conjunction of antimonotone constraints is also an antimonotone constraint. The  $Th_D(C_{sup} \wedge C_{sen}^-)$  that we are looking for, is then a downward closed theory which means that if an itemset  $X$  is part of this theory, then all subsets of  $X$  belong to this theory, too. This reminds us of the antimonotone heuristic in Apriori which will be utilised by our hiding algorithm as well, for efficiently computing the theory of itemsets satisfying the conjunction of the two constraints. The algorithm is presented in the sequel.

Following we present the Constraint-Based Mining Algorithm proposed. The constraint-based itemset mining algorithm will generate the  $STh_D(C_{sup} \wedge C_{sen}^-)$  in order to create the positive border  $Bd^+(STh_D(C_{sup} \wedge C_{sen}^-))$  of the itemsets that satisfy both constraints  $C_{sup}$  and  $C_{sen}^-$ . These itemsets will be needed in the next phase by the linear programming-based hiding solution. Initially, the mining algorithm stores in  $C_1^B$  all the candidate items from the set of items  $I$  that satisfy the Boolean formula  $B$ . Then it stores in  $L_1$  these items from  $C_1^B$  which were found to be frequent, after counting them through the first pass in the database. In the subsequent for-loop, the level-wise operation of the algorithm unfolds. Based on the specific level  $k$ , the algorithm computes in  $C_k^B$  the candidate itemsets of the  $k$ -th level by applying the Apriori-genB procedure to the set of frequent and non-sensitive itemsets of size  $k - 1$ , that comprises the collection of itemsets  $L_{k-1}$ . Apriori-genB performs a 2-way pruning changing the approach of the standard Apriori-gen. Initially, it removes from  $C_k^B$  candidates of size  $k$  with either infrequent or sensitive subsets of size  $k-1$ , and then it removes sensitive candidates of size  $k$ . A sensitive candidate of size  $k$  is a candidate generated by Apriori-genB that does not satisfy  $C_{sen}^-$ , even if all of its subsets are frequent and non-sensitive (that is, all of its subsets belong to  $L_{k-1}$ ). Next, the support of candidates in  $C_k^B$  is counted and all those candidates, that are found to be frequent, are placed into  $L_k$ . When  $L_k$  becomes empty, the algorithm terminates by returning the union of  $L_k$  set collections generated so far.

## Algorithm

**Input:**  
D: transaction database  
 $\sigma$ : the minimum support threshold  
B: Boolean formula representing  $C_{\text{sen}}$

**Output:**  
L: frequent itemsets satisfying B

Steps:

- 1: Description:
- 2:  $C_1^B$  the candidate items that satisfy B
- 3: L1 the frequent items in  $C_1^B$
- 4: for ( $k = 2$ ;  $L_{k-1} \diamond 0$ ;  $k++$ )
- 5:  $C_k^B = \text{Apriori-genB}(L_{k-1})$ ;
- 6: for each transaction  $t \in D$
- 7:  $C_t = \text{subset}(C_k^B, t)$ ;
- 8: for each candidate  $c \in C_t$
- 9:  $c:\text{count}++$ ;
- 10: end
- 11:  $L_k = \{c \in C_k^B \mid c.\text{count} \geq \sigma\}$ ;
- 12: end
- 13: return  $L = \cup_k L_k$ ;

## 6 Experimental evaluation of the constraint-based mining algorithm

In this section, we present evaluation of the proposed Constraint-based Mining Algorithm compared with the conventional Level-wise Apriori Algorithm (Verykios et al., 2019). Evaluation includes real datasets are available in the Frequent Itemset Mining Dataset Repository. The rest of the synthetic datasets were generated using the IBM Basket Data Generator. We utilised datasets with a variety of characteristics in terms of the number of transactions, number of items, and average transaction length. The mushroom dataset was prepared by Roberto Bayardo (Inan et al., 2010). The BMS1 and BMS2 datasets were used for the KDD Cup 2000. Finally, the kosarak dataset (deVries et al., 2011) contains anonymised click-stream data of a Hungarian online news portal. We summarise the details of the datasets in Table 2.

