

Topics in Information-Theoretic Cryptography

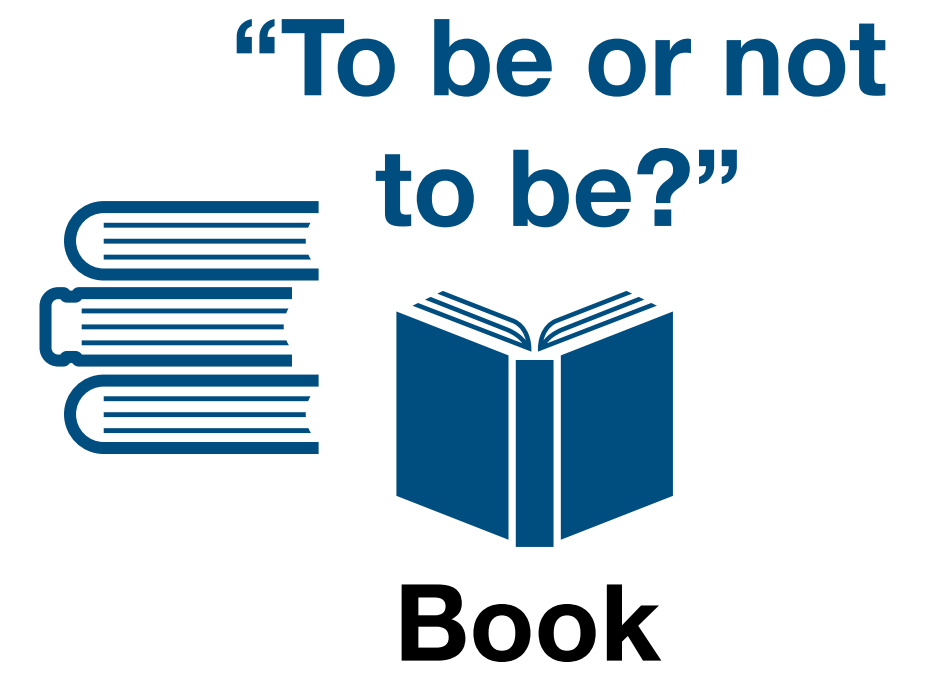
Course Wrap Up

Yanina Shkel, December 9, 2021

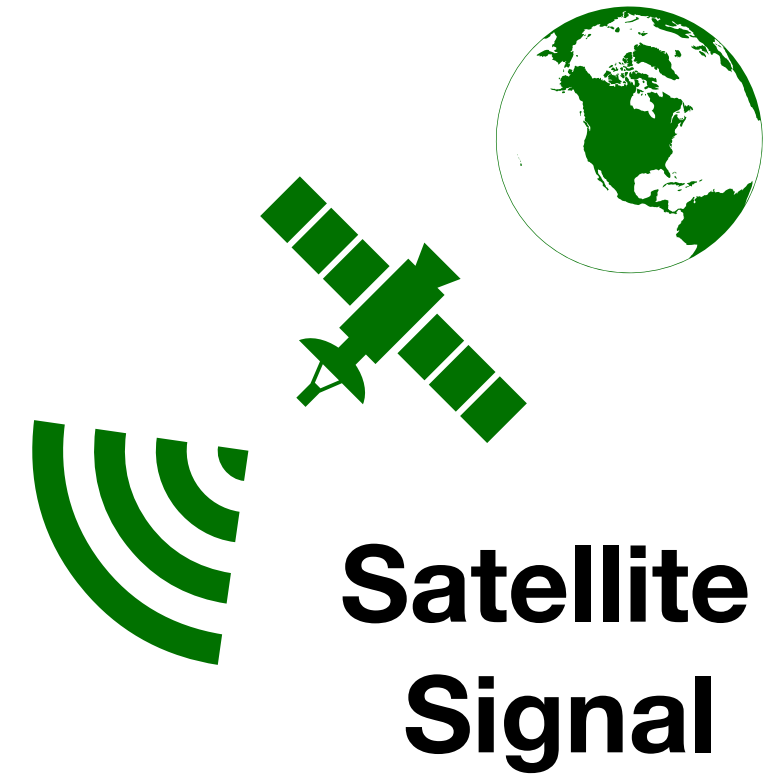
How do we model information?

**How do we measure
information?**

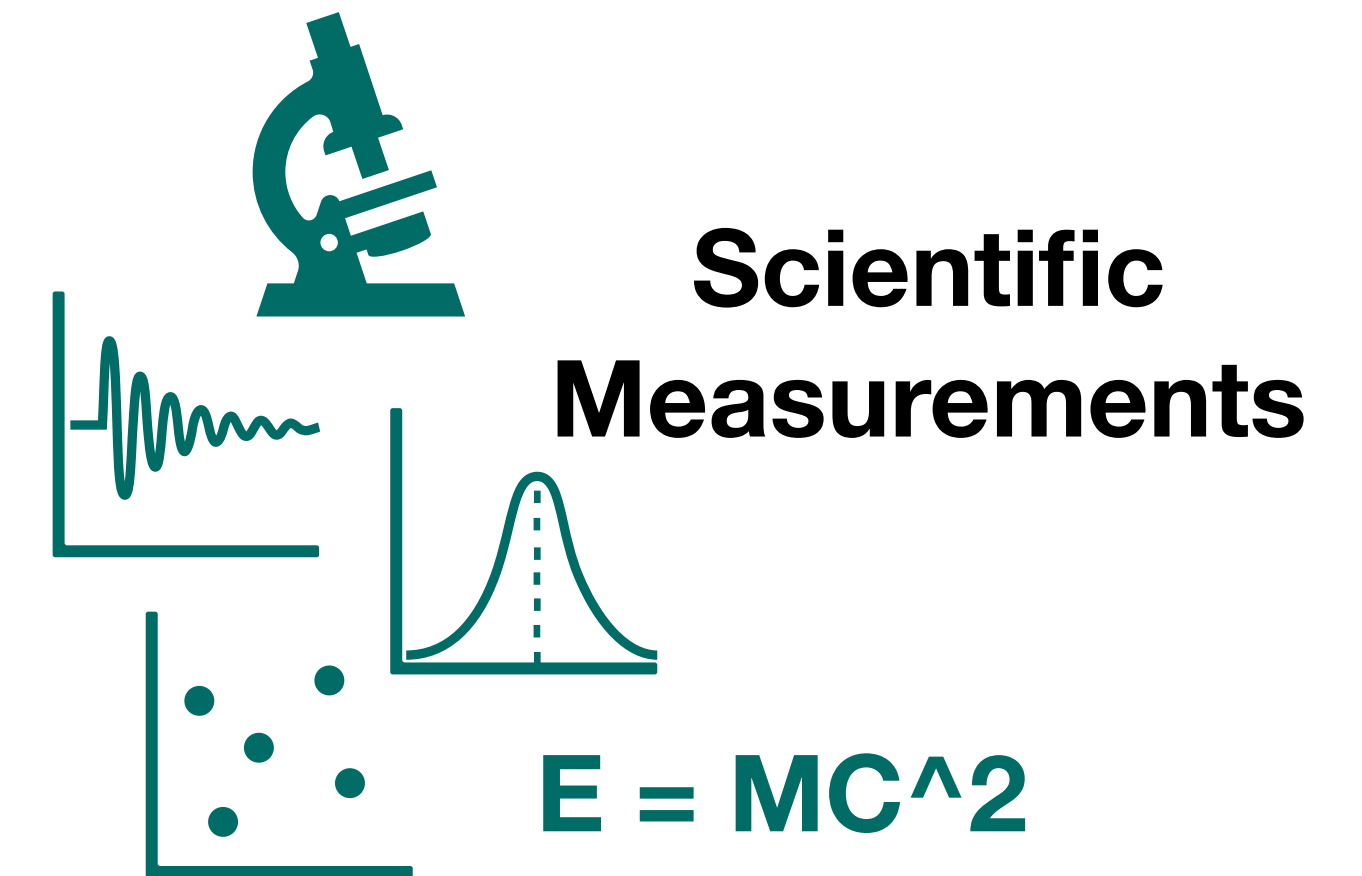
How much information is there?



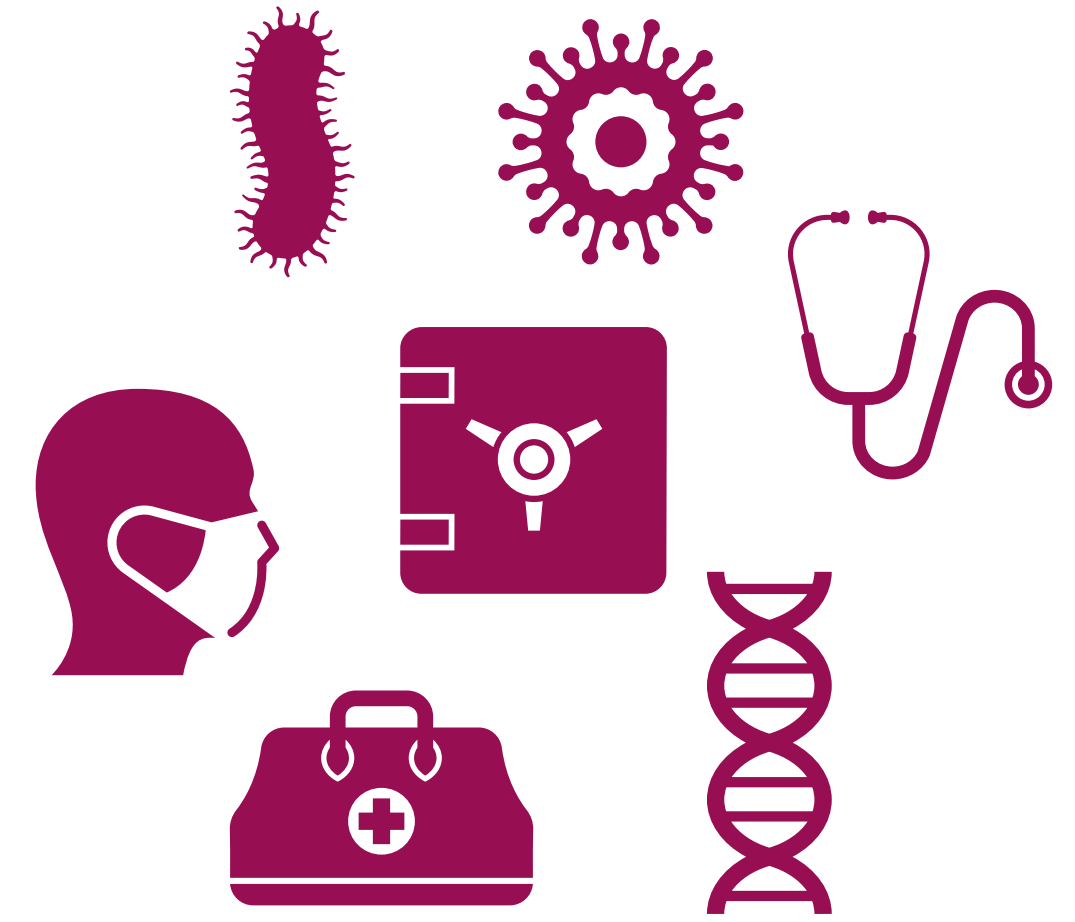
How much information is there?



How much information is there?



How much information is there?



