

Source Distinguishability under Distortion-Limited Attack: an Optimal Transport Perspective

Mauro Barni*, *Fellow, IEEE*, Benedetta Tondi, *Student Member, IEEE*

Abstract

We analyze the distinguishability of two sources in a Neyman-Pearson set-up when an attacker is allowed to modify the output of one of the two sources subject to a distortion constraint. By casting the problem in a game-theoretic framework and by exploiting the parallelism between the attacker's goal and Optimal Transport Theory, we introduce the concept of Security Margin defined as the maximum average per-sample distortion introduced by the attacker for which the two sources can be distinguished ensuring arbitrarily small, yet positive, error exponents for type I and type II error probabilities. Several versions of the problem are considered according to the available knowledge about the sources and the type of distance used to define the distortion constraint. We compute the security margin for some classes of sources and derive a general upper bound assuming that the distortion is measured in terms of the mean square error between the original and the attacked sequence.

Index Terms

Adversarial signal processing, hypothesis testing, source identification, cybersecurity, game theory, optimal transportation theory, Earth Mover Distance (EMD), Hoffman's algorithm.

M. Barni and B. Tondi are with the Department of Information Engineering and Mathematical Sciences, University of Siena, Via Roma 56, 53100 - Siena, ITALY, phone: +39 0577 234850 (int. 1005), e-mail: {barni@dii.unisi.it, benedettatondi@gmail.com}.

Source Distinguishability under Distortion-Limited Attack: an Optimal Transport Perspective

I. INTRODUCTION

Adversarial Signal Processing (Adv-SP), sometimes referred to as adversary-aware signal processing, is an emerging research field targeting the study of signal processing techniques explicitly thought to withstand the attacks of one or more adversaries aiming at system failure. Adv-SP methods can be applied to a wide variety of security-oriented applications including multimedia forensics, biometrics, digital watermarking, steganography and steganalysis, network intrusion detection, traffic monitoring, video-surveillance, just to mention a few [1]. Source identification is a common problem in Adv-SP, due to its importance in several applications. In multimedia forensics, for instance, the analyst may want to distinguish which between two sources (e.g. a photo camera and a scanner) generated a given document, or whether a document has undergone a given processing or not. In spam filtering, e-mail messages have to be classified either as spam or authentic messages. In 1-bit watermarking, the detector has to decide whether a document is watermarked or not, while it is the goal of steganalysis to distinguish between cover and stego-images. In yet other situations, the security of a system relies on the capability of distinguishing the profile of malevolent and fair users.

In [2], a game-theoretic framework is proposed to analyze the source identification problem under adversarial conditions. To be specific, [2] introduces the so called source identification game. The game is played by a Defender (D) and an Attacker (A) and is defined as follows: given two discrete memoryless sources X and Y with alphabet \mathcal{X} and probability mass functions (pmf) P_X and P_Y , and a test sequence $x^n = (x_1, x_2 \dots x_n)$, the goal of D is to decide between hypothesis H_0 that x^n has been drawn from X and hypothesis H_1 that x^n has been generated by Y . The goal of A is to take a sequence y^n generated by Y and modify it in such a way that D classifies it as being generated by X . In doing so, D must ensure that the type I error probability (usually referred to as false positive error probability P_{fp}) of deciding for H_1 when H_0 holds stays below a given threshold, whereas A has to respect a distortion constraint, limiting the amount of modifications he can introduce into y^n . The payoff of the game is the type II error probability, or false negative error probability P_{fn} , i.e., the probability of deciding for

H_0 when H_1 holds. Of course, D aims at minimizing P_{fn} , while A wishes to maximize it. The above scenario accounts for a situation in which P_X corresponds to so-to-say normal conditions and P_Y refers to an anomalous situation. It is the goal of the attacker to modify a sequence produced under anomalous conditions in such a way that the defender does not recognize that the observed system exited the normal state.

The analysis provided in [2] assumes that the defender is confined to base its analysis only on first order statistics of x^n . Under this assumption, [2] derives the asymptotic equilibrium point of the game when the length of the test sequence tends to infinity and the false positive error probability is required to tend to zero exponentially fast with decay rate at least equal to λ (λ is nothing but the error exponent of the false positive error probability). Given two pmf's P_X and P_Y , a false positive error exponent λ , and the maximum allowed distortion L_{max} , the analysis in [2] permits to determine whether, at the equilibrium, the false negative error probability P_{fn} tends to 0 or to 1 when $n \rightarrow \infty$. This, in turn, permits to define the so-called indistinguishability region $\Gamma(P_X, \lambda, L_{max})$ as the set of pmf's that can not be distinguished reliably from P_X when $n \rightarrow \infty$ due to the presence of the attacker. If $P_Y \in \Gamma(P_X, \lambda, L_{max})$, in fact, a strictly positive false negative error exponent can not be achieved and the attacker is going to *win* the game. A similar analysis is carried out in [3], [4] for a scenario in which P_X and P_Y are not known, and the statistics of the two sources are obtained through the observation of training sequences.

A. Contribution

A drawback with the analysis carried out in [2], [3], [4] is the asymmetric role of the false positive and false negative error exponents, namely λ and ε ($\varepsilon = \lim_{n \rightarrow \infty} -\frac{1}{n} \log P_{fn}$). In such works, in fact, the defender aims at ensuring a given λ , but is satisfied with any strictly positive ε . In this paper, we make a more reasonable assumption and say that the defender wins the game, i.e. he is able to distinguish between X and Y despite the presence of the adversary, if - at the equilibrium - both error probabilities tend to zero exponentially fast, regardless of the particular values assumed by the error exponents. More precisely, by mimicking Stein's lemma [5], we analyze the behavior of $\Gamma(P_X, \lambda, L_{max})$ when $\lambda \rightarrow 0$ to see whether, given a maximum allowable distortion L_{max} , it is possible for D to simultaneously attain strictly positive error exponents for the two kinds of error, hence permitting to reliably distinguish between P_X and P_Y . Having done so, we will adopt a different perspective and introduce a new distinguishability measure, called Security Margin (\mathcal{SM}), defined as the *maximum distortion allowed to the attacker, for which two sources can be reliably distinguished*. As we will see, this is a powerful concept that permits to summarize in a single quantity the distinguishability of two sources X and Y under adversarial conditions.