Algorithms were implemented by using or extending the PyFIM extension module of Borgelt (2012), so as to compute efficiently the non-sensitive frequent itemsets and the corresponding positive border. All experiments were performed on a personal computer with an Intel Core i5 at 3.4 GHz processor.

Both algorithms take as input the dataset  $D$ , the set of sensitive itemsets  $S_D^\sigma$  and the support threshold  $\sigma$ . The evaluation metric we used is the execution time required by each algorithm in order to give output, which is the set of non-sensitive frequent itemsets  $\tilde{F}_D^\sigma$ , and in extension its positive border  $Bd^+(\tilde{F}_D^\sigma)$ . The number of sensitive itemset was 500 for all experiments.

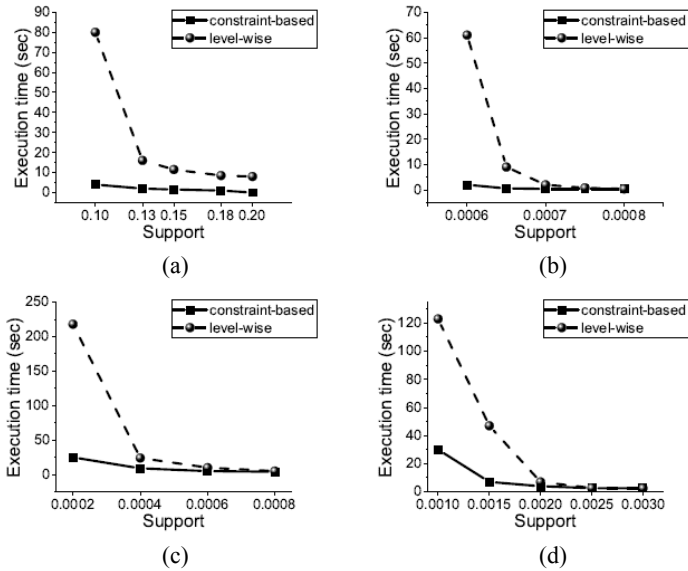
To show the potentials of the approach, we present the execution time (in seconds) for both Level-wise and Constraint-based algorithms, for the baseline datasets of

mushroom, BMS1, BMS2, and kosarak dataset respectively (Verykios et al., 2019), see Figure 1. The baseline datasets has been used widely as this is encouraging researchers to make their work as reproducible as possible and compare the potential algorithm in setups with different diversity of transactions volume and the itemsets length (<https://archive.ics.uci.edu/ml/datasets.php>). Observe that in all scenarios, the Constraint-based algorithm outperforms the Level-wise algorithm.

**Table 2** Dataset descriptions

<i>Dataset</i>	<i>Transactions #</i>	<i>Items #</i>	<i>Avg. transactions length</i>
Mushroom	8,124	119	23
BMS1	59,602	497	2.5
BMS2	77,512	3,340	5.6
kosarak	990,002	41,270	8.1

**Figure 1** Sensitive hiding using constraint based algorithm vs level-wise time performance: (a) execution time for *mushroom*; (b) execution time for *BMS1*; (c) execution time for *BMS2* and (a) execution time for *kosarak*



## 7 Conceptual framework towards a general anti-fraud approach

### 7.1 Case of a simple anti-fraud personal data multi source integration

The former sensitive data hiding technique is shown to have real life potentials in protecting private data fields and in parallel allowing data mining in order to be used within an anti-fraud approach. In order to explain further its application we provide below a conceptual framework towards a general anti-fraud approach. Initially we ll

provide the necessary steps to achieve multi-source data alignment. To enable mining is a complex process, consisting of the following steps.

The first step is data pre-processing. It is now common for database entries to contain errors and inconsistencies (Rahm and Do, 2000; Hernandez and Stolfo, 1998; Widlak and Peeters, 2020). A pre-processing is required in order for the records to acquire a well-defined format. The organisations that are going to participate in the record integration process must proceed with similar pre-processing actions.

The second step concerns the indexing of the records, which applies to all records of the databases to be processed/mined. The purpose is to protect or reduce the number of sensitive fields while including as many records as possible that are referring to the same entity.

The third step concerns the integration of the compared records.

A fourth step concerns the technical review that is the non-automated examination of them with the intervention of an experienced specialist to verify that.

