



UNICA

UNIVERSITÀ
DEGLI STUDI
DI CAGLIARI



Università di Cagliari

UNICA IRIS Institutional Research Information System

This is the author's accepted version of:

Nergiz, M.E., Atzori, M., Clifton, C.W. (2025). δ -Presence . In: Jajodia, S., Samarati, P., Yung, M. (eds) Encyclopedia of Cryptography, Security and Privacy. Springer, Cham.

The publisher's version is available at:

https://doi.org/10.1007/978-3-030-71522-9_1570

Users may view, print, copy, download and text and data-mine the AM content, for the purposes of academic research, subject always to these Terms of Use. Any AM content downloaded for text and data-mining purposes must be deleted or destroyed when the analysis is complete.

Use must not be for Commercial Purposes.

All use must be fully attributed.

When citing this work, please cite the original published paper

This full text was downloaded from UNICA IRIS <https://iris.unica.it/>

δ -Presence

Mehmet Ercan Nergiz¹, Maurizio Atzori², and Christopher W. Clifton³

¹Idea Teknoloji Cozumleri, Eclipse Maslak, Istanbul, Turkey

²DMI, University of Cagliari, Cagliari, Italy

³CS Department, Purdue University, West Lafayette, IN, USA

Definition

δ -presence is a probabilistic privacy standard that protects against membership attacks, attacks identifying an individual as a member of an exposed group. The probabilistic definition allows risk-cost-benefit analysis on real-world applications letting a data owner choose appropriate privacy parameters.

Background

You just realize that a recent article published in a local newspaper refers to your university. The article has the title “What are the odds? The National Lottery reports jackpots won by two Italian academicians from the same university.” As an Italian academician, should you be concerned about your privacy? Obviously, the answer depends on how many Italian academicians are working in your institution. If the answer is two, congratulations, you are officially super-rich. The bad news is, however, everybody now knows about it. What if the answer is three? There is now only two out of three chances that you are rich. Privacy wise, the odds are certainly more assuring but should you feel safe enough considering that the percentage of “super-rich” people within the population is less than 1%? What would you lose for every 1% increase in others’ beliefs on your wealth? Certainly, our perception of a privacy loss due to an information exposure depends on how much of an increase/decrease in belief beyond normal takes place plus how much we lose for every unit of such an increase/decrease.

The above scenario is an example showing that revealing even statistics (e.g., count) regarding a group of people sharing a sensitive characteristics (e.g., jackpot winners, patients of a sensitive disease, minorities, members of a particular race or religion, people with a specific sexual orientation, etc.) *might* be a threat to privacy. The privacy standard δ -presence is designed for such scenarios where the privacy risk is from identifying that an individual is a member of an exposed group. Two main properties of the metric allow effective privacy protection against membership disclosure:

- δ -presence probabilistically limits an adversary’s ability to identify individuals as being a member of a group. The parameter δ controls the trade-off between privacy and utility.

- The probabilistic definition allows setting a good δ -parameter via a risk-cost-benefit analysis based on the underlying application. This forms the necessary but often neglected bridge between human-understandable policy and mathematically sound standards for anonymity.

δ -presence has been adopted in the industry, being now one of the four privacy metrics available on the Google Cloud Data Loss Prevention service (Google Inc. 2019).

Theory

δ -Presence: Generic Definition

In our setting, the data owner has knowledge on a group of people G within a population P . Every individual in G shares a sensitive characteristic; thus knowledge of G is private while that of P is public. Due to privacy concerns, the owner chooses to release $f(G)$ (e.g., statistics regarding G). Thus we have

Data owner knowledge: G, P

Adversary knowledge: $f(G), P$

Given this setting, δ -presence is a privacy standard for deciding if the disclosed information $f(G)$ is privacy-preserving.

Definition 1 (Generic δ -presence) (The original definition of δ -presence enforces also a lower bound limit on the membership probability, providing a range (δ_{\min} - δ_{\max}) of allowed membership probabilities (Nergiz et al. 2007). Since the upper bound limit is typically a bigger concern, we drop the lower bound requirement ($\delta_{\min} = 0$) here for a concise presentation.) Let $G \subset P$ be a group of people sharing sensitive characteristics drawn from a population P . Let function f be any function computed on G . We say $f(G)$ satisfies δ -presence if and only if $\forall p \in P$

$$P(p \in G \mid f(G), P) \leq \delta$$

We call $P(p \in G \mid f(G), P)$ the *membership probability* of p .

Note that in our setting the perceived privacy risk is from discovering that an individual is a member of G . δ -presence limits the membership probability, preventing an adversary from identifying, with high probability, that any particular individual within the population is a member of G .