In order to derive our main results, we look at the optimum attacker's strategy already derived in [2] and [4] from a new perspective, i.e. by paralleling it to *optimal transport theory* [6]. Doing so, in fact, allows to derive a very intuitive and insightful interpretation of the optimum attacker's strategy, and will permit us to derive the \mathcal{SM} for a wide class of pmf's in both the discrete and the continuous case. A fast numerical algorithm for the computation of the security margin between any two discrete pmf's will also be presented.

In the framework depicted above, the main results proven in this paper can be summarized as follows

- 1) We compute the best achievable false negative error exponent for a given distortion L_{max} and for a strictly positive, yet arbitrarily small, value of the false positive error exponent λ (Theorem 3, Section IV-B);
- 2) We introduce the security margin (\mathcal{SM}) concept as the maximum allowed distortion for which two sources X and Y can be distinguished (in the adversarial setup defined in the paper) by ensuring strictly positive error exponents of the two kinds and show that \mathcal{SM} corresponds to the Earth Mover Distance (EMD) between P_X and P_Y (Definition 1, Section IV-B);
- 3) We extend the analysis to a version of the source identification game in which P_X and P_Y are known only through training sequences (Source Identification game with training data, SI_{tr}) and show that the security margin does not change despite the fact that the SI_{tr} is in general more favorable to the attacker than the SI_{ks} game where the exact statistics of the sources are perfectly known to D and A (Theorem 6, Section V-B);
- 4) By relying on some results in the field of optimal transport theory, we present a number of ways whereby the \mathcal{SM} can be computed efficiently for both discrete and continuous sources (Section VI);
- 5) We introduce a new version of the game in which the distortion constraint is expressed in terms of maximum absolute distance between the sequence y^n and the attacked sequence z^n (Theorem 7, Section VII). We then extend our analysis to the new version of the game. This is a very interesting, yet not trivial, scenario, since in many practical applications the quality of the attacked sequence is judged in terms of maximum distance (L_∞ norm), rather than in terms of average distance.

It is worth stressing that point 4) complements and generalizes some recent studies in the field of Multimedia security, namely [7], [8] regarding image counterforensics, and [9] related to perfect steganography. As a matter of fact, all the solutions proposed in those papers can be seen as particular instances of the general optimal transport problem addressed (and solved) in Section VI. Finally, point

5) relies on a generalization of the results proven in [2], [3], [4], where the analysis was restricted to the case of additive distortion measures. As we will show, such an analysis can be extended to the case of L_∞ distortion, opening the way to the application of our methodology to all the scenarios in which the distortion constraint is applied uniformly to the elements of y^n .

Some of the results presented in this paper have already been stated in [10]. With respect to [10], however, the current paper contains a complete proof of all the main theorems, the extension to the case of source identification with training data, the derivation of a fast numerical methodology to compute the security margin between any two discrete sources, and the extension of the analysis to the case of L_∞ distortion.

The rest of this paper is organized as follows. In Section II, we introduce the notation used throughout the paper, give some definitions and review some basic concepts in game theory. In Section III, we give a rigorous definition of the addressed problem and summarize the main results proven in [2]. Section IV is the core of the paper: we use optimal transport to shed new light on the addressed problem and introduce the security margin concept. In Section V, we extend the analysis to cover the case of source identification with training data. In Section VI, we derive the security margin for several classes of sources, and provide an efficient algorithm to compute it when a close form solution can not be found. Section VII extends the analysis to a situation in which the allowed distortion is defined in terms of L_∞ distance. The paper ends in Section VIII, with some conclusions and highlights for future research. The most technical proofs are given in the appendices to avoid interrupting the flow of ideas in the main body of the paper.

II. NOTATIONS AND DEFINITIONS

In this section we introduce the notation and definitions used throughout the paper. We will use capital letters to indicate discrete memoryless sources (e.g. X). Sequences of length n drawn from a source will be indicated with the corresponding lowercase letters (e.g. x^n); accordingly, x_i will denote the i -th element of a sequence x^n . The alphabet of an information source will be indicated by the corresponding calligraphic capital letter (e.g. \mathcal{X}). The probability mass function (pmf) of a discrete memoryless source X will be denoted by P_X , while the cumulative mass function will be indicated with C_X . For the sake of simplicity, the same notation will be adopted to denote the probability density function (pdf) of a continuous random variable X . The calligraphic letter \mathcal{P} will be used to indicate the class of all the probability density functions. In addition, the notation P_X will be also used to indicate the probability measure ruling the emission of sequences from a source X , so we will use the expressions $P_X(a)$ and

$P_X(x^n)$ to indicate, respectively, the probability of symbol $a \in \mathcal{X}$ and the probability that the source X emits the sequence x^n , the exact meaning of P_X being always clearly recoverable from the context wherein it is used. Finally, we will use the notation $P_X(A)$ to indicate the probability of the event A (be it a subset of \mathcal{X} or \mathcal{X}^n) under the probability measure P_X .

Our analysis relies extensively on the concepts of type and type class defined as follows (see [5] and [11] for more details). Let x^n be a sequence with elements belonging to a finite alphabet \mathcal{X} . The type P_{x^n} of x^n is the empirical pmf induced by the sequence x^n , i.e. $\forall a \in \mathcal{X}, P_{x^n}(a) = \frac{1}{n} \sum_{i=1}^n \delta(x_i, a)$, where $\delta(x_i, a) = 1$ if $x_i = a$ and zero otherwise. In the following we indicate with \mathcal{P}_n the set of types with denominator n , i.e. the set of types induced by sequences of length n . Given $P \in \mathcal{P}_n$, we indicate with $T(P)$ the type class of P , i.e. the set of all the sequences in \mathcal{X}^n having type P .

The Kullback-Leibler (KL) divergence between two distributions P and Q on the same finite alphabet \mathcal{X} is defined as:

$$\mathcal{D}(P||Q) = \sum_{a \in \mathcal{X}} P(a) \log \frac{P(a)}{Q(a)}, \quad (1)$$

where, according to usual conventions, $0 \log 0 = 0$ and $p \log p/0 = \infty$ if $p > 0$.

