# Measuring the Value of Privacy and the Efficacy of PETs

Kimmo Halunen
VTT Technical Research Centre of Finland
kimmo.halunen@vtt.fi

Anni Karinsalo
VTT Technical Research Centre of Finland
anni.karinsalo@vtt.fi

## ABSTRACT

Privacy is a very active subject of research and also of debate in the political circles. In order to make good decisions about privacy, we need measurement systems for privacy. Most of the traditional measures such as $k$-anonymity lack expressiveness in many cases. We present a privacy measuring framework, which can be used to measure the value of privacy to an individual and also to evaluate the efficacy of privacy enhancing technologies. Our method is centered on a subject, whose privacy can be measured through the amount and value of information learned about the subject by some observers. This gives rise to interesting probabilistic models for the value of privacy and measures for privacy enhancing technologies.

## CCS CONCEPTS

•**Security and privacy → Privacy protections;**

## KEYWORDS

Privacy, measurement, metric, probability, value

## 1 INTRODUCTION

Privacy has become a focal point of research in the information security community and also in the society at large. The revelations of Edward Snowden and the realisation that many large companies now have and collect more and more information about us as individuals has made the erosion of privacy visible to many stakeholders. Some are proposing new regulation, e.g. [22], to tackle the impact of new technology and business models on our privacy.

The collection of data of almost every transaction that happens online has profound effects on privacy both at an individual and societal level. Some claim, that data is the pollution of the internet age [1] and argue that it should be regulated as such. Thus, there is a need for tools to help regain privacy.

A key problem with privacy is that it is hard to measure accurately and correctly with a uniform scale in all situations. Measures of anonymity such as $k$-anonymity [21], entropy and anonymity

[1] https://www.schneier.com/blog/archives/2008/01/data_as_polluti.html

sets are used to measure the privacy level. These are well-developed, some even widely used methods, but in our viewpoint anonymity is not the same as privacy and measuring privacy needs to consider also other aspects than mere anonymity. In addition, anonymity can be mixed with pseudonymity, leading to false sense of privacy.

We survey existing definitions, measures and metrics for privacy and anonymity and a new approach to measure privacy. Our measure is centered around a subject, whose privacy is under scrutiny. We develop this measure as a probabilistic expectation over certain functions that relate to the observed value of the private information. Our measures give several ways to evaluate the effectiveness of privacy enhancing technologies (PETs). In conclusion we discuss and analyze our findings and lay ground for future work.

## 2 PRIVACY AND ANONYMITY

The concepts of privacy and anonymity are often difficult to distinguish even in the scientific context. Many studies do not explicitly separate these and use them interchangeably.

Anonymity is defined in [18] as the state of being not identifiable within a set of subjects, the anonymity set, or in another words, unlinkability of an item of interest and an identifier of a subject. The concept of *unlinkability* of two or more items is defined as "within this system, these items are no more and no less related than they are related concerning the a-priori knowledge" [18]. *Connection anonymity* [5] is about hiding the identities of source and destination during the actual data transfer, so that connections between individuals remain private (in addition to the actual data). *Data anonymity* [5] is about filtering any identifying information out of the data that is exchanged in a particular application.

*Pseudonymity* is the use of pseudonyms as IDs (subjects) [18]. There are many examples, where pseudonymous data has been de-anonymized, when such data has been analyzed more thoroughly, e.g., [8, 10, 16, 17].

*Unobservability* is the state of items of interest being indistinguishable from any items of interest at all [18]. In *identity disclosure*, identity of an individual is associated with a record containing confidential information in the released dataset, whereas in *attribute disclosure* an attribute value is associated with an individual, not necessary linked to a specific record [2].

## 3 MEASURES OF PRIVACY AND ANONYMITY

*Differential privacy* [6]: the model ensures that small changes (one record) in the original database has a limited impact on the outcome of any statistical analysis on the data.

*Bayes-optimal privacy* [14] refers to the situation in which the adversary and the data holder both have all available background information.

*k-anonymity* [21] protection is provided in the release, if the information for each person contained in the release cannot be distinguished from at least $k - 1$ individuals whose information

also appears in the release. The problem with *k*-anonymity is that it cannot prevent attribute disclosure.

*Effective anonymity set size* [19] uses entropy as a measure for describing their information theoretic approach to privacy metrics.

The *degree of anonymity* is expressed and quantified in various ways in literature. In [5], the information theoretic approach to the degree of anonymity is described as the measure the information the attacker gets, taking into account the whole set of users and the probabilistic information the attacker obtains about them.