Following our example in the previous section, P is the set of faculty members in the university, G is the set of people hitting jackpots, and $f(G)$ is the number of “lucky” Italian faculty members in the university. Let n be the number of Italians in P . Assuming each Italian has the same prior probability of winning a jackpot, the membership probability for each Italian faculty member is $f(G)/n$; δ -presence is satisfied if $f(G)/n < \delta$.

Data Publishing Scenario: Enforcing δ -Presence

In a data publishing scenario, the data owner has a private table TCP of a group of people sharing sensitive characteristics (e.g., a diabetes institute having demographics of diabetic patients) where P is a public table of the whole population (e.g., the census dataset of the region). Although the use of a public table P has been exploited in the so-called weak k -anonymity (Atzori 2006) to protect a private table T from linking attacks, in many applications the data owner wants to protect against membership attacks described earlier; thus the owner wants to ensure that any information shared on T satisfies δ -presence.

An interesting and practical property of δ -presence is that it can be enforced over an existing table by applying standard generalization techniques such as those used in k -anonymity (Sweeney 2002; Samarati 2001; Atzori 2006). Thus, the owner can create a generalization T^* of T such that $\forall p \in P$

$$P(p \in T \mid T^*, P) \leq \delta$$

It is worth noting that publishing a generalized table T^* obtained by enforcing k -anonymity over T may not prevent membership attacks. To exemplify this, suppose we want to share information on T to promote research (Table 1). In case T is shared as it is, an attacker, by accessing publicly

δ -Presence, Table 1
Public dataset P and research subset T

P					T			
Publicly known data					Research subset			
	Name	Zip	Age	Nationality	Zip	Age	Nationality	
<i>a</i>	Alice	47906	35	USA	<i>b</i>	47903	59	Canada
<i>b</i>	Bob	47903	59	Canada	<i>c</i>	47906	42	USA
<i>c</i>	Christine	47906	42	USA	<i>f</i>	47633	63	Peru
<i>d</i>	Dirk	47630	18	Brazil	<i>h</i>	48972	47	Bulgaria
<i>e</i>	Eunice	47630	22	Brazil	<i>i</i>	48970	52	France
<i>f</i>	Frank	47633	63	Peru				
<i>g</i>	Gail	48973	33	Spain				
<i>h</i>	Harry	48972	47	Bulgaria				
<i>i</i>	Iris	48970	52	France				

(Initial “key” columns for clarity only)

δ -Presence, Table 2 T_1^* :
5-anonymization of T

P					T_1^*		
Publicly known data					Research subset		
	Name	Zip	Age	Nationality	Zip	Age	Nationality
<i>a</i>	Alice	47906	35	USA	4*	> 40	*
<i>b</i>	Bob	47903	59	Canada	4*	> 40	*
<i>c</i>	Christine	47906	42	USA	4*	> 40	*
<i>d</i>	Dirk	47630	18	Brazil	4*	> 40	*
<i>e</i>	Eunice	47630	22	Brazil	4*	> 40	*
<i>f</i>	Frank	47633	63	Peru			
<i>g</i>	Gail	48973	33	Spain			
<i>h</i>	Harry	48972	47	Bulgaria			
<i>i</i>	Iris	48970	52	France			

available dataset P , can easily identify Bob as being in T as Bob is the only Canadian in P . Publishing a k -anonymization of T does not help.

In Table 2, T_1^* is a 5-anonymization for T ($k=5$ is the most strict k -anonymity setting for a table of 5 rows). Despite most of the information being removed, T_1^* is still prone to membership attacks. In fact, by noting that there are exactly 5 people in P with age > 40 (in gray), the attacker can derive with certainty that Bob, Christine, Frank, Harry, and Iris are all in T , while Alice is not in T . In other words, the membership probabilities for the five people are all 1.

Nevertheless, a generalization that preserves $2/3$ -presence is possible as shown in Table 3. Given the generalization T_3^* and the public table P , the attacker can at best identify each person in P as in T with probability at most $2/3$, as shown by groups in light and dark gray in P and T^* . For instance, Gail (as well as Harry and Iris) matches one of the two dark gray rows in T_3^* . Assuming initially each of these three people is equally likely to be in the research set, the posterior membership probability for each person after observing T_3^* is $2/3$.