A. Game theory in a nutshell

A 2-player game is defined as a 4-uple $G(\mathcal{S}_1, \mathcal{S}_2, u_1, u_2)$, where $\mathcal{S}_1 = \{s_{1,1} \dots s_{1,n_1}\}$ and $\mathcal{S}_2 = \{s_{2,1} \dots s_{2,n_2}\}$ are the set of actions (usually called strategies) the first and the second player can choose from, and $u_l(s_{1,i}, s_{2,j}), l = 1, 2$, is the payoff of the game for player l , when the first player chooses the strategy $s_{1,i}$ and the second chooses $s_{2,j}$. A pair of strategies $(s_{1,i}, s_{2,j})$ is called a profile. When $u_1(s_{1,i}, s_{2,j}) + u_2(s_{1,i}, s_{2,j}) = 0$, the game is said to be a competitive (or zero-sum) game. In the set-up adopted in this paper, $\mathcal{S}_1, \mathcal{S}_2$ and the payoff functions are assumed to be known to the two players. In addition, we assume that the players choose their strategies before starting the game without knowing the strategy chosen by the other player (strategic game).

A common goal in game theory is to determine the existence of equilibrium points, i.e. profiles that in *some way* represent a *satisfactory* choice for both players [12]. The most famous equilibrium notion is due to Nash. Intuitively, a profile is a Nash equilibrium if each player does not have any interest in changing his choice assuming the other does not change his strategy. Despite its popularity, the practical meaning of Nash equilibrium is doubtful, since there is no guarantee that the players will end up playing at the equilibrium. A notion with a more practical meaning is that of *dominant equilibrium*. A strategy is said to be strictly dominant for one player if it is the best strategy for the player, regardless of the

strategy chosen by the other player. In many cases dominant strategies do not exist, however when one such strategy exists for one of the players, he will surely adopt it (at least under the assumption of rational behavior). The other players, in turn, will choose their strategies anticipating that the first player will play the dominant strategy. As a consequence, in a two-player game, if a dominant strategy exists the players have only one rational choice called the only rationalizable equilibrium of the game [13]. Games with the above property are called *dominance solvable* games.

III. THE SOURCE IDENTIFICATION GAME WITH KNOWN SOURCES

In this section, we give a rigorous definition of the problem considered in the paper. In order to make our treatment self-contained and ease the understanding of subsequent derivations, we also summarize the main results proven in [2]. With respect to [2], however, we adopt a different perspective that facilitates the interpretation of the attacker's optimal strategy as the solution of an optimal transport problem. As a matter of fact, this can be considered as an important contribution of this paper, since the new perspective opens the way to the adoption of a new, more insightful, methodology to analyze the structure of the game and the achievable performance.

A. Definition of the SI_{ks} game and equilibrium point

We start with the definition of the source identification game with known sources (SI_{ks}). Given a test sequence x^n , we indicate with H_0 the hypothesis that x^n has been generated by P_X and with H_1 the alternative hypothesis that x^n has been generated by P_Y . In order to define the SI_{ks} game, we need to define the set of strategies of D and A and the payoff function.

Defender's strategies. The set of strategies of the Defender (\mathcal{S}_D) consists of all possible acceptance regions for H_0 . More precisely, by following [2], we require that D bases its analysis only on the first order statistics of x^n . This is equivalent to ask that the acceptance region for hypothesis H_0 , hereafter referred to as Λ^n , is a union of type classes¹. Since a type class is univocally defined by the empirical pmf of the sequences it contains, Λ^n can be seen as a union of types $P \in \mathcal{P}_n$. We consider an asymptotic version of the game and require that the false positive error probability P_{fp} decreases exponentially with decay rate at least equal to λ . Under the above assumptions, the space of strategies of D is given by:

$$\mathcal{S}_D = \{\Lambda^n \in 2^{\mathcal{P}_n} : P_{fp} \leq 2^{-\lambda n}\}, \quad (2)$$

where $2^{\mathcal{P}_n}$ indicates the power set of \mathcal{P}_n .

¹We use the superscript n to indicate explicitly that Λ^n refers to n -long sequences.

Attacker's strategies. Given a sequence y^n drawn from Y , the goal of A is to transform it into a sequence z^n belonging to the acceptance region chosen by D. Let us indicate by $n(i, j)$ the number of times that the i -th symbol of the alphabet is transformed into the j -th one as a consequence of the attack. Similarly, we indicate by $S_{YZ}^n(i, j) = n(i, j)/n$ the relative frequency with which the i -th symbol of the alphabet is transformed into the j -th one. In the following, we refer to S_{YZ}^n as *transportation map*. Once again, we explicitly indicate that S_{YZ}^n refers to n -long sequences by adding the superscript n . For any additive distortion measure, the overall distortion introduced by the attack can be expressed in terms of $n(i, j)$; in fact we have:

$$d(y^n, z^n) = \sum_{i,j} n(i, j)d(i, j), \quad (3)$$

where $d(i, j)$ is the distortion introduced when the symbol i is transformed into the symbol j . Similarly, the average per-sample distortion depends only on S_{YZ}^n :

$$\frac{d(y^n, z^n)}{n} = \sum_{i,j} S_{YZ}^n(i, j)d(i, j). \quad (4)$$

S_{YZ}^n determines also the empirical pmf (i.e. the type) of the attacked sequence. In fact, by indicating with $P_{z^n}(j)$ the relative frequency of symbol j into z^n , we have:

$$P_{z^n}(j) = \sum_i S_{YZ}^n(i, j) \triangleq S_Z^n(j). \quad (5)$$

Finally, we observe that the attacker can not change more symbols than there are in the sequence y^n ; as a consequence a map S_{YZ}^n can be applied to a sequence y^n only if:

$$S_Y^n(i) \triangleq \sum_j S_{YZ}^n(i, j) = P_{y^n}(i). \quad (6)$$

Equations (5) and (6) suggest an interesting interpretation of S_{YZ}^n , which can be seen as the joint empirical pmf between the sequences y^n and z^n . In the same way, S_Y^n and S_Z^n correspond, respectively, to the empirical pmf of y^n and z^n .