*l-diversity* [14]: a block is *l*-diverse if it contains at least *l* "well-represented" values for the sensitive attribute S. A weakness in *l*-diversity is the limited assumption of adversarial knowledge [12].

*t-closeness* [12] formalizes the idea of global background knowledge by requiring that the distribution of a sensitive attribute in any equivalence class is close to the distribution of the attribute in the overall table. This model has some weaknesses, related for example in the flexibility of specifying different protection levels for different sensitive values [7].

*Entropy* as an information metric [20] provides a measure of the uncertainty of a random variable. This can be seen as a high level metric for evaluating privacy related information. The notion of *one-symbol information* [2] (i.e., the contribution to mutual information by a single record) is to "allow to express and compare the disclosure risk metrics."

## 4 PREVIOUS WORK

*Security metrics* have been studied quite extensively and matured in some fields such as software vulnerabilities with the Common Vulnerability Scoring System (CVSS). Privacy metrics however have been centered around anonymity, or how much entropy a given PET yields in different settings. Such measures lack the nuance that privacy in general holds.

Various anonymity metrics are discussed in [4]. They define anonymity metric as to "indicate how uncertain, from the attacker's point of view, the identity of a fixed entity is." According to [13], privacy is an individual concept and should be measured separately for each individual.

A method called "Privacy-preserving data publishing" is presented in [7], in which recipients of the data are unknown in advance to the data publisher. It classifies recent privacy models into two categories according to their attack principles.

In [9] a framework to assess the function and properties of various PETs is presented, but some of the measures are somewhat prone to subjective interpretation. Also, the model does not offer information of how much the PET in question improves privacy.

In [11] a privacy measuring framework is provided with finite state machines. Privacy is divided into nine distinct states and then different transitions between these states, in which the state of privacy can be derived from a finite state of factors for a given situation. The division of privacy is intriguing, but fairly simplified.

In [1] the authors present a method for calculating the exposure of data items in social networks. Their work shares some similarities with our approach, but is focused on social networks.

It is interesting to note that many times privacy is being measured through surveys and questionnaires [15] and this can lead to some biased results [3].

## 5 PROPERTIES OF A GOOD PRIVACY METRIC

In order to build a good metric for privacy, we need to evaluate the properties that the metric should possess, and make clear distinction between anonymity and privacy. In our view, privacy is a more subjective feature of an individual, whereas anonymity considers sets of subjects with similar level of privacy.

Like anonymity, privacy can be presented as unlinkability of some information to a specific person, but other quantities need to be considered as well.

The expectation of privacy and the level of privacy protection offered can depend on the context of *location*. For example, one might use a screen overlay over a laptop in public, but not in the office, even though one expects only the same level of privacy in all cases. *The temporal factor* is also significant. For example, having your name published as an applicant to some job might be sensitive before the choice is made, but insignificant when you are chosen.

Another example of context-dependency of privacy is the *amount of intent (and potential) to cause harm* to the data subject, not only what is revealed but to whom. *Connectivity of data* is another dependency: how can one piece of information be connected with other information to erode privacy.

Traditionally, privacy metrics concentrating on anonymity, measure the "amount of privacy" in bits as in entropy. This usually works with a static setting, where privacy is measured with respect to some data set and some set of people. In a dynamic environment, this kind of "snapshot" measure might not be adequate.

Privacy can be measured as exposure, indicating how much and possibly what kind of sensitive information about the individual is available to some attackers. This can be a probabilistic metric, measuring the probability of a given piece of information becoming public and/or linked to the individual, such as in [5]. Given enough information, such a metric can give an expected impact of different breaches to privacy and even an expected value of losses.

The problem with this approach is that many of the valuations and variables are subjective. Also, a single metric might not be able to incorporate all the variables that are required for comprising a comprehensive level of privacy for an individual. Thus, there is a need to balance between a simple and effective metric, that is not necessarily accurate or applicable in all cases, and a more complex and comprehensive metric, that will require the individual using the metric to give valuations etc. in order to measure privacy.

We conclude with some *ideal* properties of a privacy metric:

- Comprehensive - cover all facets of privacy including user privacy and data privacy
- Objective - independent of the interpreter of the measurement
- Unambiguous - again, the end result independent of the interpreter
- Comparable - presented in an uniform and homogeneous measure units
- Quantifiable - presented in measures that can be associated to numbers
- Individual - measured separately for an individual
- Probabilistic - taking into account the uncertainty of information gathered about a subject

- Computable - the metric is expressed in formulas that can be effectively computed

## 6  THE VALUE OF PRIVACY

In order to measure privacy we look into probabilities and expected values in certain distributions. We assume that information about individuals in the system is available as *facts*. A fact is a single piece of correct information that relates to an individual, e.g., Mark is 6'4" tall. Facts can have differing levels of detail and each fact has an associated weight vector, which contains the privacy value of the fact to the subject itself and all possible observers in the system. Facts can also be links between other facts, e.g., "x and y are related".

