

Information–Theoretic Analysis of Privacy Protection for Noisy Identification Based on Soft Fingerprinting

Vladimir B. Balakirsky, Svyatoslav Voloshynovskiy, Oleksiy Koval, Taras Holotyak

Data Security Association “Confident”, Russia
University of Geneva, Switzerland

e-mail: v.b_balakirsky@rambler.ru, {svolos, Oleksiy.Koval, Taras.Holotyak}@unige.ch

ABSTRACT

Identification of contents or objects based on some data that are stored/distributed in public domain is required in various applications. At the same time, these data should not reveal any information about original content or object that may be missused in terms of privacy leakage. We consider a privacy protection strategy based on reliable components of data and investigate the performance of this scheme with respect to achievable identification rate and privacy leak. The data stored/distributed in the public domain are binary, while the encoder and the decoder operate with real data. The advocated strategy is referred to as soft fingerprinting.

Keywords

Information theory, Soft fingerprinting, Identification rate, Privacy leak.

1. INTRODUCTION

Many problems of modern multimedia management (content filtering, content retrieval/search, content tagging and recommendation), multimedia security (copyright protection, broadcast monitoring, etc.) and physical object security such as biometrics and anti-counterfeiting require efficient tool providing content identification. To find the reasonable trade-off between accuracy, privacy leak, complexity and memory storage, most identification techniques use binary digital fingerprinting. Usually a binary fingerprint represents a short, robust and distinctive content description that allows to overcome fundamental sensitivity restrictions of classical cryptographic encryption and hashing to minor noise in input data [1], [2], [3].

The binary fingerprint is typically constructed based on the dimensionality reduction followed by binarization [4]. Mostly, the soft information about the magnitudes of transformed components is neglected and some privacy amplification procedure is applied to binary data to avoid the recovery of the original data based on its binary counterpart. The overview of the state-of-the-art of privacy amplification based on encryption and randomization/compression is presented in [5], while privacy protection using data hiding approach is proposed in [6]. It is shown that the latter strategy is more efficient, when information about the fingerprint bit reliability is used in terms of achievable identification rate-privacy leak trade-off.

Contrarily to randomization/compression based privacy amplification, which blindly flips certain fraction of fingerprint bits, the privacy amplification based on data hiding uses soft

information about the bit reliability [4] to randomize only the least reliable bits while keeping the most reliable bits unchanged. Additionally, the positions of the most reliable bits in the fingerprint vector are secret and defined by the soft information that is only available to the authorized encoder/decoder pair and is not stored in the public domain.

The selection of the reliable components can be achieved based on either thresholding of magnitudes of projected components or order statistics by selecting the fraction of the largest components [6], [7]. The thresholding approach is an element-wise operation that ensures the independence of other vector components. However, it leads to the variable cardinality sets of reliable components that might represent certain challenges for the construction of practical codes. Alternatively, the order statistics approaches guarantees the fixed cardinality sets and leads to simple implementation. Since the order statistics are based on the entire vector, the resulting components can not be considered independent that should be properly analysed in the context of achievable rate-privacy leak trade-off.

We will consider the problem for Gaussian data, but the obtained results can be extended to other probability distributions.

2. PROBLEM FORMULATION

Let us introduce the following notation.

Let $w \leq n$ be a fixed integer and let $\mathbf{s} = (s_1, \dots, s_n) \in \{0, 1\}^n$ be a binary vector of the Hamming weight w , i.e.,

$$\left| \left\{ j \in \{1, \dots, n\} : s_j = 1 \right\} \right| = w.$$

Given a float-valued vector $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$, let

$$\text{bin}(\mathbf{x}) = (\text{bin}(x_1), \dots, \text{bin}(x_n)) \in \{0, 1\}^n$$

denote the binary vector constructed according with the rules

$$\text{bin}(x_j) = \begin{cases} 0, & \text{if } x_j < 0, \\ 1, & \text{if } x_j \geq 0 \end{cases}$$

for all $j = 1, \dots, n$. Furthermore, let

$$\text{abs}(\mathbf{x}) = (|x_1|, \dots, |x_n|) \in (\mathbb{R}^+)^n$$

In other words, the vectors $\text{bin}(\mathbf{x})$ and $\text{abs}(\mathbf{x})$ contain information about the signs and the magnitudes of components of the vector \mathbf{x} , respectively.

