

Presence Patterns and Privacy Analysis

Citation for published version (APA):

Roubtsova, E. E., Roubtsov, S. A., & Alpár, G. (2018). Presence Patterns and Privacy Analysis. In W. van der Aalst, J. Mylopoulos, M. Rosemann, M. J. Shaw, & C. Szyperski (Eds.), *Business Modeling and Software Design: 8th International Symposium, BMSD 2018. Vienna, Austria, July 2-4, 2018. Proceedings* (Vol. 319, pp. 298-307). Springer. https://doi.org/10.1007/978-3-319-94214-8_21

DOI:

[10.1007/978-3-319-94214-8_21](https://doi.org/10.1007/978-3-319-94214-8_21)

Document status and date:

Published: 30/06/2018

Document Version:

Publisher's PDF, also known as Version of record

Document license:

Taverne

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

<https://www.ou.nl/taverne-agreement>

Take down policy

If you believe that this document breaches copyright please contact us at:

pure-support@ou.nl

providing details and we will investigate your claim.

Downloaded from <https://research.ou.nl/> on date: 10 Oct. 2024

Open Universiteit
www.ou.nl





Presence Patterns and Privacy Analysis

Ella Roubtsova¹(✉), Serguei Roubtsov², and Greg Alpar^{1,3}

¹ Open University of the Netherlands, Heerlen, The Netherlands
{Ella.Roubtsova,Greg.Alpar}@ou.nl

² Technical University Eindhoven, Eindhoven, The Netherlands
s.roubtsov@tue.nl

³ Radboud University, Nijmegen, The Netherlands
g.alpar@cs.ru.nl

Abstract. Business applications often use such data structures as Presence Patterns for presentation of numbers of customers in service-oriented businesses including education, retailing and media. Presence Patterns contain open data derived from internal data of organizations. In this paper, we investigate different ways of defining Presence Patterns and possible privacy consequences dependent on the chosen definition. The first contribution of the paper is a definition of a family of Presence Patterns. The second contribution is a method for privacy analysis of Presence Patterns.

Keywords: Presence Patterns · Privacy requirements
Method for privacy analysis of presence patterns

1 Introduction

Businesses nowadays provide valuable Internet services to their customers or employees which require registration of customers or employees. The registration data contain information about login and logout that can correspond to the ‘presence’ of registered users. In this paper, we assume that ‘presence’ online means actual or physical ‘presence’¹. This information about ‘presence’ can be used for different business analytics purposes. Business applications often use such data structures as Presence Patterns to analyse participation of customers in service-oriented businesses, including education, retailing and media.

However, businesses are not allowed to use personal data without limitations. The use of personal data is regulated by laws and, in particular, by the General

¹ The login and logout do not always correspond to the physical presence of the logged persons. Nevertheless, the degree of correspondence can be always established by manual control of presence during a chosen time period by sampling. Usually, this degree of correspondence is high, as many people use their mobile phones or PC to connect to the organizational network because this connection provides useful information.

Data Protection Regulation (GDPR)². The major privacy rule, coming from the GDPR, allows businesses to use only the minimum data needed for the given task and does not allow them to relate business indicators to personal data of clients. The way to meet this privacy requirement is to pay attention to new business concepts used by business analytics, carefully design the data structures for all new business concepts and avoid relations with personal data in their definitions.

Because of the GDPR, companies (and organisations in general) have to comply with the privacy-by-design principles [2, 7]. Businesses have a hard time to figure out how to do this. They also experience pressure from the technology side. Ever more possibilities and new opportunities emerge every day. However, these new technologies are often coupled with extensive data usage making hard to assess how privacy invasive their deployment can be.

In this paper, we investigate different ways of defining so-called Presence Patterns and possible privacy violations caused by their coupling with the organisation's files and also with, possibly relevant, publicly-available information. Section 2 presents the work related to protection of privacy. Section 3 defines a family of Presence Patterns. Section 4 defines a method for privacy analysis of Presence Patterns. Section 5 illustrates the use of the method with a case study. Section 6 presents conclusions and future work.