Medical Data

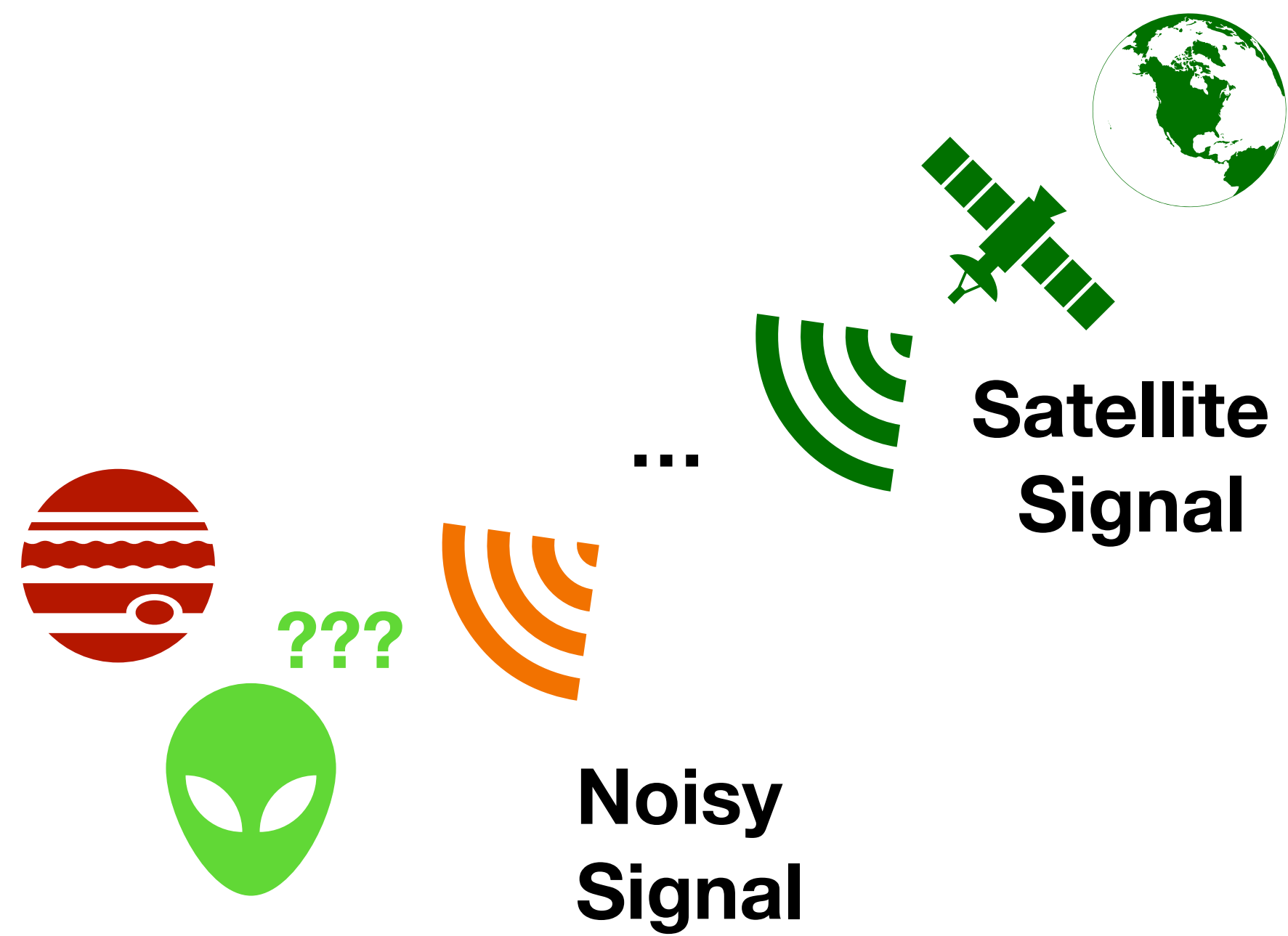
How much information is there?

#H*kehtfw20e;sp0s0gsl

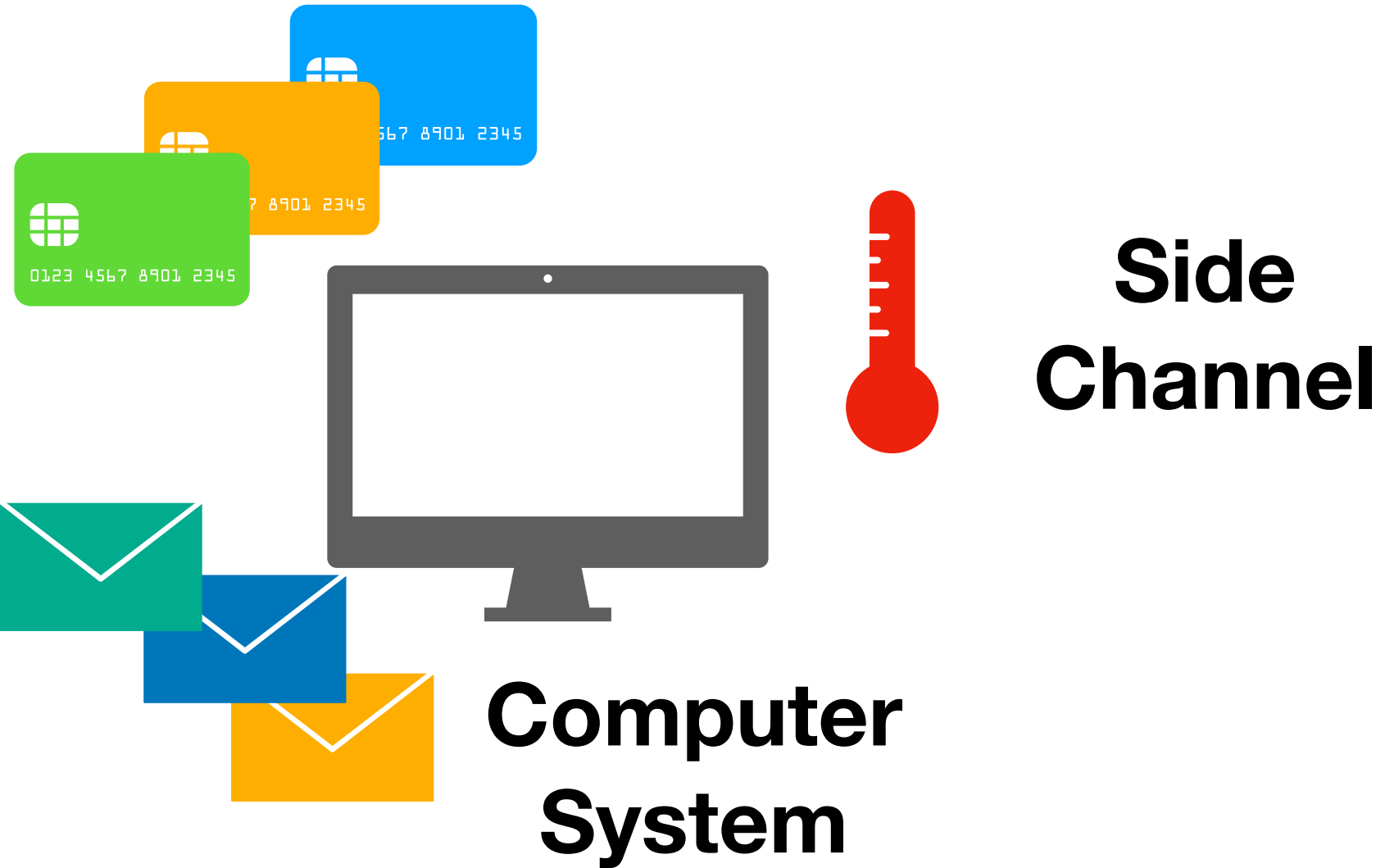


Random Noise

How much relevant information is there in a related observation?



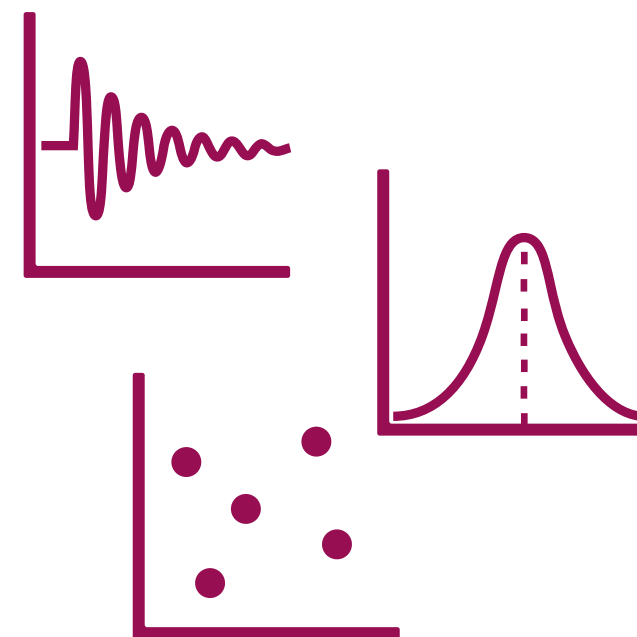
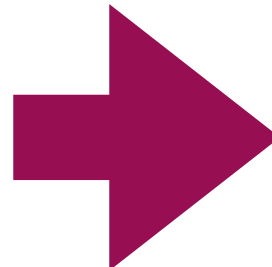
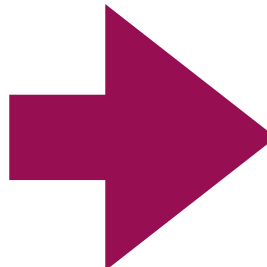
How much relevant information is there in a related observation?



How much relevant information is there in a related observation?



Medical Data



Data Statistics

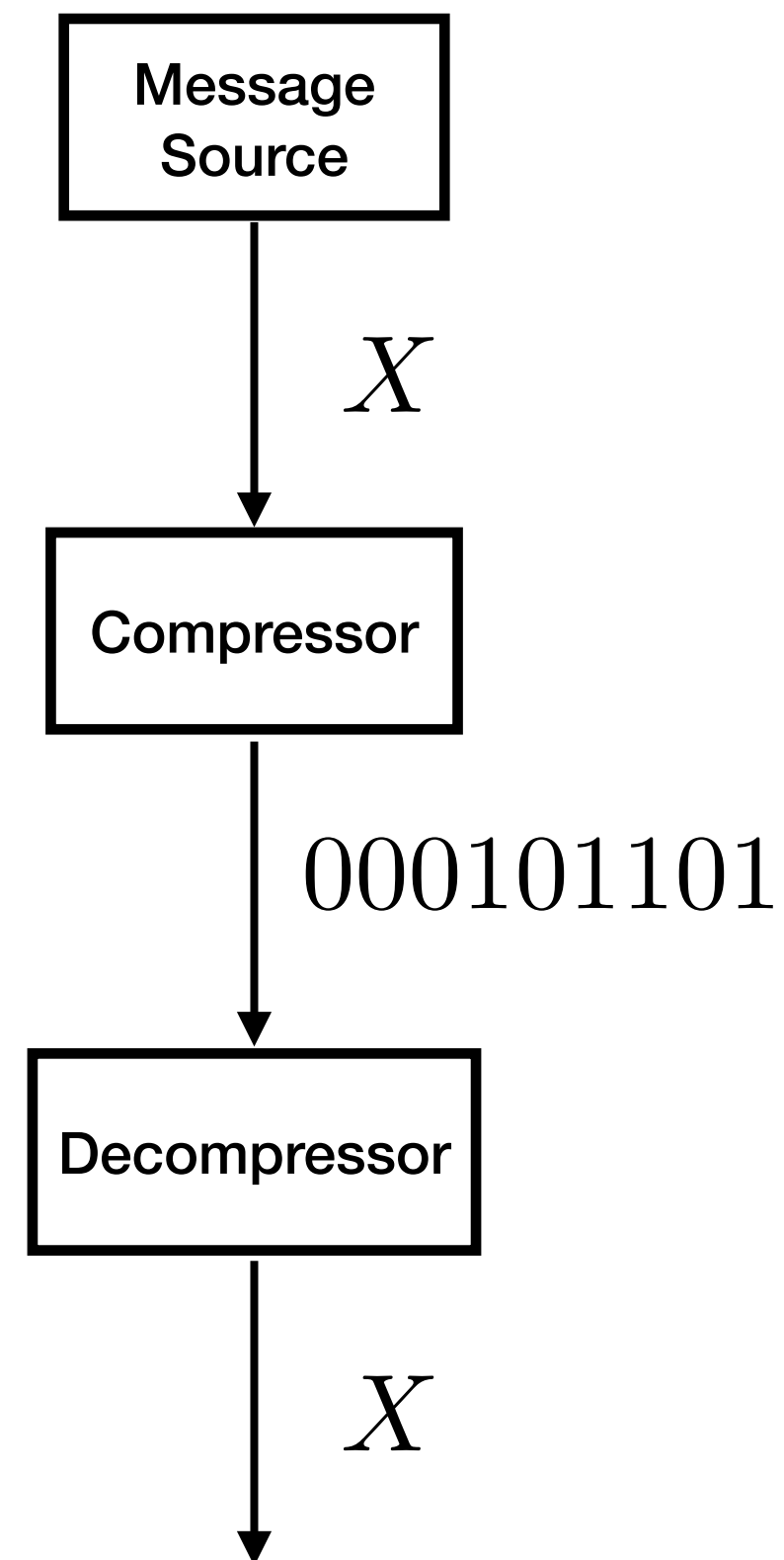
What works well...