The encoder, described below, transforms the vector \mathbf{x} to a binary vector $\mathbf{b} = (b_1, \dots, b_n)$. It keeps w components of the vector $\text{bin}(\mathbf{x})$, located at positions j such that $s_j = 1$, and replaces $(n - w)$ components, located at positions such that $s_j = 0$, with arbitrary chosen bits. The result of such a

transformation is declared as a privacy protected version of the vector \mathbf{x} . The idea of privacy amplification is an assignment of the vector \mathbf{s} depending on the vector \mathbf{x} . Namely, we select positions containing w maximum magnitudes of components of the vector \mathbf{x} : if (j_1, \dots, j_n) is a sorted list of components of the vector $(1, \dots, n)$ such that

$$|x_{j_1}| \geq \dots \geq |x_{j_n}|,$$

then

$$s_j = \begin{cases} 1, & \text{if } j \in \{j_1, \dots, j_L\}, \\ 0, & \text{if } j \notin \{j_1, \dots, j_L\} \end{cases}$$

and the vector \mathbf{b} belongs to the set

$$\mathcal{B}(\mathbf{x}, \mathbf{s}) = \left\{ \mathbf{b} \in \{0, 1\}^n : \mathbf{b} \wedge \mathbf{s} = \text{bin}(\mathbf{x}) \wedge \mathbf{s} \right\}, \quad (1)$$

where \wedge denotes the component wise AND operation over a pair of binary vectors.

Example. Let $n = 7$ and $w = 2$. Then the following vectors are consistent,

$$\begin{aligned} & \begin{bmatrix} \mathbf{x} \\ \text{bin}(\mathbf{x}) \\ \mathbf{s} \\ \mathbf{b} \end{bmatrix} \\ = & \begin{bmatrix} +1.0 & -1.5 & +2.0 & +0.5 & -3.0 & +0.2 & -1.0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ * & * & 1 & * & 0 & * & * \end{bmatrix}, \end{aligned}$$

where the “*” sign denotes an arbitrary chosen bit.

The scheme under our considerations is motivated by applications to noisy identification, when the vector \mathbf{x} contains information about multimedia or a physical object. Without loss of generality we assume that the vector \mathbf{x} is extracted from a feature vector obtained from the object at the pre-processing stage. Moreover, we let the dimensionality of \mathbf{x} be smaller than the dimensionality of the original feature vector. This information has to be converted to a binary vector and stored in a database under the object’s ID at the enrolment stage. If the same object appears at the identification stage, then a vector \mathbf{y} , which represents a noisy version of the vector \mathbf{x} , is offered. In this case, the verifier has to accept the identity claim. The scheme is protected against an attacker, who wants to receive the acceptance decision by presenting an artificially constructed vector \mathbf{y} , if the probability distribution over the vectors stored in the database is uniform. In our setup, the attacker might intend to reconstruct the original vector \mathbf{x} based on the available counterpart \mathbf{b} . Since the designed scheme does not explicitly contain an external randomness, like a cryptographic key, one cannot eliminate some information leakage about the vector \mathbf{x} . Therefore, we want to design a scheme, which simultaneously maximizes the identification rate, quantified by the mutual information between X^n and its noisy observation Y^n , for a given information leak, quantified by the mutual information between X^n and its binary counterpart B^n stored in the database. We develop our analysis based on the Gaussian assumptions about the statistics of X^n and identification channel that can be extended to a class of symmetric unimodal distributions. Our choice has the following twofold justification. First, it corresponds to the statistical assumptions used in the analysis of certain practical identification schemes [6] that demonstrate high identification accuracy. Secondly, the knowledge of a fixed Hamming weight vector \mathbf{s} brings common randomness that is approximately shared between enrolment and identification stages and poses additional difficulties to the attacker in inferring information about the original vector \mathbf{x} .