The fifth step concerns the evaluation of the whole process in terms of the degree of complexity, quality of the interface and its completeness.

At step four, though, it is usually possible that privacy preservation is practically impossible obstructing further processing for authorities.

## *7.2 Case of sensitive data hiding – privacy-preserving approach*

The organisations that own PII in question are often reluctant to share information, as there are significant restrictions on data protection and confidentiality. The process of integrating database records in the possession of different organisations must in this case result into a privacy preservation resource that will protect sensitive data and will not disclose critical information unnecessarily to the case file reviewers. However, the need for investigation and cross-referencing of data is great in the context of anti-fraud investigation. Investigations shall only go further only if there are a minimum number of indications before accessing even more sensitive information. So the process of private sensitive data handling is particularly delicate.

The stage of data pre-processing remains the same as in the case of simple data integration. However, once the pre-processing has taken place, the sensitive data is hidden following the approach discussed in Section 4 before attempting to further mine them, to ensure the protection and confidentiality of the data.

The proposed approach protects privacy while the data are being processed keeping the accuracy as high as possible limiting the noise introduced. In the same use case, differential privacy has been employed however, heavily limiting accuracy. Most federated learning systems therefore use differential privacy to introduce noise to the parameters. This adds uncertainty to any attempt to reveal private client data, but also reduces the accuracy of the shared model, limiting the useful scale of privacy-preserving noise. Attempts to compensate the lost accuracy demand multi-party secret sharing that may lead to potential privacy leakage. Such a system could reduce the coordinating server's ability to recover private client information, without additional accuracy loss, by also including secure multiparty computation (SMC) (Truex et al., 2019).

SMC is, however, inherently susceptible to inference. Furthermore, SMC can be heavily attacked using DDoS approach (Moustakides et al., 2014). There is implementation where Hardware Security Modules (HSM) running Trusted Execution Environments (TEE) are used to protect computation for each party. However, all such

approaches demand high maintenance secure key ceremony to be in place to initialise the HSM & TEE system securely in the first place. The proposed approach has brought the key advantage that the multi-party datasets are prepared on the relying party side and it does not need any complex trust environment to run on each side. Furthermore, in the anti-fraud use case retrieving the data in a timely manner is critical therefore high maintenance and long installation procedures may lead to loss of fraud detection potential.

## **8 Discussion on the conceptual framework for sensitive data hiding**

The overall approach on the anti-fraud value of this research proposal can be found in the following points:

- a speeding up the process of searching and analysing data records that include sensitive fields with reduced complexity without endangering privacy
- b the challenge of Parallel sensitive data hiding, where mining into multiple sensitive records is attempted in an environment of different and multiple distributed large databases with the help of parallel processors.

We particularly focus on the first point where the proposed algorithm must reliably ensure the protection of sensitive data and ultimately the quality privacy preservation without being complicated.

There are a lot of potential use cases. Anti-fraud control in foreign investments is usual in a lot of countries in several parts of the world. The particular process includes multiple data sources that need to be included. Namely, to perform just a typical real estate investment there are a lot of different databases and fields that need to be mined from systems such as

- i border control personal traveller's data (photo, personal details, passport, visa etc.)
- ii bank account data of representative escrow accounts under a different delegate name (i.e., representation lawyer company name)
- iii amount transaction data that indicate source bank account details (IBAN, SWIFT, Routes etc.). Schengen entry-exit data for the case of EU that include extremely sensitive private data such as fingerprints for identification purposes
- iv national cadastre details that provides the details of the real estate which also may include the full list of all previous owner for up to the last 80 years revealing a lot of information for additional third parties private data.

It is clear that in such a use case there are a lot of sensitive data that need to be protected and at the moment protection is applied by trust. On the other hand, there is a need to extract several reports to monitor and check for potential fraud in such transactions for cross-country investments are several times under scrutiny to prevent fraud and money laundering. Our proposed framework would mark as sensitive data all personal data such as lawyer company name, bank name, investor's country name, real estate locality and may exclude most of deeply sensitive personal data.