Use tools from probability theory and statistics

- A ‘source of information’ is a random variable
 - e.g. X , Y , Z , ...
- ‘Information measures’ quantify the amount of information
 - e.g. mutual information, Shannon entropy, relative entropy,...
- Justifications for information measures:
 - pops up as an answer to an engineering problem
 - satisfies some nice properties (axioms, operational notions, etc.)

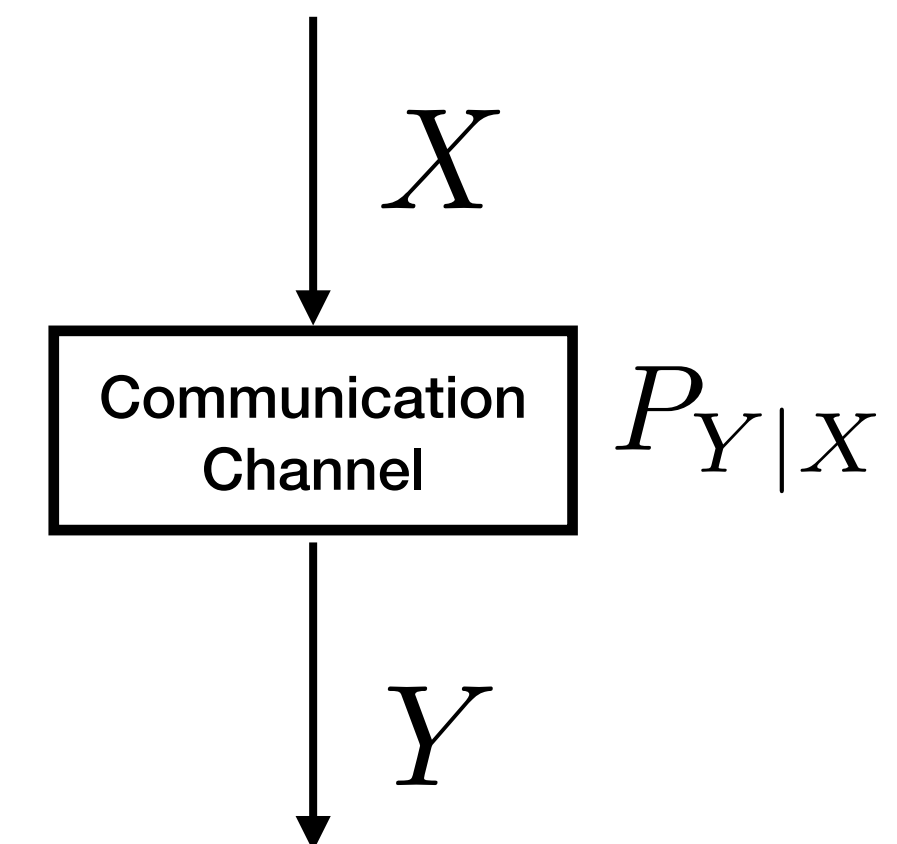
Communications

Example: Data Compression



$$H(X) = \sum_{x \in \mathcal{X}} P_X(x) \log \frac{1}{P_X(x)}$$

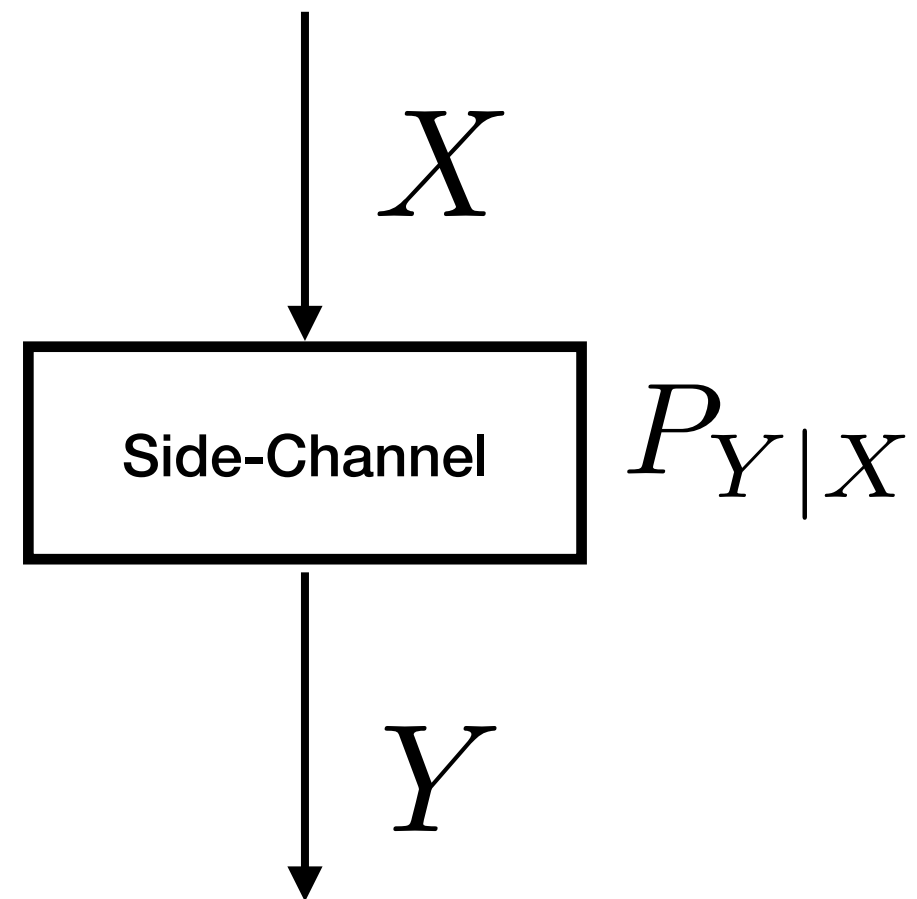
Example: Data Transmission



$$C = \max_{P_X} I(X; Y)$$

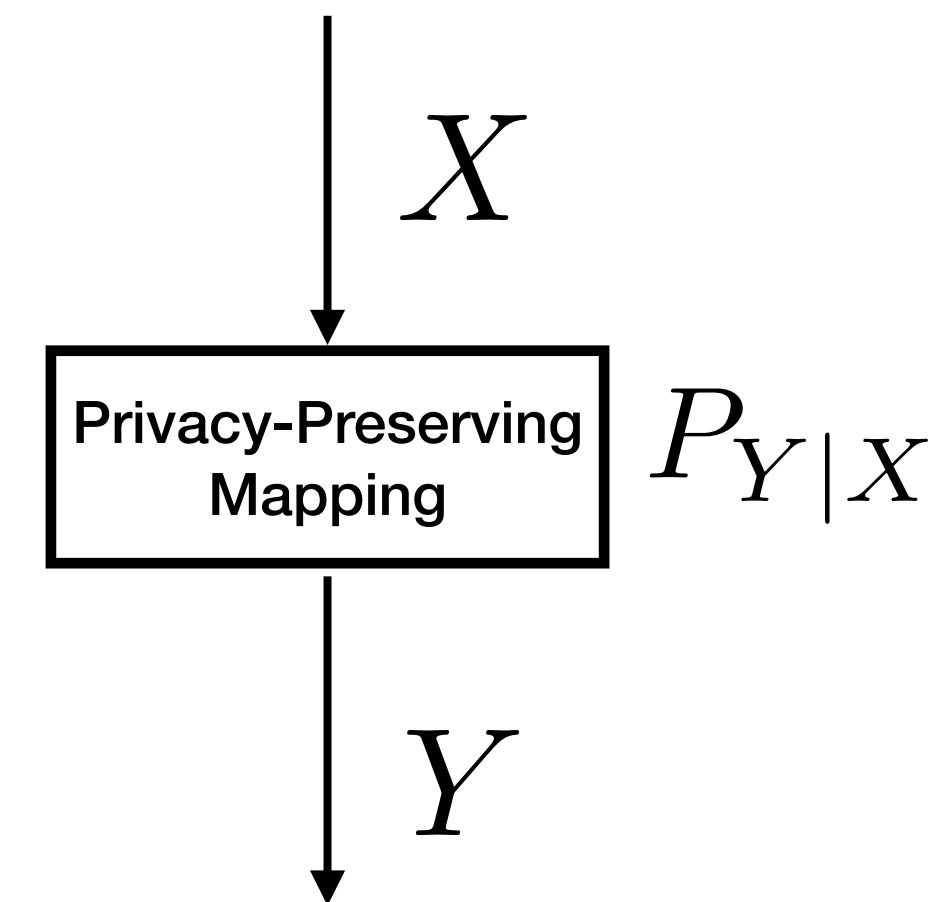
Privacy and Security

Example: Side-Channel Leakage



$$\mathcal{L}(X \rightarrow Y) = \sup_{U-X-Y-\hat{U}} \log \frac{\mathbb{P}[U = \hat{U}]}{\max_{u \in \mathcal{U}} P_U(u)}$$

Example: Database Privacy



$$\mathcal{DP}(X \rightarrow Y) = \sup_{x_1, x_2 \in \mathcal{X} : d(x_1, x_2) = 1, \mathcal{T} \subset \mathcal{Y}} \log \frac{\mathbb{P}[Y \in \mathcal{T} | X = x_1]}{\mathbb{P}[Y \in \mathcal{T} | X = x_2]}$$

COURSE REVIEW

A Mathematical Theory of Communication and Beyond (Information Theory)

- Has been extremely successful in addressing problems like communication and data compression
- Concerned with measures like Shannon entropy and mutual information
 - these measures arise as answers to specific engineering problems
- These measures have been used in other applications with much more mixed results

Reprinted with corrections from *The Bell System Technical Journal*, Vol. 27, pp. 379–423, 623–656, July, October, 1948.

A Mathematical Theory of Communication

By C. E. SHANNON

INTRODUCTION

THE recent development of various methods of modulation such as PCM and PPM which exchange bandwidth for signal-to-noise ratio has intensified the interest in a general theory of communication. A basis for such a theory is contained in the important papers of Nyquist¹ and Hartley² on this subject. In the present paper we will extend the theory to include a number of new factors, in particular the effect of noise in the channel, and the savings possible due to the statistical structure of the original message and due to the nature of the final destination of the information.

The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point. Frequently the messages have *meaning*; that is they refer to or are correlated according to some system with certain physical or conceptual entities. These semantic aspects of communication are irrelevant to the engineering problem. The significant aspect is that the actual message is one *selected from a set* of possible messages. The system must be designed to operate for each possible selection, not just the one which will actually be chosen since this is unknown at the time of design.

If the number of messages in the set is finite then this number or any monotonic function of this number can be regarded as a measure of the information produced when one message is chosen from the set, all choices being equally likely. As was pointed out by Hartley the most natural choice is the logarithmic function. Although this definition must be generalized considerably when we consider the influence of the statistics of the message and when we have a continuous range of messages, we will in all cases use an essentially logarithmic measure.