The problem under our analysis is formulated as follows. Given a float-valued vector \mathbf{x} and a binary vector \mathbf{s} , introduce the probability distribution

$$\Omega(\mathbf{x}, \mathbf{s}) = \left(\Omega(\mathbf{b}|\mathbf{x}, \mathbf{s}), \mathbf{b} \in \{0, 1\}^n \right),$$

where

$$\Omega(\mathbf{b}|\mathbf{x}, \mathbf{s}) = \begin{cases} 2^{-(n-w)}, & \text{if } \mathbf{b} \wedge \mathbf{s} = \text{bin}(\mathbf{x}) \wedge \mathbf{s}, \\ 0, & \text{otherwise,} \end{cases}$$

Furthermore, let $\mathbf{s} = \mathbf{s}(\mathbf{x})$, where the vector

$$\mathbf{s}(\mathbf{x}) = (s_1(\mathbf{x}), \dots, s_n(\mathbf{x}))$$

is determined by the rules (to avoid technical difficulties, we assume that all components of the vector \mathbf{x} are different):

$$s_j(\mathbf{x}) = 1 \Leftrightarrow \left| \left\{ j' \in \{1, \dots, n\} : |x_{j'}| \geq |x_j| \right\} \right| \leq w.$$

Define the stochastic encoding of the float-valued vector $\mathbf{x} \in \mathbb{R}^n$ as generating of a binary vector $\mathbf{b} \in \{0, 1\}^n$ according with the probability distribution $\Omega(\mathbf{x}, \mathbf{s}(\mathbf{x}))$.

3. PROBABILISTIC ENSEMBLES

In the Gaussian assignments below, we use the following notation. Let

$$\text{Gaus}(x|m, \rho^2) = \frac{1}{\sqrt{2\pi\rho^2}} \exp\left\{-\frac{(x-m)^2}{2\rho^2}\right\}, \quad x \in \mathbb{R},$$

denote the probability density function (PDF) of the Gaussian distribution having the mean m and the variance ρ^2 . Furthermore, let

$$\text{erf}(a) = \frac{2}{\sqrt{\pi}} \int_{b \in (0, a)} \exp\{-b^2\}, \quad a \in \mathbb{R}.$$

Let

$$\begin{aligned} P(x) &= \text{Gaus}(x|0, \rho^2), \\ V(y|x) &= \text{Gaus}(y|x, \sigma^2) \end{aligned}$$

specify the marginal PDF of the random variable X and the conditional PDF of the random variable Y , given $X = x$. Then

$$\begin{aligned} Q(y) &= \text{Gaus}(y|0, \rho^2 + \sigma^2), \\ W(x|y) &= \text{Gaus}\left(x \mid \frac{\rho^2}{\rho^2 + \sigma^2} y, \frac{\rho^2 \sigma^2}{\rho^2 + \sigma^2}\right) \end{aligned}$$

are the marginal PDF of the random variable Y and the conditional PDF of the random variable X , given $Y = y$.

Suppose that $X^n = (X_1, \dots, X_n)$, where X_1, \dots, X_n are independent identically distributed (i.i.d.) continuous random variables generated by the PDF $(P(x), x \in \mathbb{R})$. Furthermore, let $Y^n = (Y_1, \dots, Y_n)$, where Y_1, \dots, Y_n are generated by the conditional PDF’s $(V(y|x), y \in \mathbb{R}), x \in \mathbb{R}$. Then

$$\begin{aligned} P(\mathbf{x}) &= \prod_{j=1}^n P(x_j), \\ V(\mathbf{y}|\mathbf{x}) &= \prod_{j=1}^n V(y_j|x_j). \end{aligned}$$

We also have memoryless models for generating the random vector Y^n and for generating the random vector X^n , given $Y^n = \mathbf{y}$. Namely,

$$\begin{aligned} Q(\mathbf{y}) &= \prod_{j=1}^n Q(y_j), \\ W(\mathbf{x}|\mathbf{y}) &= \prod_{j=1}^n W(x_j|y_j). \end{aligned}$$

However, the channel $X^n - B^n$, determined by the conditional PDF's $\Omega(\mathbf{x}, \mathbf{s}(\mathbf{x}))$, has memory, as the vector $\mathbf{s}(\mathbf{x})$ depends on all components of the vector \mathbf{x} . Therefore the joint probability distribution of the pair (B^n, X^n) , computed for the Markov chain

$$B^n - X^n - Y^n$$

also has memory.