2 Related Work

2.1 Privacy, Anonymity and k -anonymity

Privacy studies conclude that there are quite a few wrong beliefs about privacy protection. Sweeney [14] reports that it is believed that removing explicit identifiers, such as name, address, telephone number makes data anonymous and protects privacy. However, the studies show that combinations of a few publicly available characteristics “uniquely or nearly uniquely identify some individuals”.

The Internet contains many independent data sources. They may contain the same personal information with different extensions and, as a result, may partially release this personal information. In this sense, the Internet is similar to a multi-level database. It is formally proven by Su and Ozsoyoglu [13], that in general, it is impossible to guarantee privacy in a multi-level database due to functional dependencies and multi-level dependencies of data.

The public data can come from registers (phonebook-like databases), social networks (LinkedIn, Facebook, etc.), public ratings like on IMDb, Amazon or other websites and webshops. Moreover, we increasingly rely on infrastructures that collect, store and (possibly) publish data. For example, “smart” driver assisting services are aware of the current vehicle position and can potentially emit lots of data about your vehicle to the outside world. These data items may not be identifying in themselves, but when being aggregated, they can cause almost unforeseeable privacy violations.

² General Data Protection Regulation (GDPR)(EU) 2016/679 GDPR official: <http://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A32016R0679>.

Sweeney [14] proposed a model “for understanding, evaluating and constructing computational systems that control inferences” about private information. The model is based on the definition of a **quasi-identifier**.

“Given a population of entities (a universe) U , an entity-specific table $T(A_1, \dots, A_n)$, $f_c: U \rightarrow T$ and $f_g: T \rightarrow U'$, where $U \subseteq U'$. A quasi-identifier of T , written Q_T , is a subset of attributes $\{A_i, \dots, A_j\} \subseteq \{A_1, \dots, A_n\}$ where: $\exists p_i \in U$ such that $f_g(f_c(p_i)[Q_T]) = p_i$ ”.

In other words, a quasi-identifier Q_T of a table T is a subset of table’s columns named by attributes (A_1, \dots, A_j) . A tuple of values of the attributes of the quasi-identifier, that appears in a row of table T , is a value or an occurrence of the quasi-identifier. This value of the quasi-identifier is used for re-identification of an entity from the universe U .

Sweeney gives an example, of a voter-specific table (universe U) and a health patient-specific table (universe U'). In the voter-specific table, there are fields $\{\text{Name, Address, ZIP, Birth date, Sex}\}$. The medical data contain $\{\text{Visit date, Diagnosis, ZIP, Birth date, Sex}\}$. The quasi-identifier $\{\text{ZIP, Birth date, Sex}\}$ can be used to re-identify the medical information of the voters. He concludes that “the attributes that appear in the private data and also appear in the public data are the candidates for linking: therefore, these attributes constitute the quasi-identifier and the disclosure of these attributes must be controlled”. For such a control, he proposes an additional requirement of k -anonymity for any released table (publicly available table) T with an associated quasi-identifier.

The k -anonymity of any table T means that each tuple of attribute-values of the table’s quasi-identifier Q_T appears at least in k -rows of this table $T[Q_T]$.

Let it be two tables, a private table PT and a public (released or open) table T . It is proven that, “if the released data T satisfies k -anonymity with respect to the quasi-identifier Q_{PT} of a private table PT , then the combination of the released data T and the external sources on which Q_{PT} was based, cannot link on Q_{PT} or a subset of its attributes to match fewer than k individuals” [14].

2.2 Inference Attacks and Open Data

The k -anonymity of open (and released) data cannot protect from specific data mining techniques, called inference attacks, which use data queries, aggregation of data, sorting, etc. However, the attacks demand resources. The open data should not make the life of an attacker easier.

The most representative sources of open data are the social networks. The attacks on social networks are described in [12]. This work states that only a lack of resources can stop attackers from massive crawling via API or “screen-scraping” [10]. Gross and Acquisti [6] demonstrate that the attributes of the nodes (rows in a entity-specific data table), such as social security numbers and other profile attributes [3], can be predicted with higher accuracy than random guessing.

Besides the nodes (rows in a data table), the social network expose edges, being relations between nodes. Edges provide additional information about nodes (persons) and their behaviour patterns [9].