The logarithmic measure is more convenient for various reasons:

1. It is practically more useful. Parameters of engineering importance such as time, bandwidth, number of relays, etc., tend to vary linearly with the logarithm of the number of possibilities. For example, adding one relay to a group doubles the number of possible states of the relays. It adds 1 to the base 2 logarithm of this number. Doubling the time roughly squares the number of possible messages, or doubles the logarithm, etc.
2. It is nearer to our intuitive feeling as to the proper measure. This is closely related to (1) since we intuitively measure entities by linear comparison with common standards. One feels, for example, that two punched cards should have twice the capacity of one for information storage, and two identical channels twice the capacity of one for transmitting information.
3. It is mathematically more suitable. Many of the limiting operations are simple in terms of the logarithm but would require clumsy restatement in terms of the number of possibilities.

The choice of a logarithmic base corresponds to the choice of a unit for measuring information. If the base 2 is used the resulting units may be called binary digits, or more briefly *bits*, a word suggested by J. W. Tukey. A device with two stable positions, such as a relay or a flip-flop circuit, can store one bit of information. N such devices can store N bits, since the total number of possible states is 2^N and $\log_2 2^N = N$. If the base 10 is used the units may be called decimal digits. Since

$$\begin{aligned}\log_2 M &= \log_{10} M / \log_{10} 2 \\ &= 3.32 \log_{10} M,\end{aligned}$$

¹Nyquist, H., "Certain Factors Affecting Telegraph Speed," *Bell System Technical Journal*, April 1924, p. 324; "Certain Topics in Telegraph Transmission Theory," *A.I.E.E. Trans.*, v. 47, April 1928, p. 617.

²Hartley, R. V. L., "Transmission of Information," *Bell System Technical Journal*, July 1928, p. 535.

On Measures of Entropy and Information

- Extends Shannon entropy and relative entropy to a family of **Rényi entropies** and **Rényi divergences**
- Takes an axiomatic view of entropy
- Rényi entropy is additive across independent observations
 - does not satisfy the chain rule

ON MEASURES OF ENTROPY AND INFORMATION

ALFRED RÉNYI
MATHEMATICAL INSTITUTE
HUNGARIAN ACADEMY OF SCIENCES

1. Characterization of Shannon's measure of entropy

Let $\mathcal{P} = (p_1, p_2, \dots, p_n)$ be a finite discrete probability distribution, that is, suppose $p_k \geq 0 (k = 1, 2, \dots, n)$ and $\sum_{k=1}^n p_k = 1$. The amount of uncertainty of the distribution \mathcal{P} , that is, the amount of uncertainty concerning the outcome of an experiment, the possible results of which have the probabilities p_1, p_2, \dots, p_n , is called the *entropy* of the distribution \mathcal{P} and is usually measured by the quantity $H[\mathcal{P}] = H(p_1, p_2, \dots, p_n)$, introduced by Shannon [1] and defined by

$$(1.1) \quad H(p_1, p_2, \dots, p_n) = \sum_{k=1}^n p_k \log_2 \frac{1}{p_k}.$$

Different sets of postulates have been given, which characterize the quantity (1.1). The simplest such set of postulates is that given by Fadeev [2] (see also Feinstein [3]). Fadeev's postulates are as follows.

- (a) $H(p_1, p_2, \dots, p_n)$ is a symmetric function of its variables for $n = 2, 3, \dots$.
- (b) $H(p, 1 - p)$ is a continuous function of p for $0 \leq p \leq 1$.
- (c) $H(1/2, 1/2) = 1$.
- (d) $H(tp_1, (1 - t)p_1, p_2, \dots, p_n) = H(p_1, p_2, \dots, p_n) + p_1 H(t, 1 - t)$

for any distribution $\mathcal{P} = (p_1, p_2, \dots, p_n)$ and for $0 \leq t \leq 1$.

The proof that the postulates (a), (b), (c), and (d) characterize the quantity (1.1) uniquely is easy except for the following lemma, whose proofs up to now are rather intricate.

LEMMA. Let $f(n)$ be an additive number-theoretical function, that is, let $f(n)$ be defined for $n = 1, 2, \dots$ and suppose

$$(1.2) \quad f(nm) = f(n) + f(m), \quad n, m = 1, 2, \dots$$

Let us suppose further that

$$(1.3) \quad \lim_{n \rightarrow +\infty} [f(n+1) - f(n)] = 0.$$

Then we have

$$(1.4) \quad f(n) = c \log n,$$

where c is a constant.

Paper Highlights

Rényi Entropy

THEOREM 2. *If $H[\mathcal{P}]$ is defined for all $\mathcal{P} \in \Delta$ and satisfies postulates 1, 2, 3, 4, and 5' with $g(x) = g_\alpha(x)$, where $g_\alpha(x)$ is defined by (2.13), $\alpha > 0$, and $\alpha \neq 1$, then $H[\mathcal{P}] = H_\alpha[\mathcal{P}]$, where, putting $\mathcal{P} = (p_1, p_2, \dots, p_n)$, we have*

$$(2.14) \quad H_\alpha[\mathcal{P}] = \frac{1}{1 - \alpha} \log_2 \left[\frac{\sum_{k=1}^n p_k^\alpha}{\sum_{k=1}^n p_k} \right].$$

The quantity (2.14) will be called the *entropy of order α* of the generalized distribution \mathcal{P} . Clearly if \mathcal{P} is an ordinary distribution, (2.14) reduces to (1.21). It is also easily seen that

$$(2.15) \quad \lim_{\alpha \rightarrow 1} H_\alpha[\mathcal{P}] = H_1[\mathcal{P}],$$

where $H_1[\mathcal{P}]$ is defined by (2.5).

Rényi Divergence

THEOREM 3. *Suppose that the quantity $I(\mathcal{Q}|\mathcal{P})$ satisfies the postulates 6, 7, 8, 9, and 10. Then the function $g(x)$ in 10 is necessarily either a linear or an exponential function. In the first case $I(\mathcal{Q}|\mathcal{P}) = I_1(\mathcal{Q}|\mathcal{P})$, where*

$$(3.7) \quad I_1(\mathcal{Q}|\mathcal{P}) = \frac{\sum_{k=1}^n q_k \log_2 \frac{q_k}{p_k}}{\sum_{k=1}^n q_k},$$

while in the second case $I(\mathcal{Q}|\mathcal{P}) = I_\alpha(\mathcal{Q}|\mathcal{P})$ with some $\alpha \neq 1$, where

$$(3.8) \quad I_\alpha(\mathcal{Q}|\mathcal{P}) = \frac{1}{\alpha - 1} \log_2 \frac{\sum_{k=1}^n \frac{q_k^\alpha}{p_k^{\alpha-1}}}{\sum_{k=1}^n q_k}.$$

An Operational Approach to Information Leakage

- An information leakage measure motivated by ‘side channels’
- Measures adversary’s improvement in ability to estimate a function of data
- Has many nice properties like composition, data processing, and various robustness properties

An Operational Approach to Information Leakage

Ibrahim Issa¹, Aaron B. Wagner², and Sudeep Kamath³

Abstract—Given two random variables X and Y , an operational approach is undertaken to quantify the “leakage” of information from X to Y . The resulting measure $\mathcal{L}(X \rightarrow Y)$ is called *maximal leakage*, and is defined as the multiplicative increase, upon observing Y , of the probability of correctly guessing a randomized function of X , maximized over all such randomized functions. A closed-form expression for $\mathcal{L}(X \rightarrow Y)$ is given for discrete X and Y , and it is subsequently generalized to handle a large class of random variables. The resulting properties are shown to be consistent with an axiomatic view of a leakage measure, and the definition is shown to be robust to variations in the setup. Moreover, a variant of the Shannon cipher system is studied, in which performance of an encryption scheme is measured using maximal leakage. A single-letter characterization of the optimal limit of (normalized) maximal leakage is derived and asymptotically-optimal encryption schemes are demonstrated. Furthermore, the sample complexity of estimating maximal leakage from data is characterized up to subpolynomial factors. Finally, the *guessing* framework used to define maximal leakage is used to give operational interpretations of commonly used leakage measures, such as Shannon capacity, maximal correlation, and local differential privacy.

Index Terms—Guessing, information leakage, security, Sibson mutual information.

I. INTRODUCTION

HOW much information does an observation “leak” about a quantity on which it depends? This basic question arises in many secrecy and privacy problems in which the quantity of interest is considered sensitive and an observation is available to an adversary. The observation could be intentionally provided to the adversary, as occurs when a curator publishes statistical information about a given population. Or the observation could be an inevitable, if undesirable, consequence of a design. In the latter case, which is the focus

Manuscript received July 20, 2018; revised August 2, 2019; accepted December 10, 2019. Date of publication December 27, 2019; date of current version February 14, 2020. This work was supported in part by the U.S. National Science Foundation under Grant CCF-1704443 and in part by the U.S. Army Research Office under Grant W911NF-18-1-0426. This article was presented in part at the 2016 Annual Conference on Information Sciences and Systems and in part at the 2016 and 2017 IEEE International Symposium on Information Theory.

Ibrahim Issa was with the School of Electrical and Computer Engineering, Cornell University, Ithaca, NY 14853 USA, and also with the School of Computer and Communication Sciences, École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland. He is now with the Department of Electrical and Computer Engineering, American University of Beirut, Beirut 1107 2020, Lebanon (e-mail: ii19@aub.edu.lb).

Aaron B. Wagner is with the School of Electrical and Computer Engineering, Cornell University, Ithaca, NY 14853 USA (e-mail: wagner@cornell.edu).

Sudeep Kamath was with the Electrical Engineering Department, Princeton University, Princeton, NJ 08542 USA. He is now with PDT Partners, New York, NY 10019 USA (e-mail: sudeep.kamath@gmail.com).

Communicated by S. Watanabe, Associate Editor for Shannon Theory. Digital Object Identifier 10.1109/TIT.2019.2962804

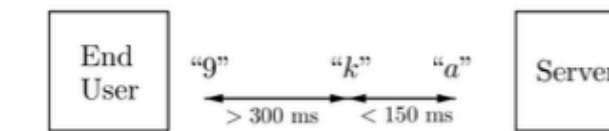


Fig. 1. The Secure Shell: each keystroke is sent immediately to the remote machine.

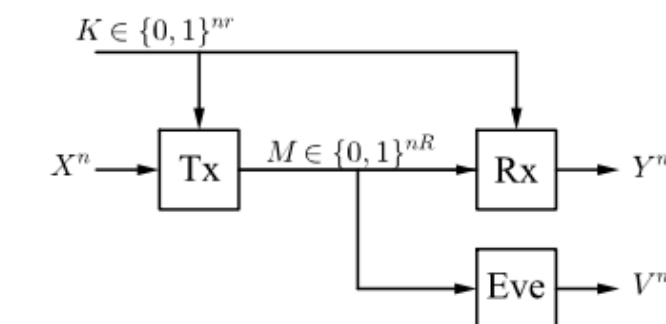


Fig. 2. The Shannon cipher system.

of this paper, we call the observation the output of a *side channel*. Some examples of side channels include:

- When using the Secure Shell (SSH), after the initial handshake, each keystroke is sent immediately to the remote machine, as shown in Figure 1. When communicating over a wireless network, an eavesdropper can observe the timing of the packets which are correlated with the timing of the keystrokes, and hence with the input of the user (e.g., the inter-keystroke delay in ‘ka’ is significantly smaller than that in ‘9k’ [1]).
- Consider an on-chip network that has several processes running simultaneously, one of which is malicious. Because resources such as memory and buses are shared on the chip, the timing characteristics (e.g., memory access delays) observed by the malicious application are affected by the behavior of the other applications (e.g., memory access patterns) and can leak sensitive information such as keys. Similar phenomena occur when users share links or buffers in a communication network [2].
- Consider the Shannon cipher system (shown in Figure 2) in which a transmitter and a receiver are connected through a public noiseless channel and share a secret key. Unless the key rate is very high, the public message depends on the message [3].
- An adversary could “wiretap” a communication channel to intercept transmissions. The wiretap channel is typically noisier than the main channel, but its output nevertheless depends on the transmitted message [4], [5].
- Suppose one would like to anonymously transmit a message through a given network (say, a call for protest on a social network). A powerful adversary (say, a government) could learn the spread of the message (i.e., who

Paper Highlights

Definition 1 (Maximal Leakage): Given a joint distribution P_{XY} on alphabets \mathcal{X} and \mathcal{Y} , the *maximal leakage* from X to Y is defined as

$$\mathcal{L}(X \rightarrow Y) = \sup_{U-X-Y-\hat{U}} \log \frac{\Pr(U = \hat{U})}{\max_{u \in \mathcal{U}} P_U(u)}, \quad (1)$$

where the supremum is over all U and \hat{U} taking values in the same finite, but arbitrary, alphabet.