Regarding the second point, the parallel protected sensitive data is a new scenario that is under research and initial results show that it can be a realistic approach for modern

information ecosystems. An already implemented solution to the challenge is to leverage Hadoop's MapReduce technologies. Initial research shows that protected data can be scaled up in the case of parallel systems (Krasadakis et al., 2020).

The innovation of the research proposal is particularly interesting in the case of anti-fraud as it falls directly into key aims:

- a its inclusion in the framework of the revised research priorities of the National Strategy RIS3, as determined by the General Secretariat for Research and Technology in Greece as well as in key pillars of EU Research Priorities
- b the fact that it starts with the fight against financial crime, but covers and finds application in a wide range of forms of fraud such as insurance fraud and theft
- c the fact that it opens up new perspectives and opportunities for the fight against insurance and financial fraud.

Regarding the first point, the definition of the revised research priorities of the National Strategy RIS3 for the Informatics Communications and Telecommunications Sector, intervention areas are prioritised in particular as far as it concerns:

- a in content and information management technologies
- b in privacy and security of personal data
- c in electronic identification of persons (e-ID), objects and electronic information.

For these areas of intervention, the following directions of innovation are favoured using the proposed approach:

- creation and development of cross-sectoral solutions through the utilisation of large volumes of data and learning techniques in heterogeneous data
- secure data management and sharing
- efficient data encryption and anonymisation algorithms and/or data masking/obfuscation techniques
- electronic identification interoperability techniques
- recording interconnection technologies between third party systems to speed up and ensure security checks.

The presented approach is within the above thematic directions, because:

- a it provides privacy preserving solution for sensitive and multidisciplinary nature of financial and insurance fraud
- b it promotes the adoption, improvement and practical utilisation of the data hiding approach to the implementation of secure data management and sharing with respect for the principles of personal data protection
- c it poses the base for the improvement of data anonymisation algorithms and their scaling in the case of Big Data in cloud and parallel infrastructure.

Regarding the second point (fight against a wide range of forms of fraud), the Electronic Government of Social Security, whose mission is to support the public registers to strengthen social solidarity services and citizens, is faced with cases of insurance fraud.

The policy of consolidation of the funds of former Social Security Institutions in a Single Social Insurance Institution, as well as the effort to streamline the multiple and overlapping benefit policies, intensify the need for an integrated sensitive data management in order to carry out procedures and controls to capture, at a central level, the necessary information e.g., for the provision of the medical history of each patient in the nursing units or the calculation and liquidation of the insurance contributions.

Moreover, the social security stakeholders participate, in collaboration within EU, in the development of cross-border applications in the fields of health and social security (HealthID, eIDAS Cross Border Services, etc.), constituting the National Node of the European System for Accessing and Exchanging Social Security Data (EESI).

Regarding the third point (opening new perspectives and opportunities in the fight against fraud) the following are noted:

First, the research approach removes significant barriers to the use of interoperable solutions, which are related to serious privacy issues. The sharing of data using hiding for the integration of records allows the full utilisation of interoperable solutions in accordance with the EU General Data Protection Regulation.

Second, the fight against fraud is being facilitated and accelerated. For example, the fight against financial crime is slow, as, among other things, a prosecutor's permit is required in order for the competent authorities to access databases of third parties. The competent prosecutorial authorities then and only then should be summoned, in order to reveal more sensitive values of the records under investigation. The benefit is therefore twofold, as the task of cross-checking information is accelerated and the prosecution is summoned at the last and most critical stage.

Third, the competent bodies acquire techniques that allow them to modernise their actions in accordance with international recommendations and standards. The Intra-European Organization of Tax Administrations (IOTA) (<https://www.iota-tax.org/about-us>) notes that a comprehensive strategic approach to combating financial crime, based on data mining, needs to be developed, including:

- a the graphical representation of networks (social network analysis), as they are formed through the analysis of transactions (transactions analysis) recorded in databases of the interbank and stock exchange sector, the telecommunications sector, the sector of state social benefits for welfare
- b the development of a risk analysis system for case selection and prioritisation of investigations
- c the use of machine learning algorithms.

The adoption of secure anti-fraud sensitive data hiding practices will allow the authorities to remedy any information shortages. The enrichment of the database of law enforcement authorities will allow the next step, namely the construction of a complex anti-fraud strategy, with features as described in the IOTA reports.