Social networks provide some protection from attackers. Attackers need to create many dummy nodes helping to re-identify social network members. The online social networks “check the uniqueness of e-mail addresses, and deploy other methods for verifying accuracy of supplied information making creating of dummy nodes relatively difficult” [12].

Nevertheless, on small scale, for defined targets found in the released (open) data, the attacks can be feasible. So, the awareness about quasi-identifiers in any released data should be raised. In this paper, we talk about the released data in form of Presence Patterns.

3 Presence Pattern

In this section we define a family of Presence Patterns: a Simple Presence Pattern and several extensions of it for real-world applications. In the next section we use these definitions to propose a method for privacy analysis of such patterns.

3.1 Simple Presence Pattern

Presence of someone or something is the state of existing of someone or something in a given place in a given time interval. The given place can be described as a class room, a shop, but may also be seen as the state “online” in a media. Presence Patterns are often released for public. For example, a polyclinic may release a presence pattern as patients may choose a less visited day and time interval for their visit.

A Simple Presence Pattern in business usually concerns with numbers of present customers. As businesses are interested in performance and revenues, they are interested in the number of customers who are in the state “In” during a given time interval. They can estimate how many human and other resources are needed for service of customers. We define a Simple Presence Pattern as a mapping of numbers of present customers onto natural time intervals defined during a week.

Simple Presence Pattern = (Days, Hours, NumberOfPresent,
function-times, function-presence)

- Days={Monday, Tuesday, Wednesday, Thursday, Friday, Saturday, Sunday}.
- Hours={1,2,3,24}.
- NumberOfPresent is a finite natural number (positive integer number or zero).
- function-times: Hours \rightarrow Days divides days of the week into time intervals.
- function-presence:NumberOfPresent \rightarrow (Hours \rightarrow Days) assigns the numbers of customers to the time intervals for each day of the week.

3.2 Presence Pattern Extended with a Schedule

Any real Presence Pattern used for a given business purpose is usually an extension of the Simple Presence Pattern. The most widely used sources of extension are schedules.

```
Schedule = (Days, Times, OpeningStatus, function-times,
            function-opening)
```

- Times= (Morning work hours 9:00-12:00,
Evening work hours 13:00-16:00,
Closing times 16:00-9:00);
- OpeningStatus = {Open,Closed}.
- function-times: Times \rightarrow Days;
- function-opening = OpeningStatus \rightarrow (Times \rightarrow Days).

A schedule is easily included into the Presence Pattern.

```
Presence Pattern = (Days, Times, NumberOfPresent, OpeningStatus,
                   function-times, function-opening, function-presence).
```

3.3 Presence Pattern and Schedule with Business Information

In practice, schedules contain also sets and relations [1] with business related information. For instance, in an educational institution, an address with the room number (say, "R206"), and a course name (for example, "Physics") and a teacher name usually correspond to some time interval.

```
Presence Pattern Education =(Days, Times, NumberOfPresent,
Locations, Courses, Teachers, function-times, function-location,
function-course, function-teacher, function-presence).
```

- function-location = Address \rightarrow (Times \rightarrow Days)
- function-course = Courses \rightarrow (Times \rightarrow Days)
- function-teacher = Teachers \rightarrow (Times \rightarrow Days)

This Presence Pattern Education contains private information about the service provider: addresses of an educational institute and the names of teachers.

3.4 Presence Pattern, Event Log and User Profile

Another data structure that is closely related to Presence Patterns is an event log. Industrial event logs are emitted by operating systems and applications to record, among other information, online events. The record of any occurrence of a login or logout of a network device always contain a time stamp and may contain private information like the user IP address or the MAC address of a mobile device or a PC, activity (associated or disassociated) with the time stamp. An event record often contains additional information such as the profile (or the access right) of the user (guest, employee, etc), protocol and so on.

A time stamp, an activity (state) and a User Identifier are often considered as the main fields of any industrial log.

Event log = (Time stamp, Activity, User, function-activity-time,
function-user-time)

- Time stamp is a set of time stamps, for example, 12.07.2017.10:25.
- Activity = {In, Out}. For example, In means login or check-in; Out means Logout or check-out.
- User is a set of user unique IDs such as names or addresses of personal devices.