Theorem 1: For any joint distribution P_{XY} on finite alphabets \mathcal{X} and \mathcal{Y} , the maximal leakage from X to Y is given by the Sibson mutual information of order infinity, $I_\infty(X; Y)$. That is,

$$\mathcal{L}(X \rightarrow Y) = \log \sum_{y \in \mathcal{Y}} \max_{\substack{x \in \mathcal{X}: \\ P_X(x) > 0}} P_{Y|X}(y|x) = I_\infty(X; Y).$$

Tunable Measures for Information Leakage

- A Rényi-like extension of Maximal leakage
- Interpolates between mutual information and maximal leakage

Tunable Measures for Information Leakage and Applications to Privacy-Utility Tradeoffs

Jiachun Liao[✉], *Student Member, IEEE*, Oliver Kosut[✉], *Member, IEEE*, Lalitha Sankar[✉], *Senior Member, IEEE*, and Flavio du Pin Calmon[✉], *Member, IEEE*

Abstract—We introduce a tunable measure for information leakage called *maximal α -leakage*. This measure quantifies the maximal gain of an adversary in inferring any (potentially random) function of a dataset from a release of the data. The inferential capability of the adversary is, in turn, quantified by a class of adversarial loss functions that we introduce as α -loss, $\alpha \in [1, \infty) \cup \{\infty\}$. The choice of α determines the specific adversarial action and ranges from refining a belief (about any function of the data) for $\alpha = 1$ to guessing the most likely value for $\alpha = \infty$ while refining the α^{th} moment of the belief for α in between. Maximal α -leakage then quantifies the adversarial gain under α -loss over all possible functions of the data. In particular, for the extremal values of $\alpha = 1$ and $\alpha = \infty$, maximal α -leakage simplifies to mutual information and maximal leakage, respectively. For $\alpha \in (1, \infty)$ this measure is shown to be the Arimoto channel capacity of order α . We show that maximal α -leakage satisfies data processing inequalities and a sub-additivity property thereby allowing for a weak composition result. Building upon these properties, we use maximal α -leakage as the privacy measure and study the problem of data publishing with privacy guarantees, wherein the utility of the released data is ensured via a *hard distortion* constraint. Unlike average distortion, hard distortion provides a deterministic guarantee of fidelity. We show that under a hard distortion constraint, for $\alpha > 1$ the optimal mechanism is independent of α , and therefore, the resulting optimal tradeoff is the same for all values of $\alpha > 1$. Finally, the tunability of maximal α -leakage as a privacy measure is also illustrated for binary data with average Hamming distortion as the utility measure.

Index Terms—Mutual information, maximal leakage, maximal α -leakage, Sibson mutual information, Arimoto mutual information, f -divergence, privacy-utility tradeoff, hard distortion.

I. INTRODUCTION AND OVERVIEW

THE measure and control of private information leakage is a recognized objective in communications, information theory, and computer science. Modern cryptography [1]–[3], for example, aims at designing and analyzing security systems that are believed to be impervious to computationally bounded adversaries. Alternatively, information-theoretic security studies settings where an asymmetry of information between

an adversary and the legitimate parties (e.g., the wiretap channel [4]–[6]) can be exploited to guarantee that no private information is leaked regardless of computational assumptions. An adversary that *only* observes the output of a (computationally) secure cipher or cannot overcome the information asymmetry in a wiretap-like setting does not, for all practical purposes, pose a privacy risk.

However, modern applications such as online data sharing, social networks, cloud-based services, and mobile computing have significantly increased the number of ways in which private information can leak. Services that require a user to disclose data in order to receive utility inevitably incur a privacy risk through unwanted inferences. For example, *sensitive information* such as political preference, medical conditions, and identity can be reliably estimated from movie ratings [7], online shopping patterns, [8], and via deanonymization and tracking of interactions in social network data [9], [10], respectively. Moreover, practical implementations of cryptographic schemes are susceptible to so-called “side-channel attacks,” where sensitive information leaks through unexpected channels. For example, a malicious application may get timing characteristics [11], [12]. In these examples, an adversary that observes information leaked through a side-channel can more reliably infer private data, such as a key or a plaintext.

Several (often overlapping) definitions of privacy/information leakage have been proposed over the past decade. The most widely adopted measure is differential privacy (DP) [13], [14], which was introduced within the context of querying databases. DP seeks to ensure that changes in the database entries do not significantly influence the value of a query. A variety of information-theoretic measures have also been proposed as leakage measures. Foremost among them is mutual information (MI): its use as a privacy measure in [15]–[24] is inspired by the common appearance of MI as an operationally-meaningful quantity throughout the literature on communication systems. In a similar vein, divergence-based quantities such as total variation distance between the prior and posterior distributions [25] have also been proposed as leakage measures. Information-theoretic measures have been studied in the DP community via Rényi differential privacy which is based on Rényi divergence [26] that allow relaxing the original definition of DP in order to enable better utility guarantees. However, the gamut of information-theoretic leakage measures proposed to address the privacy problem do not *yet* have clear operational meanings or adversarial models in their definitions.

Manuscript received September 24, 2018; revised April 29, 2019; accepted August 1, 2019. Date of publication August 16, 2019; date of current version November 20, 2019. This work was supported in part by the National Science Foundation under Grants CCF-1422358, CCF-1350914, CIF-1815361, and CIF-1901243. This article was presented in part at the 2018 IEEE International Symposium on Information Theory and Information Theory Workshop.

J. Liao, O. Kosut, and L. Sankar are with Arizona State University, Tempe, AZ 85281 USA (e-mail: jiachun.liao@asu.edu).

F. du Pin Calmon is with Harvard University, Cambridge, MA USA.

Communicated by S. Watanabe, Associate Editor for Shannon Theory.

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIT.2019.2935768

0018-9448 © 2019 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

Paper Highlights

Definition 3 (α -loss). Let random variables X , Y and \hat{X} form a Markov chain $X - Y - \hat{X}$, where \hat{X} is an estimator of X . The α -loss of the strategy $P_{\hat{X}|Y}$ for estimating X from Y is

$$\ell_\alpha(x, y, P_{\hat{X}|Y}) = \frac{\alpha}{\alpha - 1} \left(1 - P_{\hat{X}|Y}(x|y)^{\frac{\alpha-1}{\alpha}} \right), \quad (19)$$

where $\alpha \in (1, \infty)$. It is defined by its continuous extension for $\alpha = 1$ and $\alpha = \infty$, respectively, and is given by

$$\ell_1(x, y, P_{\hat{X}|Y}) = \lim_{\alpha \rightarrow 1} \ell_\alpha(x, y, P_{\hat{X}|Y}) = \log \frac{1}{P_{\hat{X}|Y}(x|y)}, \quad (20)$$

$$\ell_\infty(x, y, P_{\hat{X}|Y}) = \lim_{\alpha \rightarrow \infty} \ell_\alpha(x, y, P_{\hat{X}|Y}) = 1 - P_{\hat{X}|Y}(x|y). \quad (21)$$

Definition 5 (α -Leakage). Given a joint distribution $P_{X,Y}$ and an estimator \hat{X} with the same support as X , the α -leakage from X to Y is defined as

$$\mathcal{L}_\alpha(X \rightarrow Y) \triangleq \frac{\alpha}{\alpha - 1} \log \frac{\max_{P_{\hat{X}|Y}} \mathbb{E} \left[P_{\hat{X}|Y}(X|Y)^{\frac{\alpha-1}{\alpha}} \right]}{\max_{P_{\hat{X}}} \mathbb{E} \left[P_{\hat{X}}(X)^{\frac{\alpha-1}{\alpha}} \right]}, \quad (27)$$

for $\alpha \in (1, \infty)$ and by the continuous extension of (27) for $\alpha = 1$ and ∞ .

Definition 6 (Maximal α -Leakage). Given a joint distribution $P_{X,Y}$ on finite alphabets $\mathcal{X} \times \mathcal{Y}$, the maximal α -leakage from X to Y is defined as

$$\mathcal{L}_\alpha^{\max}(X \rightarrow Y) \triangleq \sup_{U \sim X-Y} \mathcal{L}_\alpha(U; Y), \quad (28)$$

where $1 \leq \alpha \leq \infty$, and U represents any function of X and takes values from an arbitrary finite alphabet.

Note that for $\alpha \geq 1$,

$$\begin{aligned} & \max_{P_{\hat{U}|Y}} \mathbb{E} \left[P_{\hat{U}|Y}(U|Y)^{\frac{\alpha-1}{\alpha}} \right] \\ &= 1 - \frac{\alpha - 1}{\alpha} \min_{P_{\hat{U}|Y}} \mathbb{E} \left[\ell_\alpha(U, Y, P_{\hat{U}|Y}) \right]. \end{aligned} \quad (29)$$

Theorem 2. For $1 \leq \alpha \leq \infty$, the maximal α -leakage defined in (28) simplifies to

$$\begin{aligned} & \mathcal{L}_\alpha^{\max}(X \rightarrow Y) \\ &= \begin{cases} \sup_{P_{\tilde{X}}} I_\alpha^S(\tilde{X}; Y) = \sup_{P_{\tilde{X}}} I_\alpha^A(\tilde{X}; Y), & 1 < \alpha \leq \infty \\ I(X; Y), & \alpha = 1 \end{cases} \end{aligned} \quad \begin{aligned} & (31a) \\ & (31b) \end{aligned}$$

where $P_{\tilde{X}}$ is a probability distribution over the support of P_X .

Calibrating Noise to Sensitivity in Private Data Analysis (Differential Privacy)

- An information leakage measure motivated by statistical databases
- Measures the perturbation in the output due to small changes in input
- Has many nice properties and is extremely well studied

$$\mathcal{DP}(X \rightarrow Y) = \sup_{x_1, x_2 \in \mathcal{X} : d(x_1, x_2) = 1, \mathcal{T} \subset \mathcal{Y}} \log \frac{\mathbb{P}[Y \in \mathcal{T} | X = x_1]}{\mathbb{P}[Y \in \mathcal{T} | X = x_2]}$$

Calibrating Noise to Sensitivity in Private Data Analysis

Cynthia Dwork¹, Frank McSherry¹, Kobbi Nissim², and Adam Smith^{3*}

¹ Microsoft Research, Silicon Valley. {dwork,mcsherry}@microsoft.com

² Ben-Gurion University. kobbi@cs.bgu.ac.il

³ Weizmann Institute of Science. adam.smith@weizmann.ac.il

Abstract. We continue a line of research initiated in [10, 11] on privacy-preserving statistical databases. Consider a trusted server that holds a database of sensitive information. Given a query function f mapping databases to reals, the so-called *true answer* is the result of applying f to the database. To protect privacy, the true answer is perturbed by the addition of random noise generated according to a carefully chosen distribution, and this response, the true answer plus noise, is returned to the user.