We use an information-theoretic approach and characterize the scheme under considerations by the following parameters.

- The identification rate R_{id} , which is understood as the value of the mutual information function between the pair $(S^n, B^n \wedge S^n)$ and Y^n normalized by n , i.e.,

$$\begin{aligned} nR_{\text{id}} &= I(S^n, B^n \wedge S^n; Y^n) \\ &= -n \int_{y \in \mathbb{R}} Q(y) \log Q(y) - h(Y^n | S^n, B^n \wedge S^n), \end{aligned}$$

where $h(Y^n | S^n, B^n \wedge S^n)$ is the conditional differential entropy of Y^n given S^n and $B^n \wedge S^n$. The value of $2^{nR_{\text{id}}}$ is interpreted as the limit on the number of objects that can be reliably identified.

- The privacy leak, which is understood as the value of the mutual information function between B^n and Y^n normalized by n , i.e.,

$$\begin{aligned} nL_{\text{p}} &= I(B^n; X^n) \\ &= -n \int_{x \in \mathbb{R}} P(x) \log P(x) - h(X^n | B^n), \end{aligned}$$

where $h(X^n | B^n)$ is the conditional differential entropy of X^n given B^n . The value of nL_{p} characterizes the decrease of the uncertainty of an attacker about the input data after the content of the database becomes available.

We will give explicit formulas for these quantities. The idea of our analysis is an introduction of an auxiliary random variable defined as the w -th sorted magnitude of the input vector.

4. COMPUTING THE IDENTIFICATION RATE AND THE PRIVACY LEAK

For all $a \in \mathbb{R}^+$, denote

$$\varepsilon_a = \int_{|x| > a} P(x) = 1 - \text{erf}\left(\frac{a}{\sqrt{2}}\right).$$

Then

$$\varepsilon_a/2 = \int_{x > +a} P(x) = \int_{x < -a} P(x)$$

and

$$1 - \varepsilon_a = \int_{x \in (-a, +a)} P(x).$$

Let S^n be a random vector whose realizations are binary vectors of weight w chosen according to a uniform probability distribution. The entropy of S^n is equal to

$$H(S^n) = \log \binom{n}{w}.$$

Let A be a continuous random variable and let the PDF of A given $S^n = \mathbf{s}$ and $B^n = \mathbf{b}$ be defined as

$$\Psi(a | \mathbf{s}, \mathbf{b}) = 2^w \binom{w}{1} (1 - \varepsilon_a)^{n-w} (\varepsilon_a/2)^{w-1} P(a).$$

Then

$$h(A | S^n = \mathbf{s}, B^n = \mathbf{b}) = - \int_{a \in \mathbb{R}^+} \Psi(a | \mathbf{s}, \mathbf{b}) \log \Psi(a | \mathbf{s}, \mathbf{b})$$

is the differential entropy of A given $S^n = \mathbf{s}$ and $B^n = \mathbf{b}$.

One can see that the random variables S^n and A , introduced above, are also created in the random experiment, where one generates a random vector X^n according to the PDF P , finds w positions of components with maximum magnitudes, labels them by ones in the vector \mathbf{s} , and declares a as the value of the magnitude of the w -th sorted component. As a result, we conclude that S^n and A are deterministic functions of X^n and

$$h(X^n | B^n) = h(X^n, A, S^n | B^n).$$

On the other hand, by the chain rule,

$$\begin{aligned} h(X^n, A, S^n | B^n) &= H(S^n | B^n) + h(A | S^n, B^n) + h(X^n | A, S^n, B^n) \\ &= H(S^n) + h(A | S^n, B^n) + h(X^n | A, S^n, B^n). \end{aligned}$$