By remembering that Λ^n depends only on the empirical pmf of the test sequence (i.e., on its type), and given that the empirical pmf of the attacked sequence depends on S_Z^n only through S_{YZ}^n , we can define the action of the attacker as the choice of a transportation map among all *admissible* maps, a map being admissible if:

$$S_Y^n = P_{y^n} \quad (7)$$

$$\sum_{i,j} S_{YZ}^n(i, j)d(i, j) \leq L_{max},$$

where the second condition expresses the per-letter distortion constraint the attacker is subject to, and L_{max} is the maximum allowable (average) per-letter distortion. In the following, we will refer to the set of admissible maps as $\mathcal{A}^n(L_{max}, P_{y^n})$. With the above definitions, the space of strategies of the attacker is the set of all the possible ways of associating an admissible transformation map to the to-be-attacked sequence. In the following, we will refer to the result of such an association as $S_{YZ}^n(y^n)$, or $S_{YZ}^n(i, j; y^n)$, when we need to refer explicitly to the relative frequency with which the symbol i is transformed into the symbol j . In the same way, $S_Z^n(j; y^n)$ indicates the output marginal of $S_{YZ}^n(i, j; y^n)$ ². By adopting the above symbolism, the space of strategies for the attacker can be defined as:

$$\mathcal{S}_A = \{S_{YZ}^n(i, j; y^n) : S_{YZ}^n(i, j) \in \mathcal{A}^n(L_{max}, P_{y^n})\}. \quad (8)$$

The payoff. Having fixed the maximum false positive error probability, we adopt a typical Neyman-Pearson approach and let the payoff correspond to the false negative error probability, that is:

$$u_D = -u_A = - \sum_{y^n: S_Z^n(j; y^n) \in \Lambda^n} P_Y(y^n), \quad (9)$$

where $P_Y(y^n)$ is the probability that the source Y outputs the sequence y^n .

Equilibrium point. Given the above formulation of the SI_{ks} game, the main result of [2] is summarized by the following theorem.³

Theorem 1. *Let*

$$\Lambda^{n,*} = \left\{ P \in \mathcal{P}_n : \mathcal{D}(P||P_X) < \lambda - |\mathcal{X}| \frac{\log(n+1)}{n} \right\}, \quad (10)$$

and

$$S_{YZ}^{n,*}(i, j; y^n) = \arg \min_{S_{YZ}^n \in \mathcal{A}^n(L_{max}, P_{y^n})} \mathcal{D}(S_Z^n || P_X). \quad (11)$$

Then $\Lambda^{n,}$ is a dominant equilibrium for D and the profile $(\Lambda^{n,*}, S_{YZ}^{n,*}(i, j; y^n))$ is the only rationalizable equilibrium of the SI_{ks} game, which, then, is a dominance solvable game.*

²With regard to the input marginal, of course, we always have $S_Y^n(i; y^n) = P_{y^n}(i) \forall i$.

³In this paper we use a different formulation of the theorem with respect to [2] so to adapt it to the new formalism based on the concept of transportation map adopted here.

B. Payoff of the SI_{ks} game at the equilibrium

Given the optimal acceptance region $\Lambda^{n,*}$ and the optimum attacking strategy $S_{YZ}^{n,*}(y^n)$, we can introduce the indistinguishability region $\Gamma^n(P_X, \lambda, L_{max})$ as follows:

$$\begin{aligned} \Gamma^n(P_X, \lambda, L_{max}) = & \\ & \{P \in \mathcal{P}_n : \exists S_{YZ}^n \in \mathcal{A}^n(L_{max}, P) \text{ s.t. } S_Z^n \in \Lambda^{n,*}\}. \end{aligned} \quad (12)$$

The indistinguishability region defines all the type classes (with denominator n) whose sequences can be moved within $\Lambda^{n,*}$ by the attacker. The problem with the above analysis is that it applies only to types with denominator n and hence can not be used to decide whether the sequences generated by two generic sources (not necessarily belonging to \mathcal{P}_n) can be distinguished. In order to answer this question, we can rely on the density of rational numbers in \mathbb{R} , and let n tend to infinity. In this way we can define the asymptotic counterpart of Γ^n , specifying whether two sources can eventually be distinguished for increasing values of n [2]:

$$\begin{aligned} \Gamma(P_X, \lambda, L_{max}) = & \\ & \{P \in \mathcal{P} : \exists S_{YZ} \in \mathcal{A}(L_{max}, P) \text{ s.t. } S_Z \in \Lambda^*(P_X, \lambda)\}, \end{aligned} \quad (13)$$

where

$$\Lambda^*(P_X, \lambda) = \{P \in \mathcal{P} : \mathcal{D}(P||P_X) \leq \lambda\}, \quad (14)$$

and where the definitions of $S_{YZ}(i, j)$, $S_Z(j)$ and $\mathcal{A}(L_{max}, P)$ are obtained immediately from those of $S_{YZ}^n(i, j)$, $S_Z^n(j)$ and $\mathcal{A}^n(L_{max}, P)$, by relaxing the requirement that $S_{YZ}(i, j)$, $S_Z(j)$ and $P(i)$ are rational numbers with denominator n . More precisely, we can state the following theorem:

Theorem 2. *For the SI_{ks} game, the error exponent of the false negative error probability at the equilibrium is given by⁴:*

$$\varepsilon = \min_{P \in \Gamma(P_X, \lambda, L_{max})} \mathcal{D}(P||P_Y), \quad (15)$$

leading to the following cases:

- 1) $\varepsilon = 0$, if $P_Y \in \Gamma(P_X, \lambda, L_{max})$;
- 2) $\varepsilon \neq 0$, if $P_Y \notin \Gamma(P_X, \lambda, L_{max})$.

⁴Here and in the rest of the paper, the use of the minimum instead of the infimum is justified by the compactness of $\Gamma(P_X, \lambda, L_{max})$ and other similar sets defined in the following.