Previous work focused on the case of noisy sums, in which $f = \sum_i g(x_i)$, where x_i denotes the i th row of the database and g maps database rows to $[0, 1]$. We extend the study to general functions f , proving that privacy can be preserved by calibrating the standard deviation of the noise according to the *sensitivity* of the function f . Roughly speaking, this is the amount that any single argument to f can change its output. The new analysis shows that for several particular applications substantially less noise is needed than was previously understood to be the case.

The first step is a very clean characterization of privacy in terms of indistinguishability of transcripts. Additionally, we obtain separation results showing the increased value of interactive sanitization mechanisms over non-interactive.

1 Introduction

We continue a line of research initiated in [10, 11] on privacy in *statistical* databases. A statistic is a quantity computed from a sample. Intuitively, if the database is a representative sample of an underlying population, the goal of a privacy-preserving statistical database is to enable the user to learn properties of the population as a whole while protecting the privacy of the individual contributors.

We assume the database is held by a trusted server. On input a query function f mapping databases to reals, the so-called *true answer* is the result of applying f to the database. To protect privacy, the true answer is perturbed by the addition

* Supported by the Louis L. and Anita M. Perlman Postdoctoral Fellowship.

The Composition Theorem for Differential Privacy

- A recent work on Differential Privacy
- Connects Differential Privacy to Hypothesis Testing
- Proves a new composition theorem for differential private mechanisms

The Composition Theorem for Differential Privacy

Peter Kairouz, *Member, IEEE*, Sewoong Oh, *Member, IEEE*, and Pramod Viswanath, *Fellow, IEEE*

Abstract—Sequential querying of differentially private mechanisms degrades the overall privacy level. In this paper, we answer the fundamental question of characterizing the level of overall privacy degradation as a function of the number of queries and the privacy levels maintained by each privatization mechanism. Our solution is complete: we prove an upper bound on the overall privacy level and construct a sequence of privatization mechanisms that achieves this bound. The key innovation is the introduction of an operational interpretation of differential privacy (involving hypothesis testing) and the use of a data processing inequality along with its converse. Our result improves over the state of the art, and has immediate connections to several problems studied in the literature.

Index Terms—Differential privacy, hypothesis testing.

I. INTRODUCTION

DIFFERENTIAL privacy is a formal framework to quantify to what extent individual privacy in a statistical database is preserved while releasing useful aggregate information about the database. It provides strong privacy guarantees by requiring the indistinguishability of whether or not an individual is in a database based on the released information, regardless of the side information on the other aspects of the database the adversary may possess. Denoting the database when the individual is present as D_1 and as D_0 when the individual is not, a differentially private mechanism provides indistinguishability guarantees with respect to the pair (D_0, D_1) . The databases D_0 and D_1 are referred to as “neighboring” databases.

Definition 1 (Differential Privacy [10], [12]): A randomized mechanism M over a set of databases is (ϵ, δ) -differentially private if for all pairs of neighboring databases D_0 and D_1 , and for all sets S in the output space of the mechanism \mathcal{X} ,

$$\mathbb{P}(M(D_0) \in S) \leq e^\epsilon \mathbb{P}(M(D_1) \in S) + \delta.$$

Manuscript received January 20, 2014; revised December 3, 2015, May 27, 2016, and January 12, 2017; accepted March 15, 2017. Date of publication March 21, 2017; date of current version May 18, 2017. This work was supported in part by NSF CISE under Award CCF-1422278, Award CCF-1553452, NSF SaTC Award CNS-1527754, NSF CMMI Award MES-1450848, NSF ENG Award ECCS-1232257, and in part by the Google Faculty Research Award. This paper was presented at the 2015 International Conference on Machine Learning in [KOV15].

P. Kairouz is with the Department of Electrical Engineering, Stanford University, Stanford, CA 94305 USA (e-mail: kairouz2@illinois.edu).

S. Oh is with the Department of Industrial and Enterprise Systems Engineering, University of Illinois at Urbana-Champaign, Champaign, IL 61801, USA (e-mail: swoh@illinois.edu).

P. Viswanath is with the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Champaign, IL 61801, USA (e-mail: pramodv@illinois.edu).

Communicated by A. Smith, Associate Editor for Complexity and Cryptography.

Digital Object Identifier 10.1109/TIT.2017.2685505

A basic problem in differential privacy is how privacy of a fixed pair of neighbors (D_0, D_1) degrades under *composition* of interactive queries when each query, individually, meets certain differential privacy guarantees. A routine argument shows that the composition of k queries, each of which is (ϵ, δ) -differentially private, is at least $(k\epsilon, k\delta)$ -differentially private [10]–[12], [16]. A tighter bound of $(\tilde{\epsilon}_{\tilde{\delta}}, k\delta + \tilde{\delta})$ -differential privacy under k -fold adaptive composition is provided, using more sophisticated arguments, in [16] for the case when each of the individual queries is (ϵ, δ) -differentially private. Here $\tilde{\epsilon}_{\tilde{\delta}} = O(k\epsilon^2 + \epsilon\sqrt{k\log(1/\tilde{\delta})})$. On the other hand, it was not known if this bound could be improved until this work.

Our main result is the *exact* characterization of the privacy guarantee under k -fold composition. Any k -fold adaptive composition of (ϵ, δ) -differentially private mechanisms satisfies the privacy guarantee stated in Theorem 9. Further, we demonstrate a specific sequence of (nonadaptive) privacy mechanisms which when composed, degrade the privacy to the level guaranteed in Theorem 9. Our result entails a strict improvement over the state-of-the-art result in [16]. This can be seen immediately in the following approximation – using the same notation as above, the value of $\tilde{\epsilon}_{\tilde{\delta}}$ is now reduced to $\tilde{\epsilon}_{\tilde{\delta}} = O(k\epsilon^2 + \epsilon\sqrt{k\log(e + (\epsilon\sqrt{k}/\tilde{\delta}))})$. Since a typical choice of $\tilde{\delta}$ is $\tilde{\delta} = \Theta(k\delta)$, in the regime where $\epsilon = \Theta(\sqrt{k}\delta)$, this improves the existing guarantee by a logarithmic factor. The gain is especially significant when both ϵ and δ are small.

We view differential privacy as a guarantee on the two types of error (false alarm and missed detection) in a binary hypothesis testing problem involving two neighboring databases. This approach is similar to the previous work of Wasserman and Zhou [33]. Our work leverages two benefits of this *operational interpretation* of differential privacy.

- The first is conceptual. The operational setting directs the logic of the steps of the proof, makes the arguments straightforward, and readily allows for generalizations such as heterogeneous compositions.
- The second is technical. The operational interpretation of hypothesis testing brings both the natural data processing inequality and the strong converse to the data processing inequality. These inequalities, while simple by themselves, lead to surprisingly strong technical results. As an aside, we mention that there is a strong tradition of such derivations in the information theory literature: the Fisher information inequality [5], [34], the entropy power inequality [5], [31], [32], an extremal inequality involving mutual informations [28], matrix determinant

Paper Highlights

Given a random output Y of a database access mechanism M , consider the following hypothesis testing experiment. We choose a null hypothesis as database D_0 and alternative hypothesis as D_1 :

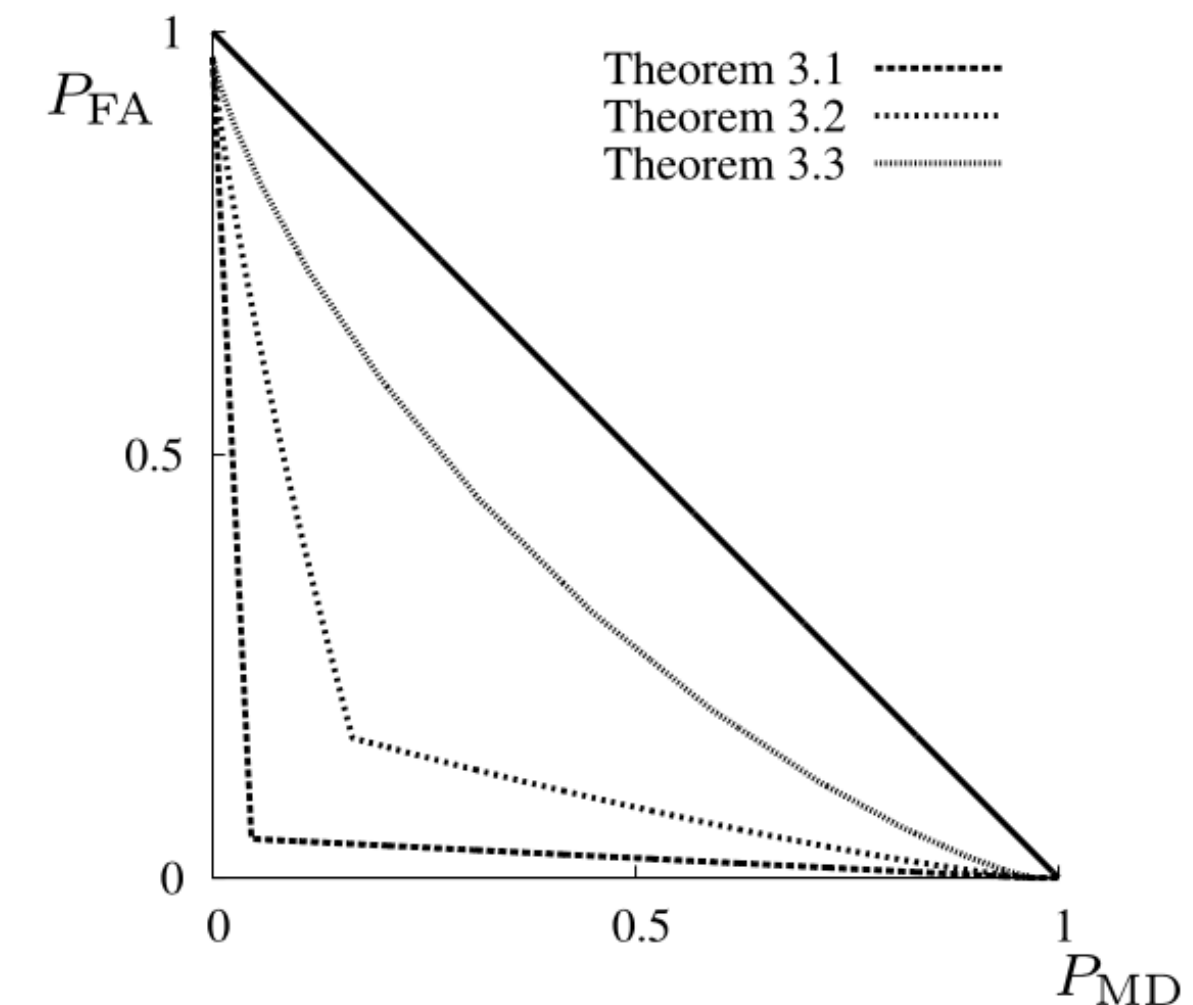
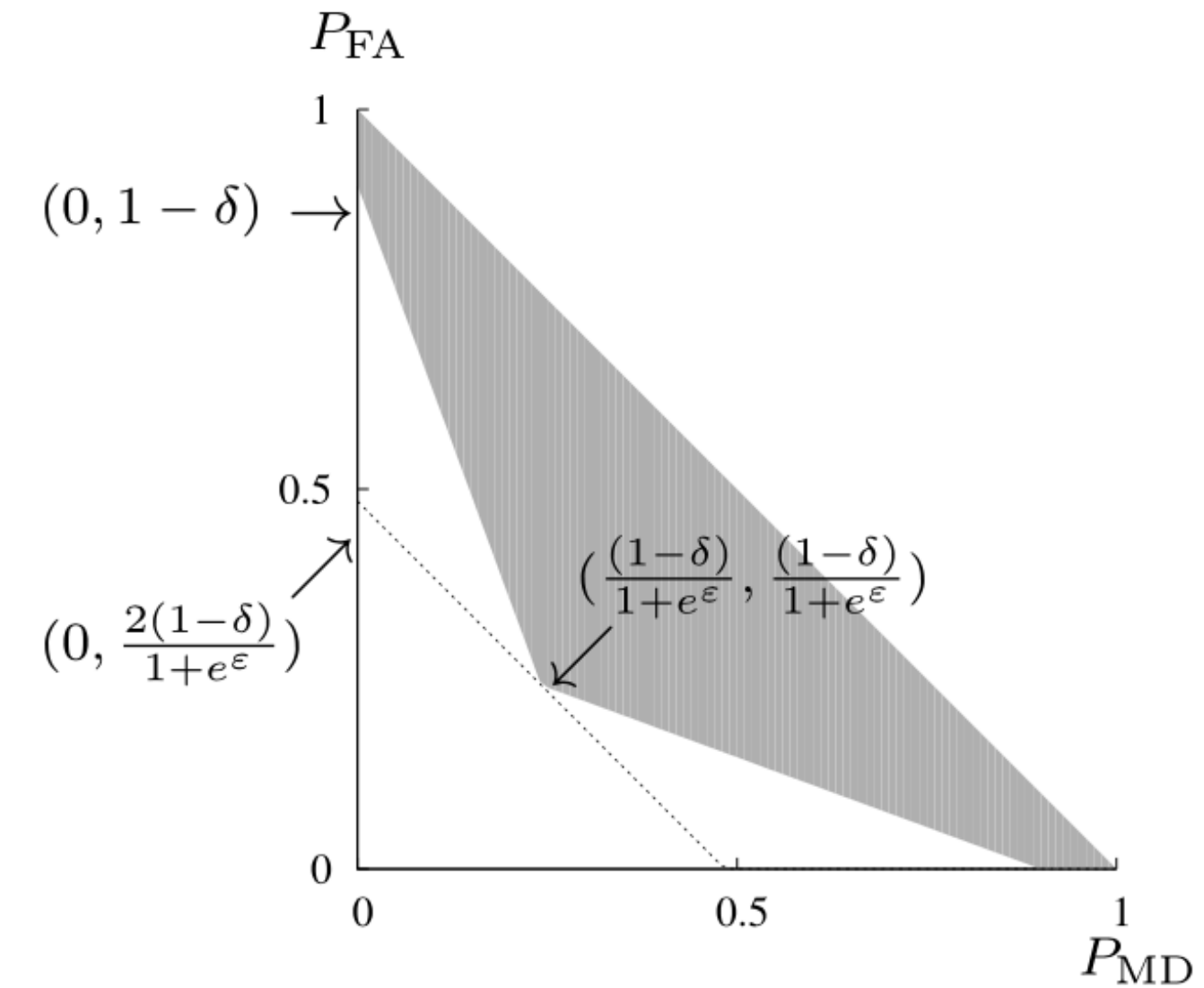
$H_0 : Y$ came from a database D_0 ,