We assume that the privacy is measured with respect to an individual, called the *subject* of privacy. In our system there may be any number of *observers* that can observe information related to the subject. Any number of observers can collude and combine their observations on any given subject. This can be still viewed as a single observer (with more information) on the subject.

The universe $\mathcal{U}$ of possible events is all the information and the links between different pieces of information and different subjects. The universe of events related to subject $s$ is $\mathcal{U}_s \subseteq \mathcal{U}$ and it contains all the information and links between different pieces of information and other subjects that consider the subject $s$.

*The effort* used by an observer $o$ is a function that determines how much effort an observer is using to learn new information about a subject. This function $f$ is dependent on the observer $o$, the subject $s$ and the fact $x$. We denote it by $f_{s,o}(x)$ for a single fact $x$.

The probability of a single fact $x \in \mathcal{U}_s$ of subject $s$ being revealed to an observer $o$ using effort $f$ is denoted as $P_{s,o}^f(x)$. If $f(s,o) = 0$, we define $P_{s,o}^f(x) = 0$ for all $x \in \mathcal{U}_s$. This means that zero effort yields no (new) information about a subject to any observer. This implies that even when a disclosure of some information is made by accident, some effort is required by the observer to remember this fact and to link it to its "database" of facts.

If we assume independent observers, the probability of a fact $x$ about a subject $s$ being disclosed to *any* observer is $\sum_{o \in O} P_{s,o}^f(x)$, where $O$ is the set of all observers in the system. This may seem a great simplification as observers can be highly linked. However, we can combine several observers into a single observer (with possibly some loss in accuracy) and thus obtain either completely independent observers or at least less connected set of observers.

The *effort threshold* $t_{s,o}(x)$ of learning a fact $x$ about subject $s$ is the minimum value of $f_{s,o}(x)$ for which $P_{s,o}^f(x) = 1$. The *cost* for observer $o$ to learn a set of facts $X$ about $s$ is $\sum_{x \in X} t_{s,o}(x)$.

We define the *conditional privacy* of fact $x$ given a set of facts $Y$ as the probability of the observer knowing the facts in set $Y$ in learning the fact $x$ and denote it by $P_{s,o}^f(x|Y)$. This notion helps in modeling more complex scenarios related to privacy. Note, that the set $Y$ can contain also facts not directly about $s$, but also about other subjects that are related to $s$. This type of inference is very common in, e.g., social media networks and related advertising.

Let $x$ and $y$ be two facts about a subject $s$. We say that $x$ is a *subfact* of $y$ if $P_{s,o}^f(x|y) = 1$ for all $f$. This is denoted by $x \prec y$. An example of this is that the fact "Sue lives in Albuquerque" has

"Sue lives in the USA" and "Sue lives on Earth" as subfacts. We also define *prime facts* as facts $x$ for which $x' \not\prec x$ for all $x' \in \mathcal{U}_s, x' \neq x$. We denote the set of prime facts of subject $s$ as $\mathcal{R}_s$.

Loss of privacy can be measured in several ways with our model. The simplest measure is the amount of information that is available to observers, e.g., the amount of public information as the amount measure of lost privacy. This type of measure is easy to construct, but fails to take into account the great variety that privacy constitutes. Not all facts are equal and not all observers are equal from the subject's point of view.

By assigning privacy values to each fact of a subject, we can calculate the expected loss of privacy due to an observer $o$ to subject $s$ as the expected value of $P_{s,o}^f(x)$ for all $x \in \mathcal{U}_s$. However, this direct calculation overestimates the loss of privacy as arguably any loss of privacy due to a subfact $y$ has already been accounted in the loss of the prime fact $p$ for which $y \prec p$.

Thus we have a formula for the expected loss of privacy (ELP) due to an observer $o$ $\mathbb{E}(P_{s,o}^f(X)) := \text{ELP}_o$. Again assuming independent observers, we can compute the total expected loss of privacy (TELP) for subject $s$ as a sum over all observers $\sum_{o \in O} \mathbb{E}(P_{s,o}^f(X))$, where $O$ is the set of all observers in the system.