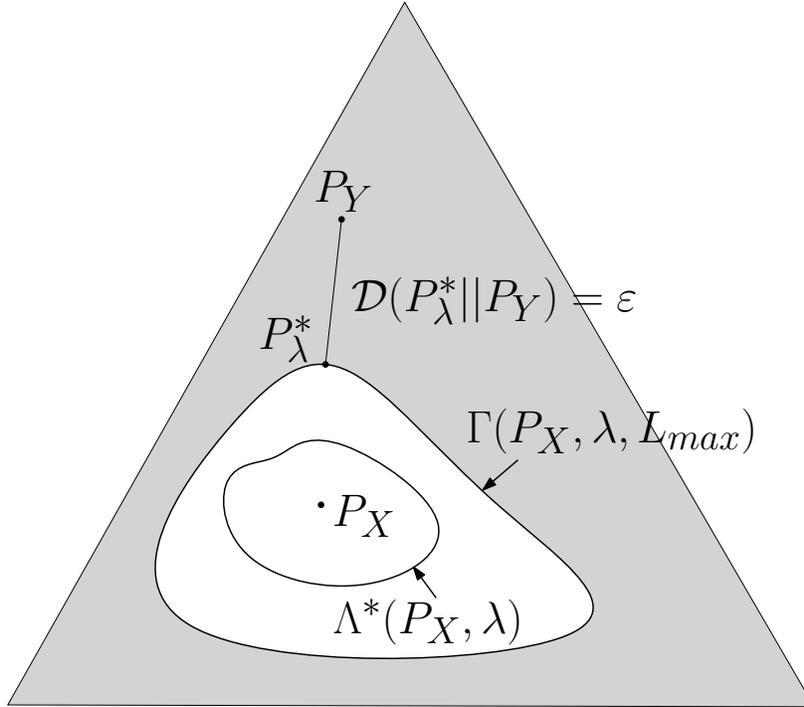


Fig. 1. Geometric interpretation of $\Gamma(P_X, \lambda, L_{max})$ and $\Lambda^*(P_X, \lambda)$ by the light of Theorem 2.

Given two pmf's P_X and P_Y , a maximum distortion L_{max} and the desired false positive error exponent λ , Theorem 2 permits to understand whether D may ever succeed to make the false negative error probability vanishingly small and thus *win* the game. Then, $\Gamma(P_X, \lambda, L_{max})$ can be interpreted as the region with the sources that cannot be reliably distinguished from P_X guaranteeing a false positive error exponent at least equal to λ in the presence of an adversary with allowed distortion L_{max} , where by *reliably distinguished* we mean distinguished in such a way to grant a strictly positive error exponent for P_{fn} . A geometric interpretation of Theorem 2 is given in Figure 1.

IV. THE SECURITY MARGIN

In this section, we use the optimal transport interpretation of the attacker's strategy to introduce a measure of source distinguishability in the set-up defined by the SI_{ks} game.

A. Characterization of the indistinguishability region using Optimal Transportation

To start with, we find it convenient to rephrase the results described in the previous section as an optimal transport problem [6].

Let P and Q be two pmf's defined over the same finite alphabet, and let $c(i, j)$ be the cost of transporting the i -th symbol into the j -th one. In one of its instances, optimal transport theory looks for the transportation map that transforms P into Q by minimizing the average cost of the transport. By using the notation introduced in the previous section, this corresponds to solving the following minimization problem:

$$\min_{S_{YZ}: S_Y=P, S_Z=Q} \sum_{i,j} S_{YZ}(i, j)c(i, j). \quad (16)$$

A nice interpretation of the problem defined by equation (16) is obtained by interpreting the pmf's P and Q as two different ways of piling up a certain amount of earth, and $c(i, j)$ as the cost necessary to move a unitary amount of earth from position i to position j . In this case, the minimum cost achieved in (16) can be seen as the minimum effort required to turn one pile into the other. Due to such a viewpoint, in computer vision applications, the minimum in equation (16) is usually known as Earth Mover Distance (*EMD*) between P and Q , [14]. However, while the definition of the *EMD* given in [14] refers in general to signatures (non-normalized distributions with unequal masses), here the pilings of earth P and Q are probability mass functions. In this case, when $c(i, j) = d(i, j)^p$ for some distance measure d (with $p \geq 1$), the *EMD* has a more general statistical meaning. Given two random variables with probability distributions P_X and P_Y , the *EMD* between P_X and P_Y corresponds to the minimum expected p -th power distance between the random variables X and Y taken over all joint probability distributions P_{XY} with marginal distributions respectively equal to P_X and P_Y :

$$EMD_{d^p}(P_X, P_Y) = \min_{P_{XY}: \sum_y P_{XY} = P_X, \sum_x P_{XY} = P_Y} E_{XY}[d(X, Y)^p]. \quad (17)$$

In transport theory terminology, expression (17) is the p -th power of the Wasserstein distance [15], [6] (or the Monge-Kantorovich metric of order p [16], [17]). In particular, when $c(i, j) = |i - j|^2$ (i.e. $d(i, j) = |i - j|$ and $p = 2$) the earth mover distance is equivalent to the squared Mallows distance between P_X and P_Y [18], that is

$$EMD_{L_2^2}(P_X, P_Y) = \min_{P_{XY}: \sum_y P_{XY} = P_X, \sum_x P_{XY} = P_Y} E_{XY}[|X - Y|^2]. \quad (18)$$

In the following, we will continue to refer to (16) as $EMD(P, Q)$. We also observe that even if we introduced the *EMD* by considering finite-alphabet sources, there is no need to restrict the definition in (17) and (18) to discrete random variables. In fact, in the second part of the paper, we will extend our analysis and use the *EMD* to measure the distinguishability of continuous sources.

Optimal transport theory permits us to rewrite the indistinguishability region in a more compact and easier-to-interpret way. In fact, it is immediate to see that equation (13) can be rewritten as:

$$\Gamma(P_X, \lambda, L_{max}) = \{P \in \mathcal{P} : \exists Q \in \Lambda^*(P_X, \lambda) \text{ s.t. } EMD(P, Q) \leq L_{max}\}, \quad (19)$$

where in the definition of the *EMD* $c(i, j)$ corresponds to the distortion metric used to constraint the strategies available to the attacker.

B. Security Margin definition

We now study the behavior of $\Gamma(P_X, \lambda, L_{max})$ when $\lambda \rightarrow 0$. Doing so will allow us to investigate whether two sources X and Y are ultimately distinguishable in the setting defined by the SI_{ks} game.

The rationale behind our analysis derives directly from equations (13) and (14). In fact, it is easy to see that decreasing λ in the definition of \mathcal{S}_D leads to a more favorable game for the defender, since he can adopt a smaller acceptance region and obtain a larger payoff. Stated in another way, from D's perspective, evaluating the behavior of the game for $\lambda \rightarrow 0$ corresponds to exploring the best achievable false negative error exponent, when P_{fp} tends to 0 exponentially fast.

More formally, we start by proving the following property.

Property 1. *For any two values λ_1 and λ_2 such that $\lambda_2 < \lambda_1$, $\Gamma(P_X, \lambda_2, L_{max}) \subseteq \Gamma(P_X, \lambda_1, L_{max})$.*

Proof: The property follows immediately from equation (19) by observing that $\Gamma(P_X, \lambda, L_{max})$ depends on λ only through the acceptance region $\Lambda^*(P_X, \lambda)$, for which we obviously have $\Lambda^*(P_X, \lambda_2) \subseteq \Lambda^*(P_X, \lambda_1)$ whenever $\lambda_2 < \lambda_1$. ■

Thanks to Property 1, we can compute the limit of the false negative error exponent when λ tends to zero, as summarized in the following theorem (somewhat resembling Stein's Lemma [5]).