Event logs do not belong to any Presence Pattern, but they are used to derive the calculated information for Presence Patterns. For this, the event logs can be automatically processed in order to produce the `NumberOfPresent`. `NumberOfPresent` has a *one-to-many* relation with the users in the log. For every time interval in the schedule, the automated query derives the users logged in within this time interval and logged out within or after this time interval. If for a user there is a time stamp In before or within the given Time Interval i and there is a stamp Out within or after it, then this user is present in the given time interval and the counter `NumberOfPresent` is increased by 1. By counting, the private information is transformed into the statistical information, i.e. the numbers of customers per a time interval of the schedule.

The internal user (customer, employee) profiles of companies can be used to extend Presence Patterns for particular business purposes. For example, if a user profile contains fields `User`, `Gender`, `Address`, `Date of birth`, then the event log and the customer profile can be used to derive the `NumberOfPresent` people of a given gender, or of a given age or registered in a given city. In the extended form, a Privacy Pattern contains information about a subset of entities from a private universe (women, students following a given education, cars, citizens of a chosen city), that are present at a given time at a given place.

For example, the Presence Pattern `Gender` may be used to indicate the presence of women:

Presence Pattern Gender = (Days, Times, NumberOfPresent, Gender,
Locations, Courses, Teachers, function-times, function-location,
function-course, function-teacher, function-presence).

- Gender = {Male, Female};
- `NumberOfPresentGender` is a finite natural number (positive integer number or zero).
- `function-presence`: `NumberOfPresentGender` \rightarrow (Times \rightarrow Days) assigns the numbers of present woman to the time intervals for each day of the week.

Another example, Presence Pattern `City` (Fig.1) shows a Presence Pattern extended with the information about present entities for a chosen home-city.

Presence Pattern City = (Days, Times, City,
NumberOfPresentFromCity, function-times,
function-presence-from-city)

- City = {Eindhoven, Uden}.
- `NumberOfPresentFromCity` is a finite natural number (positive integer number or zero).

- `function-presence-from-city: NumberOfPresentFromCity → (Times → Days)` assigns the numbers of customers to the time intervals for each day of the week.

4 Method for Privacy Analysis of Presence Patterns

We propose a method for privacy analysis of Presence Patterns. It is based on

1. our definition of a family of Presence Patterns given in Sect. 3 and
2. the fact, discussed in Sect. 2, that k -anonymity of a quasi-identifier in released data can be validated on internal data. The combination of the released data and the external data cannot link on this quasi-identifier to match fewer than k entities.

The value of k for k -anonymity is defined in advance, depending on the importance of privacy. A general ‘rule of thumb’ should be that an individual person could be determined from this k -subset with the effort more significant to the perpetrator compared to the privacy violation benefit.

Our method uses the following assumptions about public data:

- (1) The Internal data universe contains more than k entities.
- (2) Publicly available data structures may contain data fields with private information with the same data types as the company Internal data used for design of extended Presence Patterns;
- (3) Publicly available data structures do not contain event logs of any other company-related activities of any entity.

Simple Presence Pattern does not have any Privacy Issues. A simple Presence Pattern (SPP) consists of the time schedule fields and the field, `NumberOfPresent`. The values of `NumberOfPresent` are derived from the `event log` that is not public. This field, `NumberOfPresent` is always related to more than k entities of the Internal data.

Privacy Analysis of Extended Presence Patterns. Let an Analysed Presence Pattern (APP) consist of the same fields as the SPP and a set E of some other fields.

1. Initially, the set of fields for privacy analysis is empty, $SF = \emptyset$. The set of fields causing privacy violation is empty, $SV = \emptyset$.
2. Add a set E of fields in APP to SF , $SF = SF \cup E$.
 - for each value sf of SF ,
 - find in the Internal data the maximal value of `NumberOfPresent` and
 - find out if the maximal `NumberOfPresent` for a given value sf identifies less than k entities.
 - If YES, then add the quasi-identifier E to SV : $SV = SV \cup E$.
3. If $SV \neq \emptyset$, then the release of this Presence Pattern can cost privacy violation, otherwise the Presence Pattern respects privacy and can be released.