Privacy is many times a trade-off between a perceived value of possibly privacy infringing transaction and keeping this information private. Choosing to keep some facts private can lead to false impression of an individual for example through inference from the social network. This can lead into situations, where a subject is better off revealing some fact about himself, as the false impression of not having the fact revealed can have greater impact on the subject.

An observer has an *impression* on subject $s$ denoted by $I_{s,o}$, which contains all the facts about $s$ known to $o$. The set of *false impressions* $F_{s,o}$ of $o$ on $s$ is $F_{s,o} = \{i \in I_{s,o} : i \notin \mathcal{U}_s\}$. That is, an observer can have wrong information about a subject (we will still call these bits of information facts although they are not accurate).

The *public information* about a subject $s$ (denoted by $\text{Pub}_s$) is $\text{Pub}_s = \bigcap_{o \in O} I_{s,o}$ Thus, $\text{Pub}_s$ is the set of facts that all observers know about $s$. Of course, this is a very limiting definition as having only single observer, without knowledge of a fact $x$ makes the fact *not* public information.

It is worthwhile to note that the public information contains also possible false impressions of the observers. *Private information* of $s$ is defined as $\mathcal{U}_s \setminus \bigcup_{o \in O} I_{s,o}$. This definition is again very narrow as having a single observer that has knowledge of a fact $x$ excludes this fact from private information. Thus, most of some subject's facts fall between these two extremes of public and private information. We define the *spread* of a fact $x$ about a subject $s$ as the number of observers with $x \in I_{s,o}$ and we denote it by $\langle x \rangle_s$. For a set of facts $X$ we will use the same notation $\langle X \rangle_s$ with $X \subset \mathcal{U}_s$.

Our privacy measures can be utilised to measure the success of different PETs in maintaining or improving the privacy of their users. The main goal of PETs is to improve the privacy of the user that is utilising them in some context. From the above definitions, we can see several different measurable quantities that can be used as metrics of efficiency for PETs.

First of all, is the spread of a fact. If a PET can be used to assure that a fact $x$ will only have a spread under a certain threshold,

this limit can be used as a metric of the efficiency of this PET. For example, encryption between two parties should ensure that the when information about the fact $x$ is transmitted from Alice to Bob, only these two entities learn $x$. Thus the spread of $x$ is at most 2, i.e., $\langle x \rangle \leq 2$. This type of measure can be seen as a measure of the *containment* of facts.

Another measure for the effectiveness of a PET is the increase in the effort for some or all observers to learn some fact $x$ about the subject. This increase can be either absolute or relative to the original cost of learning $x$. This measures how well the PET *conceals* the facts that it is used to protect.

More subjective measures for the effectiveness of PETs are the amount of increase in the (probabilistic) false impression of (some) observers and the amount of decrease in the (probabilistic) true impression of (some) observers about the subject. These measures can be used to some extent to measure PETs. However, it is highly subjective, whether for example some false impression is beneficial to a subject $s$. Thus, to use this as a measure one needs to consider also the value of the change in impression to the subject. This is especially hard, if we consider cases where some impression is given only to some observers and not all observers. This measures how well the PET *obfuscates* the facts of the subject among observers.

## 7 DISCUSSION

Our model gives tools to evaluate privacy in many cases. The challenge is to find accurate and easy-to-evaluate values for some of the parameters in our model. Thus, first instantiations could be made with some simplifications and crude measurements. The model itself can be evaluated and improved, when concrete examples can be found.

In general, it is challenging to comprehensively model privacy in mathematical terms, such as the nine classes of privacy from [11]. Precise definitions need strong assumptions (such as independent observers) and may lead to overly simplified models.

Even with our proposed methods it is hard to find the exact values for measuring privacy as a single value for an individual. It is much easier to measure technologies and their possible effectiveness in protecting privacy than quantify privacy, but this reveals only one side of the issue.

The model we have developed can be extended by modeling colluding observers that exchange information about subjects. These can be used to improve the confidence of observers by comparing impressions and having for example majority vote on a fact belonging to the true impression. By including the confidence values of each observers as weights in the majority voting, the model can be expanded even further.

Our model achieves many of the properties of ideal measure, although some properties such as objectivity and computability are not necessarily complete. An interesting piece of further work is to evaluate different privacy metrics against the ideal features in order to compare, which come closest to fulfilling these. Such evaluation could reveal the gaps in privacy metrics development.

## 8 CONCLUSION

We present a new method for measuring the value of privacy for individuals and for evaluating the efficiency of privacy enhancing technologies. Our measure is centered on a subject and the value of different data points (facts) about this subject. At this current stage, the measures provide tools at an abstract level.