Given a triple $(a, \mathbf{s}, \mathbf{b})$, let us introduce a conditional PDF

$$\left(\beta(\hat{\mathbf{x}} | a, \mathbf{s}, \mathbf{b}), \hat{\mathbf{x}} \in \mathbb{R}^n \right)$$

in such a way that

$$\beta(\hat{\mathbf{x}} | a, \mathbf{s}, \mathbf{b}) = \prod_{j=1}^n \begin{cases} \beta_a^0(\hat{x}_j), & \text{if } (s_j, b_j) = (1, 0), \\ \beta_a^*(\hat{x}_j), & \text{if } s_j = 0, \\ \beta_a^1(\hat{x}_j), & \text{if } (s_j, b_j) = (1, 1), \end{cases}$$

where

$$\begin{aligned} \beta_a^0(x) &= \frac{P(x)}{\varepsilon_a/2}, \quad x < -a, \\ \beta_a^*(x) &= \frac{P(x)}{1 - \varepsilon_a}, \quad x \in (-a, +a), \\ \beta_a^1(x) &= \frac{P(x)}{\varepsilon_a/2}, \quad x > +a. \end{aligned}$$

Notice that $\beta(\hat{\mathbf{x}} | a, \mathbf{s}, \mathbf{b})$ can be equivalently introduced as

$$\beta(\hat{\mathbf{x}} | a, \mathbf{s}, \mathbf{b} \wedge \mathbf{s}) = \prod_{j=1}^n \begin{cases} \beta_a^0(\hat{x}_j), & \text{if } (s_j, b_j \wedge s_j) = (1, 0), \\ \beta_a^*(\hat{x}_j), & \text{if } (s_j, b_j \wedge s_j) = (0, 0), \\ \beta_a^1(\hat{x}_j), & \text{if } (s_j, b_j \wedge s_j) = (1, 1). \end{cases}$$

The random vector \hat{X}^n generated according with the conditional PDF β has the same structure as the random vector X^n in a sense that $\mathbf{b} \in \mathcal{B}(\hat{\mathbf{x}}, \mathbf{s})$ (see (1)), the magnitudes of components, located at positions j with $s_j = 0$, is less than a and the magnitudes of components, located at positions j with $s_j = 1$, is greater than a . Furthermore, by the symmetric properties of the functions β_a^0 and β_a^1 ,

$$\begin{aligned} h(\hat{X}^n | A = a, S^n = \mathbf{s}, B^n = \mathbf{b}) &= -w \int_{x < -a} \beta_a^0(x) \log \beta_a^0(x) \\ &\quad - (n - w) \int_{x \in (-a, +a)} \beta_a^*(x) \log \beta_a^*(x). \end{aligned}$$

The algorithm for generating the random vector in such a way that the w -th sorted magnitude of components, located

at positions j with $s_j = 1$, is exactly equal to a requires a slight modification, and

$$\begin{aligned} & h(X^n|A = a, S^n = \mathbf{s}, B^n = \mathbf{b}) \\ = & \log w + 1 - (w-1) \int_{x < -a} \beta_a^0(x) \log \beta_a^0(x) \\ & - (n-w) \int_{x \in (-a, +a)} \beta_a^*(x) \log \beta_a^*(x). \end{aligned}$$

Since $h(X^n|B^n = \mathbf{b})$ does not depend on a particular realization \mathbf{b} , we obtain

$$h(X^n|B^n = \mathbf{b}) = h(X^n|B^n).$$

To derive the expression for $h(Y^n|S^n, B^n \wedge S^n)$, we use a similar approach and obtain

$$\begin{aligned} & h(Y^n|A = a, S^n = \mathbf{s}, B^n \wedge S^n = \mathbf{b} \wedge \mathbf{s}) \\ = & \log w + 1 - (w-1) \int_{y \in \mathbb{R}} \gamma_a^0(y) \log \gamma_a^0(y) \\ & - (n-w) \int_{y \in \mathbb{R}} \gamma_a^*(y) \log \gamma_a^*(y), \end{aligned}$$

where

$$\begin{aligned} \gamma_a^0(y|a) &= \frac{1}{\varepsilon_a/2} \int_{x < -a} P(x)V(y|x), \\ \gamma_a^*(y|a) &= \frac{1}{1 - \varepsilon_a} \int_{x \in (-a, +a)} P(x)V(y|x), \\ \gamma_a^1(y|a) &= \frac{1}{\varepsilon_a/2} \int_{x > +a} P(x)V(y|x). \end{aligned}$$