This procedure gives a systematic way for privacy analysis of Presence Patterns using their specific structure.

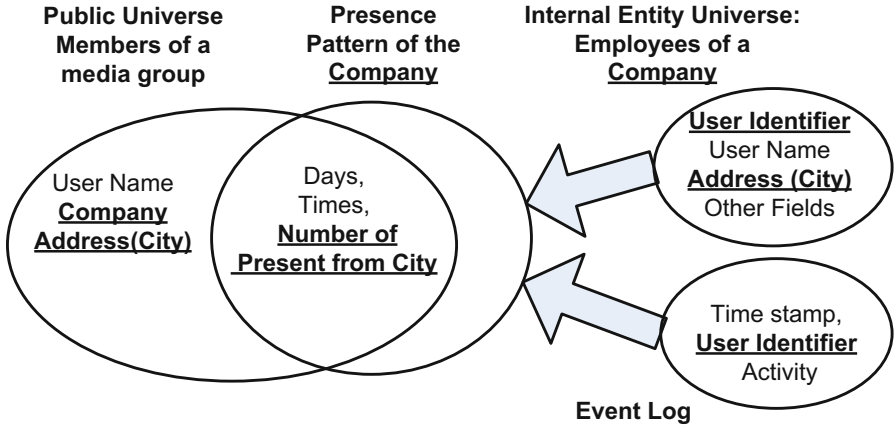


Fig. 1. Quasi-identifier in a presence pattern

5 A Case Study: Usage of Parking Space

Let us consider a case study. A company would like to monitor the use of the parking space in its premises. The entrance to the parking places has an automatic licence plate recognition system. The following internal (not public) data is available:

- an internal log of all cars that arrive to and leave from the parking. The log records include a time stamp and a licence plate number.
- an internal database of cars of employees. Any record in this database presents an employee’s name and his (her) car licence plate number.
- the database of all employee’s profiles. Among other fields, it contains the address of each employee (Fig. 1).

Further public information may be relevant to this use case:

- the LinkedIn (or other professional social network) profiles of employees of the company;
- The telephone guide, each record of which includes a person name, his (her) phone number (where the city can be identified);
- The public registry of cars containing the licence plate numbers with the corresponding brand and the model of the car, engine type, etc.³

In the investigation of the use of the parking space, the company creates new public data in the form of a series of Presence Patterns. The *k* for *k*-anonymity is set to 5 based on expert recommendations.

³ This kind of data is publicly available in many countries (for example the RDW registry in the Netherlands).

1. First, the company derives and makes public the Simple Presence Pattern with the **Number Of Present Cars at the Parking Space** during the working hours and outside working hours. This pattern does not have any privacy issues as the **Number Of Present Cars at the Parking Space** is related to all employees of the company possessing a car (say to 500 employees).
2. Next, the company desires to analyze how many electric cars are parked daily in order to arrange the necessary amount of charging sockets. The internal database of coming cars has to be extended with the information about the type of the engine in order to be able to identify cars with electric engines. The presence pattern is extended with the indicator (or data type) **ElectricOrNot**. This information can be derived from the public car registry using the license plate number. The **Number Of Present Electric Cars at the Parking Space** during and outside the working hours is presented by the new presence pattern.

The analysis reveals, that this presence pattern cannot identify less persons than the number of employees having electric cars (say, 120).

3. The employees who live nearby the company premises, unlike the ones who live far away, may be convinced to avoid using cars. To be able to monitor its anti-pollution effort, the company needs to know how many distantly leaving employees are present during and outside the working hours.

The pattern of presence is extended accordingly with the **City** taken from the addresses employees. So, the new pattern of presence shows the **Number Of Present Electric Cars at the Parking Space** from each **City**.

The analysis shows that for at least one city (say, for Uden) the maximal **Number Of Present**=2. It is less than $k=5$. This pattern has a possible privacy violation if just one or a small number of employees (less than a predefined k) lives in a particular city.

The quasi-identifier is (**Company**, **City**) (Fig. 1). From this pattern, one can derive new information, namely, a pattern of presence of persons working in the given company and living in the given city. This can be done using the combination of the public LinkedIn-like profiles of employees of the company and the public telephone guide.