$H_1 : Y$ came from a database D_1 .

Theorem 2: For any $\varepsilon \geq 0$ and $\delta \in [0, 1]$, a database mechanism M is (ε, δ) -differentially private if and only if the following conditions are satisfied for all pairs of neighboring databases D_0 and D_1 , and all rejection region $S \subseteq \mathcal{X}$:

$$P_{\text{FA}}(D_0, D_1, M, S) + e^\varepsilon P_{\text{MD}}(D_0, D_1, M, S) \geq 1 - \delta, \text{ and}$$

$$e^\varepsilon P_{\text{FA}}(D_0, D_1, M, S) + P_{\text{MD}}(D_0, D_1, M, S) \geq 1 - \delta. \quad (1)$$



Rényi Differential Privacy

- A recent work on Differential Privacy
- Does a Rényi-like extension of Differential Privacy
- Proves a new composition theorem for differential private mechanisms

Rényi Differential Privacy

Ilya Mironov
Google Brain

Abstract—We propose a natural relaxation of differential privacy based on the Rényi divergence. Closely related notions have appeared in several recent papers that analyzed composition of differentially private mechanisms. We argue that the useful analytical tool can be used as a privacy definition, compactly and accurately representing guarantees on the tails of the privacy loss.

We demonstrate that the new definition shares many important properties with the standard definition of differential privacy, while additionally allowing tighter analysis of composite heterogeneous mechanisms.

I. INTRODUCTION

Differential privacy, introduced by Dwork et al. [1], has been embraced by multiple research communities as a commonly accepted notion of privacy for algorithms on statistical databases. As applications of differential privacy begin to emerge, practical concerns of *tracking* and *communicating* privacy guarantees are coming to the fore.

Informally, differential privacy bounds a shift in the output distribution of a randomized algorithm that can be induced by a small change in its input. The standard definition of ϵ -differential privacy puts a multiplicative upper bound on the worst-case change in the distribution's density.

Several relaxations of differential privacy explored other measures of closeness between two distributions. The most common such relaxation, the (ϵ, δ) definition, has been a method of choice for expressing privacy guarantees of a variety of differentially private algorithms, especially those that rely on the Gaussian additive noise mechanism or whose analysis follows from composition theorems. The additive δ parameter allows suppressing the long tails of the mechanism's distribution where pure ϵ -differential privacy guarantees may not hold.

Compared to the standard definition, (ϵ, δ) -differential privacy offers asymptotically smaller cumulative loss under composition and allows greater flexibility in the selection of privacy-preserving mechanisms.

Despite its notable advantages and numerous applications, the definition of (ϵ, δ) -differential privacy is an imperfect fit for its two most common use cases: the Gaussian mechanism and a composition rule. We briefly sketch them here and elaborate on these points in the next section.

The first application of (ϵ, δ) -differential privacy was the analysis of the Gaussian noise mechanism [2]. In contrast with the Laplace mechanism, whose privacy guarantee is characterized tightly and accurately by ϵ -differential privacy, a single Gaussian mechanism satisfies a *curve* of $(\epsilon(\delta), \delta)$ -differential privacy definitions. Picking any one point on this curve leaves out important information about the mechanism's actual behavior.

The second common use of (ϵ, δ) -differential privacy is due to applications of advanced composition theorems. The central feature of ϵ -differential privacy is that it is closed under composition; moreover, the ϵ parameters of composed mechanisms simply add up, which motivates the concept of a *privacy budget*. By relaxing the guarantee to (ϵ, δ) -differential privacy, advanced composition allows tighter analyses for compositions of (pure) differentially private mechanisms. Iterating this process, however, quickly leads to a combinatorial explosion of parameters, as each application of an advanced composition theorem leads to a wide selection of possibilities for $(\epsilon(\delta), \delta)$ -differentially private guarantees.

In part to address the shortcomings of (ϵ, δ) -differential privacy, several recent works, surveyed in the next section, explored the use of higher-order *moments* as a way of bounding the tails of the privacy loss variable.

Inspired by these theoretical results and their applications, we propose *Rényi differential privacy* as a natural relaxation of differential privacy that is well-suited for expressing guarantees of privacy-preserving algorithms and for composition of heterogeneous mechanisms. Compared to (ϵ, δ) -differential privacy, Rényi differential privacy is a strictly stronger privacy definition. It offers an operationally convenient and quantitatively accurate way of tracking cumulative privacy loss throughout execution of a standalone differentially private mechanism and across many such mechanisms. Most significantly, Rényi differential privacy allows combining the intuitive and appealing concept of a privacy budget with application of advanced composition theorems.

The paper presents a self-contained exposition of the new definition, unifying current literature and demonstrating its applications. The organization of the paper is as follows. Section II reviews the standard definition of differential privacy, its (ϵ, δ) relaxation and its most common uses. Section III introduces the definition of Rényi differential privacy and proves its basic properties that parallel those of ϵ -differential privacy, summarizing the results in Table I. Section IV demonstrates a reduction from Rényi differential privacy to (ϵ, δ) -differential privacy, followed by a proof of an advanced composition theorem in Section V. Section VI applies Rényi differential privacy to analysis of several basic mechanisms: randomized response for predicates, Laplace and Gaussian (see Table II for a brief summary). Section VII discusses assessment of risk due to application of a Rényi differentially private mechanism and use of Rényi differential privacy as a privacy loss tracking tool. Section VIII concludes with open questions.

arXiv:1702.07476v3 [cs.CR] 25 Aug 2017

Paper Highlights

Relative Entropy

Kullback-Leibler divergence (also known as relative entropy):

$$D_1(P\|Q) = \mathbb{E}_{x \sim P} \log \frac{P(x)}{Q(x)}.$$

Definition 4 ((α, ϵ) -RDP). A randomized mechanism $f: \mathcal{D} \mapsto \mathcal{R}$ is said to have ϵ -Rényi differential privacy of order α , or (α, ϵ) -RDP for short, if for any adjacent $D, D' \in \mathcal{D}$ it holds that

$$D_\alpha(f(D)\|f(D')) \leq \epsilon.$$

Definition 3 (Rényi divergence). For two probability distributions P and Q defined over \mathcal{R} , the Rényi divergence of order $\alpha > 1$ is

$$D_\alpha(P\|Q) \triangleq \frac{1}{\alpha - 1} \log \mathbb{E}_{x \sim Q} \left(\frac{P(x)}{Q(x)} \right)^\alpha.$$

Differential Privacy

$$D_\infty(P\|Q) = \sup_{x \in \text{supp } Q} \log \frac{P(x)}{Q(x)}.$$

Further Reading

Technical Privacy Metrics: A Systematic Survey

ISABEL WAGNER, De Montfort University
DAVID ECKHOFF, TUMCREATE Ltd.

The goal of privacy metrics is to measure the degree of privacy enjoyed by users in a system and the amount of protection offered by privacy-enhancing technologies. In this way, privacy metrics contribute to improving user privacy in the digital world. The diversity and complexity of privacy metrics in the literature make an informed choice of metrics challenging. As a result, instead of using existing metrics, new metrics are proposed frequently, and privacy studies are often incomparable. In this survey, we alleviate these problems by structuring the landscape of privacy metrics. To this end, we explain and discuss a selection of over 80 privacy metrics and introduce categorizations based on the aspect of privacy they measure, their required inputs, and the type of data that needs protection. In addition, we present a method on how to choose privacy metrics based on nine questions that help identify the right privacy metrics for a given scenario, and highlight topics where additional work on privacy metrics is needed. Our survey spans multiple privacy domains and can be understood as a general framework for privacy measurement.

CCS Concepts: • **General and reference** → **Metrics**; • **Security and privacy** → **Pseudonymity, anonymity and untraceability**; **Privacy protections**; *Privacy-preserving protocols*; *Social network security and privacy*; *Usability in security and privacy*; • **Networks** → *Network privacy and anonymity*; • **Theory of computation** → *Theory of database privacy and security*;

Additional Key Words and Phrases: Privacy metrics, measuring privacy

ACM Reference format:

Isabel Wagner and David Eckhoff. 2018. Technical Privacy Metrics: A Systematic Survey. *ACM Comput. Surv.* 51, 3, Article 57 (June 2018), 38 pages.
<https://doi.org/10.1145/3168389>

1 INTRODUCTION

Privacy is a fundamental human right codified in the United Nations Universal Declaration of Human Rights, which states that “no one shall be subjected to arbitrary interference with his privacy, family, home or correspondence” [126, Art. 12]. However, it is difficult to define what exactly privacy is. As early as 1967, Westin [134] defined privacy as “the ability of an individual to control the terms under which personal information is acquired and used.” Personal information, according to the EU General Data Protection Regulation (and the OECD privacy framework [101]), is “any information relating to an [...] identifiable natural person” [45].