5. CONCLUSIONS AND POSSIBLE EXTENSIONS

The conditional entropy $H(X^n|B^n)$ under our considerations can be also presented as

$$\begin{aligned} & h(X^n|B^n) \\ = & h(X^n | \text{bin}(X^n) \oplus K^n) \\ = & h(\text{abs}(X^n), \text{bin}(X^n) | \text{bin}(X^n) \oplus K^n) \\ = & H(\text{bin}(X^n) | \text{bin}(X^n) \oplus K^n) \\ & + h(\text{abs}(X^n) | \text{bin}(X^n), \text{bin}(X^n) \oplus K^n) \\ = & H(K^n) + h(\text{abs}(X^n) | K^n), \end{aligned}$$

where \oplus denotes the component-wise sum modulo 2 of two binary vectors and K^n is an auxiliary random binary vector chosen by the designer of the scheme according to a specified probability distribution. The last equality is due to the fact that $\text{bin}(X^n)$ and K^n are generated from the Bernoulli(1/2) source. Our random vector S^n plays the same role as K^n after we set $K^n = S^n \oplus (1, \dots, 1)$. Then $H(K^n) = H(S^n)$. However

$$h(\text{abs}(X^n) | K^n) - h(\text{abs}(X^n)) \leq 0$$

with the equality if and only if K^n is assigned independently of X^n . Our analysis allows us to study the difference above as a function of n, w , and parameters of the PDF's that specify the input data and the observation channel.

If P is an arbitrary PDF, then we replace the realization $\mathbf{x} = (x_1, \dots, x_n)$ by the vector $\mathbf{x}' = (x'_1, \dots, x'_n)$, where

$$x'_j = P(x_j) - 1/2, \quad j = 1, \dots, n$$

and keep the most of the previous considerations. However, details of calculations of the identification rate and the privacy leak should be modified, since $|P(x_j) - 1/2|$ is not a

monotone function of $|x_j - x^0|$ in general case, where x^0 is defined by the equation

$$\int_{x < x^0} P(x) = \int_{x > x^0} P(x) = 1/2.$$

Nevertheless, we believe that the developed approaches have sufficient generality and can be efficiently used in identification schemes.

6. ACKNOWLEDGMENT

The paper was partially supported by the Faculty Exchange project ‘‘Information Theoretic Analysis of Large Scale Multimedia Security and Digital Forensic Systems’’, the project ‘‘Application of Security to Biometrics and Communications’’, VW: Az.: I/85 296, and SNF project 200021–132337.

REFERENCES

- [1] J. Fridrich, ‘‘Visual hash for oblivious watermarking’’. *IST/SPIE Proceedings*, vol. 3971, San Jose, California, U.S.A., 2000.
- [2] M. K. Mihcak, R. Venkatesan, ‘‘A perceptual audio hashing algorithm: A tool for robust audio identification and information hiding’’. *Proc. 4th International Information Hiding Workshop*, Pittsburgh, PA, U.S.A. 2001.
- [3] F. Lefebvre, B. Macq, ‘‘RASH : RAdon Soft Hash algorithm’’. *Proc. EUSIPCO - European Signal Processing Conference*, Toulouse, France, 2002.
- [4] S. Voloshynovskiy, O. Koval, F. Beekhof, F. Farhadzadeh, T. Holotyak, ‘‘Information-theoretical analysis of private content identification’’. *IEEE Information Theory Workshop, ITW2010*, Dublin, Ireland, 2010.
- [5] T. Ignatenko, F. M. J. Willems, ‘‘Privacy leakage in biometric secrecy systems’’. *Proc. 46th Annual Allerton Conference on Communication, Control and Computing*, Urbana, U.S.A., pp. 850–857, 2008.
- [6] S. Voloshynovskiy, O. Koval, F. Beekhof, F. Farhadzadeh, T. Holotyak, ‘‘Private content identification based on soft fingerprinting’’. *Proc. SPIE Photonics West, Electronic Imaging, Media Forensics and Security XIII*, San Francisco, USA, 2011.
- [7] V. B. Balakirsky, A. J. Han Vinck, ‘‘Biometric authentication based on significant parameters’’. *LNCS: Biometrics and ID Management*, vol. 6583, pp. 13–22, 2011.