This case study illustrates how the method identifies possible privacy violations caused by the released Presence Patterns.

6 Conclusion and Future Work

Privacy-by-design is a common goal in the field of information technologies. It is defined in [2] as seven fundamental principles: privacy as proactive, privacy by default, privacy embedded into design, full functionality (respecting other aspects of the system), end-to-end life cycle protection, respect for user privacy. These principles should be refined to methods of privacy analysis during design.

Following the principles ‘privacy as proactive’ and ‘privacy embedded into design’ we have proposed a proactive method for analysis of Presence Patterns

released by data analytics. Our method is built on our inductive definition of a family of Presence Patterns and the refinement of privacy requirement by k -anonymity for the specific structure of Presence Patterns.

Presence Patterns are widely used by education institutions [4, 11], by the television and radio companies, team sports [5], trends in tourism [8] and voting sites. Observing all these domains, we plan to investigate and evaluate the application of our method for pattern definition and privacy analysis on publicly available Presence Patterns and other released data structures derived from internal data of different organizations.

References

1. Aziz, N.L.A., Aizam, N.A.H.: A survey on the requirements of university course timetabling. *World Acad. Sci. Eng. Technol. Int. J. Math. Comput. Phys. Electr. Comput. Eng.* **10**, 236–241 (2016)
2. Cavoukian, A.: Privacy by design: the definitive workshop. A foreword by Ann Cavoukian, Ph.D. *Identity Inf. Soc.* **3**, 247–251 (2010)
3. Chew, M., Balfanz, D., Laurie, B.: (Under) mining privacy in social networks. Citeseer (2008)
4. Ganji, F., Budzisz, L., Debele, F.G., Li, N., Meo, M., Ricca, M., Zhang, Y., Wolisz, A.: Greening campus WLANs: energy-relevant usage and mobility patterns. *Comput. Netw.* **78**, 164–181 (2015)
5. Griffin, P.S.: Girls' participation patterns in a middle school team sports unit. *J. Teach. Phys. Educ.* **4**, 30–38 (1985)
6. Gross, R., Acquisti, A.: Information revelation and privacy in online social networks. In: *Proceedings of the 2005 ACM Workshop on Privacy in the Electronic Society*, pp. 71–80 (2005)
7. Hoepman, J.-H.: Privacy design strategies. In: Cuppens-Bouahia, N., Cuppens, F., Jajodia, S., Abou El Kalam, A., Sans, T. (eds.) *SEC 2014. IAICT*, vol. 428, pp. 446–459. Springer, Heidelberg (2014). https://doi.org/10.1007/978-3-642-55415-5_38
8. Kim, H., Cheng, C.-K., O'Leary, J.T.: Understanding participation patterns and trends in tourism cultural attractions. *Tour. manage.* **28**, 1366–1371 (2007)
9. Korolova, A., Motwani, R., Nabar, S.U., Xu, Y.: Link privacy in social networks. In: *Proceedings of the 17th ACM Conference on Information and Knowledge Management*, pp. 289–298 (2008)
10. Mislove, A., Marcon, M., Gummadi, K.P., Druschel, P., Bhattacharjee, B.: Measurement and analysis of online social networks. In: *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement*, pp. 29–42 (2007)
11. Mulhanga, M.M., Lima, S.R., Carvalho, P.: Characterising university WLANs within eduroam context. In: Balandin, S., Koucheryavy, Y., Hu, H. (eds.) *NEW2AN/ruSMART -2011. LNCS*, vol. 6869, pp. 382–394. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-22875-9_35
12. Narayanan, A., Shmatikov, V.: De-anonymizing social networks. In: *2009 30th IEEE Symposium on Security and Privacy*, pp. 173–187 (2009)
13. Su, T.-A., Ozsoyoglu, G.: Controlling FD and MVD inferences in multilevel relational database systems. *IEEE Trans. Knowl. Data Eng.* **3**, 474–485 (1991)
14. Sweeney, L.: k -anonymity: a model for protecting privacy. *Int. J. Uncertain. Fuzziness Knowl. Based Syst.* **10**, 557–570 (2002)