Handling of sensitive data within enclaves that are hardware backed secure computation environment is recent approach that shall also be connected to this current discussion of sensitive data hiding approach for anti-fraud. Collaborative confidential computing or trusted computing (Confidential Computing Consortium, 2020) concerns the protection of data in use. This is indeed the new frontier for mobile, desktop and server based computations. In technical terms, this is about solutions such as trusted

execution environments, mobile and server hardware security modules that provide physical protection in sensitive data computations.

In a complementary manner, sensitive data hiding actually goes further ahead and works in wider and on top of such a trusted computing. Sensitive data hiding is about techniques that will enable end-users, researchers, companies and organisations to access and work on data that may not directly and clearly appear to be private and sensitive. Nevertheless such data are deemed as such, because they can lead to implicitly revealing real identity of user/customer/holder such as location of transactions, even coarse location details, race and nationality details, product categories etc.

## 9 Conclusions and future steps

In this paper, the data privacy based complexity of anti-fraud procedures is discussed in order to highlight the need to actively provide a toolset to authorities to perform their tasks under the umbrella of privacy preservation acts and laws such as GDPR, CCAP etc. We proposed an approach of privacy preserving data mining and specifically data hiding as an efficient tool to speed up complex anti-fraud processes. To meet this goal, an efficient privacy preserving hiding approach is discussed with the view to assist anti-fraud mining before requesting special data access to PII from legal authorities. Our proposal is a constraint-based hiding model that decreases the pre-processing overhead incurred by border-based techniques in the hiding of sensitive frequent itemsets. The hiding model lends itself to the efficient computation of the theory of itemsets which are both frequent and non-sensitive with a constraint-based Apriori-like algorithm. These can be used in border-based hiding algorithms to effectively conceal the sensitive knowledge. This is particularly important in anti-fraud as big data is needed to be processed and mined in an automated manner.

The mining algorithm is coupled with a conceptual framework that provides a step by step approach to show how hiding can be useful in anti-fraud activities where data are extensive though full access and analysis is not permitted at full extent without prior legal notice.

Future steps include the further evaluation of approaches that function in parallelisation on premise or in cloud based infrastructure of sensitive data mining. The research team is particularly interested to address this question in order to advance the findings by re-exploiting approaches from the field of artificial intelligence and machine learning.

## References

- Atallah, M., Elmagarmid, A., Ibrahim, M., Bertino, E. and Verykios, V. (1999) ‘Disclosure limitation of sensitive rules’, *Proceedings of the 1999 Workshop on Knowledge and Data Engineering Exchange*, IEEE Computer Society, Chicago, Illinois, pp.45–52.
- Bauer Lacerda, A. and Dias Lopes, F. (2021) Applications of blockchain technology in the Brazilian government, *International Journal of Electronic Governance*, Vol. 13, No. 2, pp.132–148
- Bonchi, F. and Ferrari, E. (2011) *Privacy-Aware Knowledge Discovery: Novel Applications and New Techniques*, Chapman and Hall/CRC Data Mining and Knowledge Discovery Series, CRC PressINC.