Theorem 3. *Given two sources $X \sim P_X$ and $Y \sim P_Y$ and a maximum average per-letter distortion L_{max} (defined according to an additive distortion measure), let us adopt the following definition:*

$$\Gamma(P_X, L_{max}) = \{P \in \mathcal{P} : EMD(P, P_X) \leq L_{max}\}; \quad (20)$$

then the maximum achievable false negative error exponent ε for the SI_{ks} game is

$$\lim_{\lambda \rightarrow 0} \lim_{n \rightarrow \infty} -\frac{1}{n} \log P_{fn} = \min_{P \in \Gamma(P_X, L_{max})} \mathcal{D}(P || P_Y). \quad (21)$$

Proof: The innermost limit in (21) defines the error exponent for a fixed λ , say it $\varepsilon(\lambda)$. Thanks to equation (15), we know that

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log P_{fn} = \varepsilon(\lambda) = \min_{P \in \Gamma(P_X, \lambda, L_{max})} \mathcal{D}(P||P_Y). \quad (22)$$

Then, according to Property 1, the sequence $\varepsilon(\lambda)$ is monotonically non decreasing as λ decreases. In addition, since $\Gamma(P_X, L_{max}) \subseteq \Gamma(P_X, \lambda, L_{max}) \forall \lambda$, for any $\lambda > 0$, we have:

$$\varepsilon(\lambda) \leq \min_{P \in \Gamma(P_X, L_{max})} \mathcal{D}(P||P_Y). \quad (23)$$

Being $\varepsilon(\lambda)$ bounded from above and non-decreasing, the limit for $\lambda \rightarrow 0$ exists and is finite. We must now prove that the limit is indeed equal to $\min_{P \in \Gamma(P_X, L_{max})} \mathcal{D}(P||P_Y)$. Let P_0^* be the point achieving the minimum in (21) and P_λ^* the point achieving the minimum on the set $\Gamma(P_X, \lambda, L_{max})$, i.e. the point achieving the minimum in equation (15) (see Figure 1 for a pictorial representation of P_λ^*). Due to Lemma 1 (Appendix A), for any arbitrarily small τ , we can choose a small enough λ such that, for any P in $\Gamma(P_X, \lambda, L_{max})$, a pmf P' in $\Gamma(P_X, L_{max})$ exists whose distance from P is lower than τ . By taking $P = P_\lambda^*$ and exploiting the continuity of the \mathcal{D} function, we have

$$\mathcal{D}(P'||P_Y) \leq \min_{P \in \Gamma(P_X, \lambda, L_{max})} \mathcal{D}(P||P_Y) + \delta(\tau), \quad (24)$$

for some $P' \in \Gamma(P_X, L_{max})$ and some value $\delta(\tau)$ such that $\delta(\tau) \rightarrow 0$ as $\tau \rightarrow 0$. A fortiori, relation (24) holds for $P' = P_0^*$ and then we can write

$$\begin{aligned} \varepsilon(\lambda) &= \min_{P \in \Gamma(P_X, \lambda, L_{max})} \mathcal{D}(P||P_Y) \\ &\geq \min_{P \in \Gamma(P_X, L_{max})} \mathcal{D}(P||P_Y) - \delta(\tau). \end{aligned} \quad (25)$$

where $\delta(\tau)$ can be made arbitrarily small by decreasing λ . Equation (25), together with equation (23), shows that we can get arbitrarily close to $\min_{P \in \Gamma(P_X, L_{max})} \mathcal{D}(P||P_Y)$, by making λ small enough, hence proving that $\min_{P \in \Gamma(P_X, L_{max})} \mathcal{D}(P||P_Y)$ is the limit of the sequence $\varepsilon(\lambda)$ as $\lambda \rightarrow 0$. ■

Figure 2 gives a geometric interpretation of Theorem 3. The figure is obtained from Figure 1 by observing that when $\lambda \rightarrow 0$ the optimum acceptance region collapses into the single pmf P_X , i.e., $\Lambda^* = \{P_X\}$.

By the light of Theorem 3, $\Gamma(P_X, L_{max})$ is the smallest indistinguishability region for the SI_{ks} game. Moreover, from equation (20), we see that the distinguishability of two pmf's (in the SI_{ks} setting) ultimately depends on their EMD . In fact, if $EMD(P_Y, P_X) > L_{max}$, the defender is able to distinguish X from Y by adopting a sufficiently small λ . On the contrary, if $EMD(P_Y, P_X) \leq L_{max}$, there is no positive value of λ for which the sequences emitted by the two sources can be asymptotically distinguished.

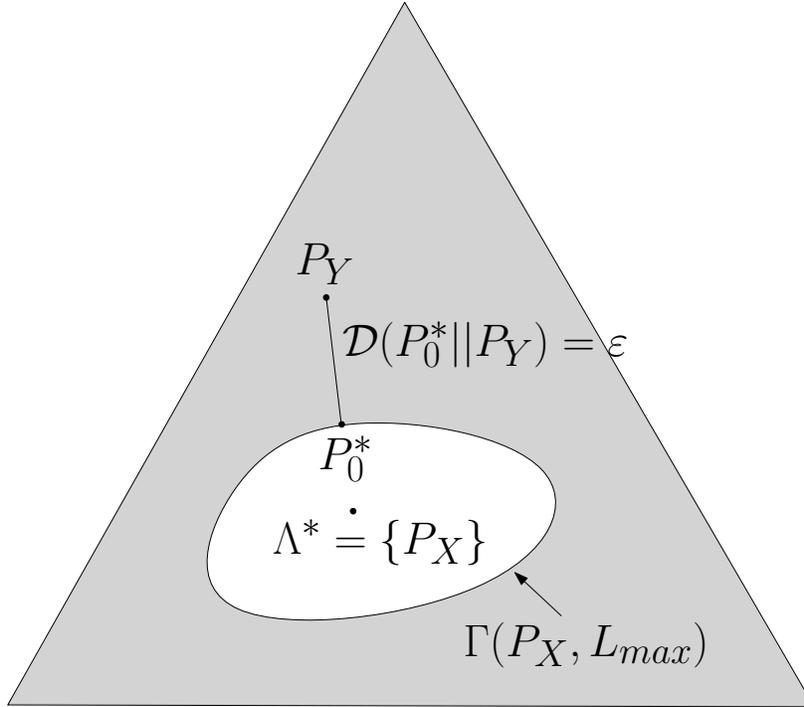


Fig. 2. Geometric interpretation of $\Gamma(P_X, L_{max})$ and P_0^* by the light of Theorem 3.