Future work should refine this method and seek to find the limits and applications of this approach as well as build a proper taxonomy, against which some appropriate targets could be validated. Analyzing case studies with concrete values and valuations would show the power, and also possible weaknesses of this model. This would help in further developing this method towards a usable and applicable metric in real world scenarios.

## REFERENCES

[1] S. Ananthula, O. Abuzaghleh, N. B. Alla, S. Chaganti, P. Kaja, and D. Mogilineedi. Measuring privacy in online social networks. *International Journal of Security, Privacy and Trust Management*, 4(2):1–9, 2015.

[2] M. Bezzi. An information theoretic approach for privacy metrics. *Trans. Data Privacy*, 3(3):199–215, 2010.

[3] A. Braunstein, L. Granka, and J. Staddon. Indirect content privacy surveys: measuring privacy without asking about it. In *Proceedings of the Seventh Symposium on Usable Privacy and Security*, page 15. ACM, 2011.

[4] S. Claußand S. Schiffner. Structuring anonymity metrics. In *Proceedings of the Second ACM Workshop on Digital Identity Management*, DIM '06, pages 55–62, New York, NY, USA, 2006. ACM.

[5] C. Diaz, S. Seys, J. Claessens, and B. Preneel. Towards measuring anonymity. In *International Workshop on Privacy Enhancing Technologies*, pages 54–68. Springer, 2002.

[6] C. Dwork. *Differential Privacy: A Survey of Results*, pages 1–19. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.

[7] B. C. M. Fung, K. Wang, R. Chen, and P. S. Yu. Privacy-preserving data publishing: A survey of recent developments. *ACM Comput. Surv.*, 42(4):14:1–14:53, June 2010.

[8] S. Gambs, M.-O. Killijian, and M. N. del Prado Cortez. De-anonymization attack on geolocated data. *Journal of Computer and System Sciences*, 80(8):1597–1614, 2014.

[9] J. Heurix, P. Zimmermann, T. Neubauer, and S. Fenz. A taxonomy for privacy enhancing technologies. *Computers & Security*, 53:1–17, 2015.

[10] H. Kataoka, Y. Ogawa, I. Echizen, T. Kuboyama, and H. Yoshiura. Effects of external information on anonymity and role of transparency with example of social network de-anonymisation. In *Availability, Reliability and Security (ARES), 2014 Ninth International Conference on*, pages 461–467. IEEE, 2014.

[11] T. A. Kosa, K. EI-Khatib, and S. Marsh. Measuring privacy. *Journal of Internet Services and Information Security (JISIS)*, 1(4):60–73, 2011.

[12] N. Li, T. Li, and S. Venkatasubramanian. t-closeness: Privacy beyond k-anonymity and l-diversity. In *2007 IEEE 23rd International Conference on Data Engineering*, pages 106–115, April 2007.

[13] T. Li and N. Li. On the tradeoff between privacy and utility in data publishing. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '09, pages 517–526, New York, NY, USA, 2009. ACM.

[14] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam. l-diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):3, 2007.

[15] K. Martin and H. Nissenbaum. Measuring privacy: an empirical test using context to expose confounding variables. *Colum. Sci. & Tech. L. Rev.*, 18:176, 2016.

[16] A. Narayanan, E. Shi, and B. I. Rubinstein. Link prediction by de-anonymization: How we won the kaggle social network challenge. In *Neural Networks (IJCNN), The 2011 International Joint Conference on*, pages 1825–1834. IEEE, 2011.

[17] A. Narayanan and V. Shmatikov. Robust de-anonymization of large sparse datasets. In *Security and Privacy, 2008. SP 2008. IEEE Symposium on*, pages 111–125. IEEE, 2008.

[18] A. Pfitzmann and M. Köhntopp. Anonymity, unobservability, and pseudonymity—a proposal for terminology. In *Designing privacy enhancing technologies*, pages 1–9. Springer, 2001.

[19] A. Serjantov and G. Danezis. Towards an information theoretic metric for anonymity. In *Privacy Enhancing Technologies*, pages 41–53. Springer, 2002.

[20] C. E. Shannon. A mathematical theory of communication. *ACM SIGMOBILE Mobile Computing and Communications Review*, 5(1):3–55, 2001.

[21] L. Sweeney. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05):557–570, 2002.

[22] THE EUROPEAN PARLIAMENT AND THE COUNCIL OF THE EUROPEAN UNION. Directive (eu) 2016/680 of the european parliament and of the council. *Official Journal of the European Union*, 2016.