- Borgelt, C. (2012) 'Frequent item set mining', *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, Vol. 2, No. 6, pp.437–456.
- Bu, S., Lakshmanan, L.V.S., Ng, R.T. and Ramesh, G. (2007) 'Preservation of patterns and input-output privacy', *ICDE*, pp.696–705.
- California Consumer Privacy Act (CCPA) (2018) <https://oag.ca.gov/privacy/ccpa>
- Canillas, R., Talbi, R., Bouchenak, S., Hasan, O., Brunie, L. and Sarraf, L. (2018) 'Exploratory study of privacy preserving fraud detection', *Middleware '18: Proceedings of the 19th International Middleware Conference Industry*, December, pp.25–31
- Clifton, C. (1999) 'Protecting against data mining through samples', *DBSec*, pp.193–207.
- Clifton, C. and Marks, D. (1996) 'Security and privacy implications of data mining', *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data (SIGMOD '96)*, February, pp.15–19.
- Confidential Computing Consortium (2020) Confidential Computing Deep Dive v1.0. 10/2020 <https://confidentialcomputing.io/wp-content/uploads/sites/85/2020/10/Confidential-Computing-Deep-Dive-white-paper.pdf>
- Dasseni, E., Verykios, V., Elmagarmid, A.K. and Bertino, E. (2001) 'Hiding association rules by using confidence and support', *Proceedings of the 4th International Workshop on Information Hiding*, Pittsburgh, PA, USA, pp.369–383.
- Deutch, D., Ginzberg, Y. and Milo, T. (2018) 'Preserving privacy of fraud detection rule sharing using intel's SGX', *CIKM '18: Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, October, Turin, Italy, pp.1935–1938.
- deVries, T., Ke, H., Chawla, S. and Christen, P. (2011) 'Robust record linkage blocking using suffix arrays and Bloom filters', *ACM Transactions on Knowledge Discovery from Data*, Vol. 5, No. 2, pp.1–27.
- DPIA (2016) Guidelines on Data Protection Impact Assessment (DPIA) and determining whether processing is "likely to result in a high risk" for the purposes of Regulation 2016/679
- Drakopoulou, S. (2018) 'Open data today and tomorrow: the present challenges and possibilities of open data', *International Journal of Electronic Governance*, Vol. 10, No. 2, pp.157–171
- Evfimievski, A.V., Srikant, R., Agrawal, R. and Gehrke, J. (2004) 'Privacy preserving mining of association rules', *Inf. Syst.*, Vol. 29, No. 4, pp.343–364.
- Gkoulalas-Divanis, A. and Verykios, V.S. (2006) 'An integer programming approach for frequent itemset hiding', *CIKM*, pp.748–757.
- Gkoulalas-Divanis, A. and Verykios, V.S. (2009) 'Exact knowledge hiding through database extension', *IEEE Trans. Knowl. Data Eng.*, Vol. 21, No. 5, pp.699–713.
- Gkoulalas-Divanis, A. and Verykios, V.S. (2009) 'Hiding sensitive knowledge without side effects', *Knowl. Inf. Syst.*, Vol. 20, No. 3, pp.263–299.
- Gritzalis, A., Tsohou, A., Lambrinoudakis, C. (2019) 'Transparency-enabling information systems: trust relations and privacy concerns in open governance', *International Journal of Electronic Governance*, Vol. 11, Nos. 3–4, pp.310–332.
- Hernandez, M.A. and Stolfo, S.J. (1998) 'Real-world data is dirty: data cleansing and the merge/purge problem', *Data Mining and Knowledge Discovery*, Vol. 2, No. 1, pp.9–37.
- Inan, A., Kantarcioglu, M., Ghinita, G. and Bertino, E. (2010) 'Private record matching using differential privacy', *EDBT*, pp.123–134.
- Ishai Y., Ostrovsky R. and Zikas V. (2014) 'Secure multi-party computation with identifiable abort', in Garay, J.A. and Gennaro, R. (Eds.): *Advances in Cryptology – CRYPTO 2014. CRYPTO 2014. Lecture Notes in Computer Science*, Vol. 8617, Springer, Berlin, Heidelberg.
- Kantarcioglu, M. and Clifton, C. (2004) 'Privacy-preserving distributed mining of association rules on horizontally partitioned data', *IEEE Transactions on Knowledge and Data Engineering*, Vol. 16, No. 9, pp.1026–1037.
- Kantarcioglu, M., Jin, J. and Clifton, C. (2004) 'When do data mining results violate privacy?', *KDD*, pp.599–604.