David Eckhoff is financially supported by the Singapore National Research Foundation under its Campus for Research Excellence And Technological Enterprise (CREATE) programme.
Authors' addresses: I. Wagner, De Montfort University, Cyber Security Centre, The Gateway, Gateway House, Leicester, LE1 9BH, UK; email: isabel.wagner@dmu.ac.uk; D. Eckhoff, TUMCREATE Ltd., 1 Create Way, #10-02 CREATE Tower, 138602 Singapore; email: david.eckhoff@tum-create.edu.sg.
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2018 ACM 0360-0300/2018/06-ART57 \$15.00
<https://doi.org/10.1145/3168389>

An Overview of Information-Theoretic Security and Privacy: Metrics, Limits and Applications

Matthieu Bloch[✉], *Senior Member, IEEE*, Onur Günlü[✉], *Member, IEEE*, Aylin Yener[✉], *Fellow, IEEE*,

Frédérique Oggier[✉], H. Vincent Poor[✉], *Life Fellow, IEEE*, Lalitha Sankar[✉], *Senior Member, IEEE*, and

Rafael F. Schaefer[✉], *Senior Member, IEEE*

Abstract—This tutorial reviews fundamental contributions to information security. An integrative viewpoint is taken that explains the security metrics, including secrecy, privacy, and others, the methodology of information-theoretic approaches, along with the arising system design principles, as well as techniques that enable the information-theoretic designs to be applied in real communication and computing systems. The tutorial, while summarizing these contributions, argues for the simultaneous pivotal role of fundamental limits and coding techniques for secure communication system design.

Index Terms—Information-theoretic security, privacy, wiretap channel, secret key agreement, coding, physical-layer security, security and privacy metrics, adversarial models.

I. INTRODUCTION

INFORMATION security, a broad umbrella term that includes attributes including secrecy, privacy, and trust, has arguably become as important as information reliability in system design, especially so, as society at large conducts most operations virtually and as future generations of applications and devices emerge that amalgamates communication, sensing, computing, and control. In current systems, information

security is largely treated as an addition to the network operations rather than a foundational design constraint at the outset. Consequently, securing information that flows over networked systems is largely guaranteed by higher network layer protocols.

While this layered approach has had undeniable success, future and emerging systems exhibit unique characteristics that challenge this prevalent view of security. The deployment of 5G, the advent of the IoT, the current and upcoming cyber-physical autonomous systems, and the envisioned all connected 6G world have all exacerbated the concerns for security and privacy in communication networks. In the next decade and beyond, tens of billions of devices are expected to be collecting and transmitting data over networks. The heterogeneity of these devices in terms of resources and capabilities, e.g., battery power, computational power, communication and storage capabilities, renders the approach to date of relying solely on computational approaches for security, e.g., cryptographic solutions, difficult. For example, networks with energy and computational power-limited IoT devices would benefit from lightweight security mechanisms that do not incur the overhead of traditional public-key infrastructures. Similarly, the stringent performance constraints of cyber-physical systems make increasingly apparent that security cannot be handled independently of other parameters, such as power consumption and latency, leading to unavoidable application dependent trade-offs. Cyber-physical systems would then benefit from bringing security closer to control in order to reduce overhead and latency in operation. Finally, all future massively connected systems would benefit from security mechanisms built into their foundation, e.g., to avoid the costly software updates required when new more powerful attacks emerge as a result of increasing computing power. Noting information security and privacy has a much larger domain of interest and impact, this tutorial focuses on communications as an exemplar to highlight recent advances in information-theoretic security.

Information-theoretic security [1], [2] aims at providing solutions to the aforementioned challenges, by offering a framework in which the security of information flows can be *measured* with quantitative information-theoretic metrics and *enforced* using a combination of signaling and coding mechanisms at the lower layers of the communication protocols. At its core, information-theoretic security embraces the observation structures inherent to communication systems. Specifically, acknowledging that legitimate users and adversaries obtain distinct signals through noisy and lossy channels, the asymmetry is harnessed through signal processing and coding mechanisms to control information

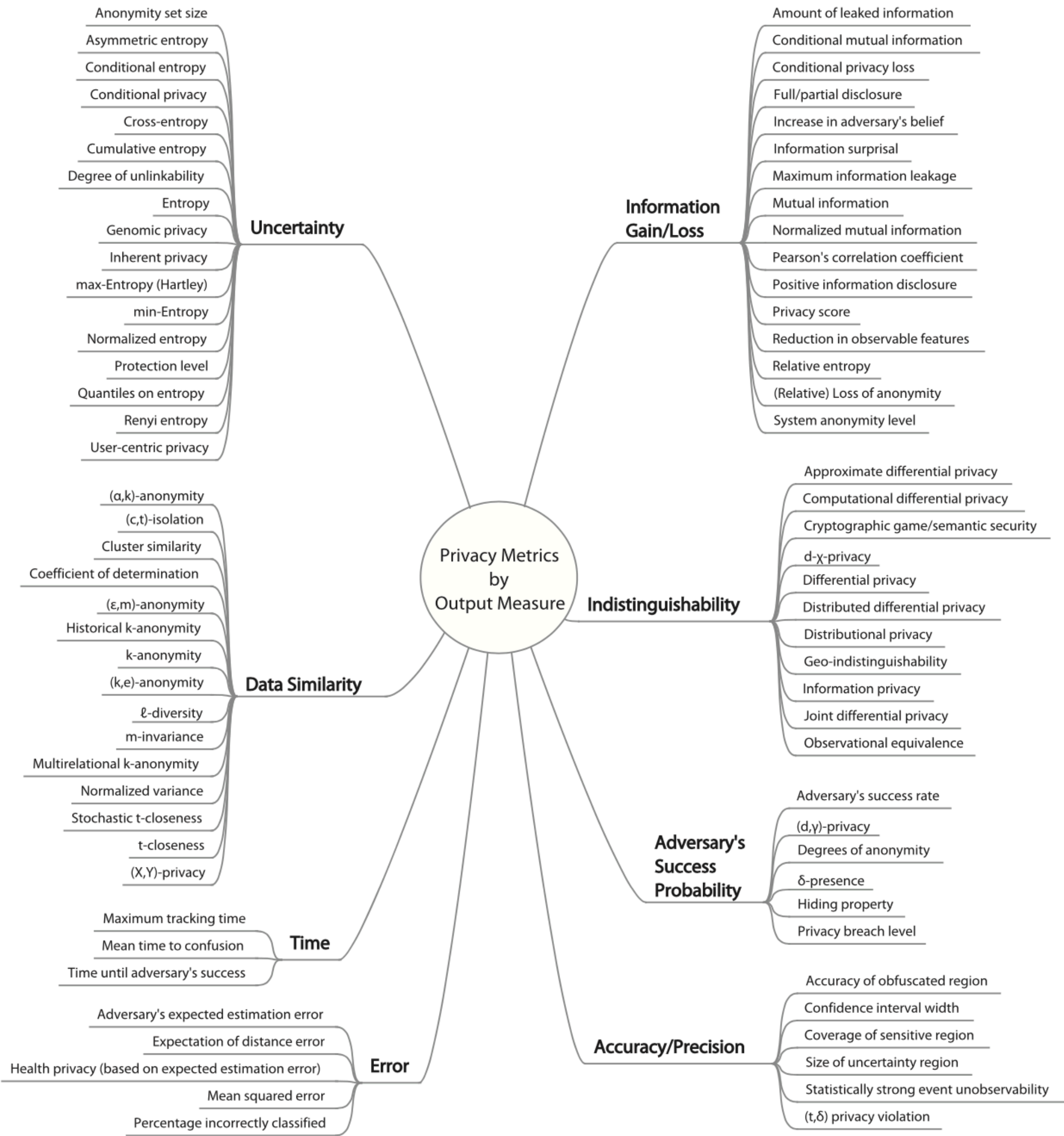


Fig. 1. Taxonomy of privacy metrics, classified by output.

Many More Security and Privacy Settings

- Location Data (correlation in space)
- Temporal Data (correlation in time)
- ...

Quantifying Location Privacy

Reza Shokri, George Theodorakopoulos, Jean-Yves Le Boudec, and Jean-Pierre Hubaux
LCA, EPFL, Lausanne, Switzerland
firstname.lastname@epfl.ch

Abstract—It is a well-known fact that the progress of personal communication devices leads to serious concerns about privacy in general, and location privacy in particular. As a response to these issues, a number of Location-Privacy Protection Mechanisms (LPPMs) have been proposed during the last decade. However, their assessment and comparison remains problematic because of the absence of a systematic method to quantify them. In particular, the assumptions about the attacker’s model tend to be incomplete, with the risk of a possibly wrong estimation of the users’ location privacy.

In this paper, we address these issues by providing a formal framework for the analysis of LPPMs; it captures, in particular, the prior information that might be available to the attacker, and various attacks that he can perform. The privacy of users and the success of the adversary in his location-inference attacks are two sides of the same coin. We revise location privacy by giving a simple, yet comprehensive, model to formulate all types of location-information disclosure attacks. Thus, by formalizing the adversary’s performance, we propose and justify the right metric to quantify location privacy. We clarify the difference between three aspects of the adversary’s inference attacks, namely their *accuracy*, *certainty*, and *correctness*. We show that correctness determines the privacy of users. In other words, the expected estimation error of the adversary is the metric of users’ location privacy. We rely on well-established statistical methods to formalize and implement the attacks in a tool: the *Location-Privacy Meter* that measures the location privacy of mobile users, given various LPPMs. In addition to evaluating some example LPPMs, by using our tool, we assess the appropriateness of some popular metrics for location privacy: entropy and k-anonymity. The results show a lack of satisfactory correlation between these two metrics and the success of the adversary in inferring the users’ actual locations.

Keywords-Location Privacy; Evaluation Framework; Location Traces; Quantifying Metric; Location-Privacy Meter

I. INTRODUCTION

Most people are now equipped with smart phones with many sophisticated sensors and actuators closely related to their activities. Each of these devices is usually equipped with high-precision localization capabilities, based for example on a GPS receiver or on triangulation with nearby base stations or access points. In addition, the environment is more and more populated by sensors and smart devices, with which smart phones interact.

The usage of these personal communication devices, although providing convenience to their owners, leaves an almost indelible digital trace of their whereabouts. A trace is not only a set of positions on a map. The contextual

information attached to a trace tells much about the individuals’ habits, interests, activities, and relationships. It can also reveal their personal or corporate secrets. It can expose the users to unwanted advertisement and location-based spams/scams, cause social reputation or economic damage, and make them victims of blackmail or even physical violence. Additionally, information disclosure breaks the balance of power between the informed entity and the entity about which this information is disclosed.

In the meantime, the tools required to analyze such traces have made tremendous progress: sophisticated data mining algorithms can leverage on fast growing storage and processing power, facilitating, for example, the analysis of multiple databases in parallel. This means that the negative side-effects of insufficient location privacy are becoming more and more threatening.

Users should have the right to control the amount of information (about themselves) that is disclosed and shared with others. This can be achieved in several ways. Users can share a minimum amount of information, or share it only with few trusted entities. Privacy policies can be put in place to force organizations to protect their users’ privacy. Finally, systems can be designed in a privacy-conscious manner, so they do not leak information to untrusted entities.

This paper refers to the last ambition. However, our goal here is not to design yet another location privacy protection mechanism (LPPM), but rather to try to make progress on the **quantification** of the performance of an LPPM. This is an important topic, because (i) human beings are notoriously bad estimators of risks (including privacy risks), (ii) it is the only way to make meaningful comparisons between different LPPMs and (iii) the research literature is not yet mature enough on the topic.

Let us develop this last reason. In specific areas, several contributions have been made to quantify privacy, be it for databases [8], for anonymity protocols [3], for anonymization networks [24], or for RFID privacy [25]. Yet, in the field of location privacy, notwithstanding many contributions from different disciplines (such as databases, mobile networks, and ubiquitous computing) for protecting location privacy, the lack of a unified and generic formal framework for specifying protection mechanisms and also for evaluating location privacy is evident. This has led to the divergence of (nevertheless interesting) contributions and, hence, has caused confusion about which mechanisms are

COURSE TAKEAWAYS

**Probability and Statistics are the
mathematical languages in which we
model information**

**There are many information measures
with different mathematical properties**

**There are many application domains
which require different measures**