By adopting a different perspective, given two sources X and Y , one may ask which is the maximum attacking distortion for which D can distinguish X and Y despite the presence of the adversary. The answer to this question follows immediately from Theorem 3 and leads naturally to the following definition.

Definition 1 (Security Margin). *Let $X \sim P_X$ and $Y \sim P_Y$ be two discrete memoryless sources. The maximum average per-letter distortion for which the two sources can be reliably distinguished in the SI_{ks} setting is called Security Margin and is given by*

$$\mathcal{SM}(P_Y, P_X) = \text{EMD}(P_Y, P_X). \quad (26)$$

Interestingly, the *EMD* is a symmetric function of P_X and P_Y [14], and hence the security margin does not depend on the role of X and Y in the test, i.e. $\mathcal{SM}(P_X, P_Y) = \mathcal{SM}(P_Y, P_X)$. The security margin is a powerful measure summarizing in a single quantity how securely two sources can be distinguished (in the SI_{ks} setting).

It is worth remarking that the security margin between two sources pertains to the *security* of the hypothesis test behind the source identification problem and not to its *robustness*, since it is measured at

the equilibrium of the game, i.e. by assuming that both the players of the game make optimal choices. To better exemplify the above concept, let us consider the simple case of two binary sources. Specifically, let X and Y be two Bernoulli sources with parameters $p = P_X(1)$ and $q = P_Y(1)$ respectively. Let also assume that the distortion constraint is expressed in terms of the Hamming distance between the sequences, that is $d(i, j) = 0$ when $i = j$ and 1 otherwise. Without loss of generality let $p > q$. The distortion associated to a transportation map S_{XY} can be written as:

$$\sum_{i,j} S_{YX}(i, j)d(i, j) = S_{YX}(0, 1) + S_{YX}(1, 0). \quad (27)$$

Since $p > q$, it is easy to conclude that the minimum of the above expression is obtained when $S_{YX}(1, 0) = 0$ (intuitively, if the source X outputs more 1's than Y , it does not make any sense to turn the 1's emitted by Y into 0's). As a consequence, to satisfy the constraint $S_X(1) = p$ we must let $S_{YX}(0, 1) = p - q$, yielding $\mathcal{SM}(P_Y, P_X) = p - q$, or more generally $|p - q|$. We can conclude that if the attacker is allowed to introduce an average Hamming distortion larger or equal than $|p - q|$, then there is no way for the defender to distinguish between the two sources. This is not the case if the output of the source Y passes through a binary symmetric channel with crossover probability equal to $|p - q|$, since the output of the channel will still be distinguishable from the sequences emitted by X . Consider, for example, a simple case in which $q = 1/2$ and $p > 1/2$. Regardless of the crossover probability, the output of the channel will still be a binary source with equiprobable symbols, which is distinguishable from X given that $p > 1/2$. In other words, in the set up defined by the SI_{ks} game, the two sources can not be distinguished securely in the presence of an attacker introducing a distortion equal to $|p - q|$, while they can be distinguished even if the output of the source Y passes through a noisy channel introducing the same average distortion introduced by the attacker.

V. EXTENSION TO SOURCE IDENTIFICATION WITH TRAINING DATA

In this section we extend the previous analysis to the case of source identification with training data (SI_{tr}), in order to provide a measure of source distinguishability, in the more general setup studied in [4]. In such a scenario, the two sources X and Y are not completely known to D and A, so they must base their actions on the knowledge of a training sequence drawn from X (the source under the null hypothesis). This is a very interesting scenario bringing the analysis closer to real applications, in which a precise statistical model of the to-be-distinguished sources is usually not available. In [4], it is proven that the source identification game with training data is more favorable to the attacker than the SI_{ks} game. Then one could argue that in the SI_{tr} setup the security margin between the two sources is smaller,

implying that a lower distortion is sufficient to the attacker to make the sources undistinguishable. The remarkable result that we will prove in this section is that this is not the case, hence showing that the ultimate distinguishability of two sources is the same for the two games.

A. The source identification game with training data (SI_{tr})

In order to present our analysis in a self-contained way, in this section we summarize the main results proven in [4]. Once again, we will do so by adopting a transportation theory perspective for the definition of the attacker's optimum strategy.

Let us start by giving a rigorous definition of the source identification game with training data.

Defender's strategies. In the SI_{tr} game the defender must decide whether a test sequence x^n has been generated by a source X with unknown pmf by relying on the knowledge of an N -sample training sequence t_D^N drawn from X . This is equivalent to deciding whether to accept or not the hypothesis H_0 that the test and the training sequences have been generated by the same source. In this framework, the acceptance region Λ is defined as the set with all the pairs of sequences (x^n, t_D^N) that D classifies as being generated by the same source. Once again, we limit the action of D to a first order analysis of x^n and t_D^N . This is equivalent to require that the acceptance region for hypothesis H_0 is a union of pairs of type classes, or equivalently, pairs of types (P, Q) , where $P \in \mathcal{P}_n$ and $Q \in \mathcal{P}_N$. As for the SI_{ks} case, the defender must ensure that the asymptotic false positive error probability tends to zero exponentially fast at least with a certain decay rate, however since P_X is not known, the constraint must be satisfied in a worst case sense, i.e. for all possible choices of P_X . More specifically, the space of strategies of D is given by:

$$\mathcal{S}_D = \{\Lambda_{tr}^n \subset \mathcal{P}_n \times \mathcal{P}_N : \max_{P_X \in \mathcal{P}} P_{fp} \leq 2^{-\lambda n}\}, \quad (28)$$

where \mathcal{P} is the class of discrete memoryless sources.