- Krasadakis, P., Verykios, V. and Sakkopoulos, E. (2020) 'Parallel based hiding of sensitive knowledge', to appear in *ICTAI 2020: International Conference on Tools with Artificial Intelligence* Baltimore, Maryland, USA, to appear.
- Liang, C., Liu, Z., Liu, B., Zhou, J., Li, X., Yang, S. and Yuan, Q. (2019) 'Uncovering insurance fraud conspiracy with network learning', *SIGIR'19: Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, July, pp.1181–1184.
- Mavriki, P. and Karyda, M. (2019) 'Big data in political communication: implications for group privacy', *International Journal of Electronic Governance*, Vol. 11, Nos. 3–4, pp.289–309.
- Moustakides, G.V. and Verykios, V.S. (2008) 'A maxmin approach for hiding frequent itemsets', *Data Knowl. Eng.*, Vol. 65, No. 1, pp.75–89.
- O'Leary, D.E. (1991) 'Knowledge discovery as a threat to database security', *Proceedings of the 1st International Conference on Knowledge Discovery in Databases*, Montréal, Québec, Canada, pp.507–516.
- Pandey, J.K. and Suri, P.K. (2020) 'Collaboration competency and e-governance performance', *International Journal of Electronic Governance*, Vol. 12, No. 3, pp.246–275.
- Rahm, E. and Do, H.H. (2000) 'Data cleaning: problems and current approaches', *IEEE Data Eng. Bull.*, Vol. 23, No. 4, pp.3–13.
- Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)
- Rizvi, S. and Haritsa, J.R. (2002) 'Maintaining data privacy in association rule mining', *VLDB*, pp.682–693.
- Ryman-Tubb, N.F., Krause, P. and Garn, W. (2018) 'How artificial intelligence and machine learning research impacts payment card fraud detection: a survey and industry benchmark', *Engineering Applications of Artificial Intelligence*, Vol. 76, No. 2018, pp.130–157.
- Sahin, Y., Bulkan, S. and Duman, E. (2013) 'A cost-sensitive decision tree approach for fraud detection', *Expert Syst. Appl.*, Vol. 40, pp.5916–5923.
- Sideri, M., Kitsiou, A., Tzortzaki, E., Kalloniatis, C. and Gritzalis, S. (2019) 'Enhancing university students' privacy literacy through an educational intervention: a Greek case-study', *International Journal of Electronic Governance*, Vol. 11, Nos. 3–4, pp.333–360.
- Truex, S., Baracaldo, N., Anwar, A., Steinke, T., Ludwig, H., Zhang, R. and Zhou, Y. (2019) 'A hybrid approach to privacy-preserving federated learning', *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security (AISec'19)*, Association for Computing Machinery, New York, NY, USA, pp.1–11.
- Verykios, V.S., Elmagarmid, A.K., Bertino, E., Saygin, Y. and Dasseni, E. (2004) 'Association rule hiding', *IEEE Trans. Knowl. Data Eng.*, Vol. 16, No. 4, pp.434–447.
- Verykios, V.S., Pontikakis, E.D., Theodoridis, Y. and Chang, L. (2007) 'Efficient algorithms for distortion and blocking techniques in association rule hiding', *Distributed and Parallel Databases*, Vol. 22, No. 1, pp.85–104.
- Verykios, V.S., Stavropoulos, E.C., Zorkadis, V. and Elmagarmid, A.K. (2019) 'A constraint-based model for the frequent itemset hiding problem', *8th International Conference e-Democracy 2019*, 12–13 December, Athens, Greece
- Widlak, A. and Peeters, R. (2020) 'Administrative errors and the burden of correction and consequence: how information technology exacerbates the consequences of bureaucratic mistakes for citizens', *International Journal of Electronic Governance*, Vol. 12, No. 1, pp.40–56.

## **Websites**

Directive 95/46/EC of the European Parliament and of the Council of 24 October 1995 on the protection of individuals with regard to the processing of personal data and on the free movement of such data, [http://ec.europa.eu/justice/policies/privacy/docs/95-46-ce/dir1995-46\\_part1\\_en.pdf](http://ec.europa.eu/justice/policies/privacy/docs/95-46-ce/dir1995-46_part1_en.pdf)

European Anti-Fraud Office (OLAF) (3333) [https://ec.europa.eu/anti-fraud/about-us/mission\\_en](https://ec.europa.eu/anti-fraud/about-us/mission_en)

Intra-European Organization of Tax Administrations (IOTA) <https://www.iota-tax.org/about-us>

Personally Identifiable Information (PII), Rules and Policies - Protecting PII - Privacy Act 2019, <https://www.gsa.gov/reference/gsa-privacy-program/rules-and-policies-protecting-pii-privacy-act>

UCI Machine Learning Repository, Center for Machine Learning and Intelligent Systems, <https://archive.ics.uci.edu/ml/datasets.php>