The probabilistic definition of δ -presence allows risk-cost-benefit analysis so that a data owner can decide on a good privacy parameter δ that achieves utility at the cost of acceptable and quantifiable privacy loss. Work in 3 presents an example of such an analysis in the context of data sharing on diabetic patients. The analysis takes into account:

- the difference between the posterior adversary belief on the membership probability, δ , and the prior adversary belief before seeing the published generalization (e.g., 0.07 – the

δ -Presence, Table 3 T_2^* :
 $\frac{2}{3}$ -present generalization
of T

P					T₃[*]		
Publicly known data					Research subset		
	Name	Zip	Age	Nationality	Zip	Age	Nationality
<i>a</i>	Alice	47906	35	USA	47* *	*	America
<i>b</i>	Bob	47903	59	Canada	47* *	*	America
<i>c</i>	Christine	47906	42	USA	47* *	*	America
<i>d</i>	Dirk	47630	18	Brazil	48* *	*	Europe
<i>e</i>	Eunice	47630	22	Brazil	48* *	*	Europe
<i>f</i>	Frank	47633	63	Peru			
<i>g</i>	Gail	48973	33	Spain			
<i>h</i>	Harry	48972	47	Bulgaria			
<i>i</i>	Iris	48970	52	France			

probability that a random individual will have diabetes).

- real-world risk incurred to an individual exposed as a result of privacy violation (e.g., a patient identified as being a diabetic might face discrimination in job interviews.)
- given the real risk, the cost to individuals per unit of increase in belief (e.g., annual cost, in dollars, of a diabetic employee to a company due to health insurance)
- acceptable cost of publishing to an individual (e.g., the maximum amount of increase, due to publishing, in expected cost, in dollars, of a potential employee that would be considered negligible compared to other costs.)

The analysis builds the necessary bridge between the mathematical privacy parameter δ and the economical and social costs that are incurred due to privacy violations. Given an acceptable upper bound on such a cost to an individual, the adversary can compute a lower bound on the privacy parameter δ and effectively control damage (in real-world units) due to sharing of anonymizations.

Open Problems

Similar to other anonymization-based privacy standards in the literature, the most notable vulnerability of the δ -presence standard is that the data owner is expected to model the adversary background knowledge accurately. In a data publishing scenario, if the adversary underestimates the adversary's prior knowledge on individuals, the resulting generalizations may not satisfy δ -presence as the membership probabilities can be estimated beyond the acceptable threshold δ . (For example, considering the public table P and generalization $T2^*$ in Table 3, Gail is considerably younger than Iris. If the adversary knows that diabetes is a disease often seen in older people, he/she can conclude that Iris is more likely to be in T compared to Gail. In such a case, the membership probability of Gail given $T2^*$ will certainly be higher than $2/3$.) Work in Nergiz and Clifton (2009) partially addresses this issue by assuming an upper bound on the adversary knowledge (e.g., the adversary knows P , but the owner only knows statistics on P) and tries to enforce δ -presence in the worst possible scenario. However, new practical techniques that can model adversary background in a more flexible fashion still stands as a future work. Another open question is how this relates to noise-based privacy protection approaches such as ϵ -differential privacy (Dwork et al. 2006). This is particularly difficult, as it can be impossible to bound the probabilities without some limitations on adversary background knowledge. Related discussion is given in work on Differential Identifiability (Lee and Clifton 2012).

Cross-References

- ▶ [Data Linkage](#)
- ▶ [Differential Privacy](#)
- ▶ [k-Anonymity](#)
- ▶ [Privacy Metrics](#)
- ▶ [Quasi-Identifier](#)

References

- Atzori M (2006) Weak k-anonymity: a low-distortion model for protecting privacy. In: Information security, 9th international conference, ISC 2006, vol 4176 of LNCS, pp 60–71. Springer
- Dwork C, McSherry F, Nissim K, Smith A (2006) Calibrating noise to sensitivity in private data analysis. In: Proceedings of the 3rd theory of cryptography conference, pp 265–284
- Google Inc (2019) Google Data Loss Prevention API 2019. Accessed 3 June 2020
- Lee J, Clifton C (2012) Differential identifiability. In: The 19th ACM SIGKDD conference on knowledge discovery and data mining, pp 1041–1049, Beijing, 12–16 Aug 2012
- Nergiz ME, Atzori M, Clifton C (2007) Hiding the presence of individuals from shared databases. In: Proceedings of the ACM SIGMOD international conference on management of data. ACM, pp 665–676
- Nergiz ME, Clifton C (2009) δ -Presence without complete world knowledge. IEEE Trans Knowl Data Eng 22:868–883
- Samarati P (2001) Protecting respondents' identities in microdata release. IEEE Trans Knowl Data Eng 13(6):1010–1027
- Sweeney L (2002) k-anonymity: a model for protecting privacy. Int J Uncertainty Fuzziness Knowledge Based Syst 10(5):557–570