Attacker's strategies. Given a sequence y^n drawn from a source $Y \neq X$, the goal of A is to transform y^n into a sequence z^n belonging to the acceptance region chosen by D while respecting a distortion constraint. Likewise the defender, all the information that the attacker has about X is a K -long training sequence t_A^K . By using the same transportation theoretic formalism used in the previous section, the set of strategies of the attacker consists of all the possible ways of choosing an admissible transportation map to transform y^n into z^n .

$$\mathcal{S}_A = \{S_{YZ}^n(i, j; y^n, t_A^K) : S_{YZ}^n(i, j) \in \mathcal{A}^n(L_{max}, P_{y^n})\}, \quad (29)$$

where we have explicitly indicated that the choice of the transportation map now depends also on t_A^K , and where the set of admissible maps is defined as in the SI_{ks} case.

Depending on the relationship between t_A^K and t_D^N , several versions of the SI_{tr} game can be defined. Here we focus on the simplest case of equal training sequences, i.e. we assume $K = N$ and $t_A^N = t_D^N \triangleq t^N$. We will see later on that our analysis can be easily extended so to cover the other cases addressed in [4]. In addition, we force N to be a linear function of n with some proportionality constant c , i.e. $N = cn$. As discussed in [4], this is the most significant case to study.

The payoff. Adopting again the Neyman-Pearson approach, the payoff corresponds to the false negative error probability, that is:

$$u_D = -u_A = - \sum_{\substack{(y^n, t^N) \in \mathcal{X}^n \times \mathcal{X}^N: \\ (S_Z^n(j; y^n, t^N), t^N) \in \Lambda_{tr}^n}} P_Y(y^n) P_X(t^N), \quad (30)$$

where $P_X(t^N)$ is the probability that the source X outputs the sequence t^N and Λ_{tr}^n is the acceptance region of the test.

Equilibrium point. The derivation of the optimum strategy for D passes through the definition of the generalized log-likelihood ratio function $h(P_{x^n}, P_{t^N})$ defined as ([19], [20], [4]):

$$h(P_{x^n}, P_{t^N}) = \mathcal{D}(P_{x^n} || P_{r^{n+N}}) + c\mathcal{D}(P_{t^N} || P_{r^{n+N}}), \quad (31)$$

where $P_{r^{n+N}}$ indicates the empirical pmf of the sequence r^{n+N} , obtained by concatenating x^n and t^N . The main result of [4] is summarized by the following theorem.

Theorem 4. *Let*

$$\Lambda_{tr}^{n,*} = \{(P, Q) \in \mathcal{P}_n \times \mathcal{P}_N : h(P, Q) < \lambda - \kappa(n, c)\}, \quad (32)$$

$$S_{YZ}^{n,*}(i, j; y^n, t^N) = \arg \min_{S_{YZ}^n \in \mathcal{A}^n(D_{max}, P_{y^n})} h(S_Z^n, P_{t^N}). \quad (33)$$

where $\kappa(n, c) = |\mathcal{X}|^{\frac{\log(n+1)(N+1)}{n}}$. Then $\Lambda_{tr}^{n,*}$ is a dominant equilibrium for D and the profile $(\Lambda_{tr}^{n,*}, S_{YZ}^{n,*}(i, j; y^n, t^N))$ is the only rationalizable equilibrium of the SI_{tr} game with equal training sequences, which, then, is a dominance solvable game [13].

As for the SI_{tr} game, by letting n tend to infinity and by exploiting the density of rational numbers in the real line, we can study the asymptotic distinguishability of sequences emitted by any two sources. To express the final result of the above procedure, we need to introduce some definitions. First of all we

need to extend the h function so to make it work on general pmf's. We let:

$$\begin{aligned} h_c(P, Q) &= \mathcal{D}(P||U) + c\mathcal{D}(Q||U); \\ U &= \frac{1}{1+c}P + \frac{c}{1+c}Q, \end{aligned} \quad (34)$$

which permits us to define the following sets:

$$\Lambda_{tr}^*(Q, \lambda) = \{P \in \mathcal{P} : h_c(P, Q) \leq \lambda\}, \quad (35)$$

and

$$\begin{aligned} \Gamma_{tr}(Q, \lambda, L_{max}) &= \{P \in \mathcal{P} : \exists R \in \Lambda_{tr}^*(Q, \lambda) \\ &\text{s.t. } EMD(P, R) \leq L_{max}\}. \end{aligned} \quad (36)$$

The following theorem, proved in [4], states that the indistinguishability region of the SI_{tr} game is given by $\Gamma_{tr}(P_X, \lambda, L_{max})$, where P_X is the true distribution of the source X .

Theorem 5. *For the SI_{tr} game with equal training sequences available to the players, the error exponent of the false negative error probability at the equilibrium is given by:*

$$\varepsilon_{tr}(\lambda) = \min_R \left[c \cdot \mathcal{D}(R||P_X) + \min_{P \in \Gamma_{tr}(R, \lambda, L_{max})} \mathcal{D}(P||P_Y) \right] \quad (37)$$

leading to the following cases:

- 1) $\varepsilon_{tr}(\lambda) = 0$, if $P_Y \in \Gamma_{tr}(P_X, \lambda, L_{max})$;
- 2) $\varepsilon_{tr}(\lambda) \neq 0$, if $P_Y \notin \Gamma_{tr}(P_X, \lambda, L_{max})$.

From the above theorem we see that the sources that cannot be asymptotically distinguished from P_X are those inside $\Gamma_{tr}(P_X, \lambda, L_{max})$. The geometrical interpretation is similar to the one given in Figure 1 for Theorem 2 where now the acceptance region is given by $\Lambda_{tr}^*(P_X, \lambda)$ and the indistinguishability region is $\Gamma_{tr}(P_X, \lambda, L_{max})$.

We point out that the only difference with respect to the case of known sources consists in the asymptotic acceptance region $\Lambda_{tr}^*(P_X, \lambda)$, which is strictly larger than $\Lambda^*(P_X, \lambda)$, given that h_c function is always lower than \mathcal{D} (see [4] for the proof). As a consequence, it is straightforward to argue that $\Gamma_{tr}(P_X, \lambda, L_{max}) \supset \Gamma(P_X, \lambda, L_{max})$.

B. Security margin for the SI_{tr} game

We now study the behavior of the SI_{tr} game when $\lambda \rightarrow 0$ so to investigate the best achievable performance for the defender in the case of training-based decision. To start with, we observe that the divergence and the h_c function share a similar behavior, in that they are convex functions and both $\mathcal{D}(P||Q)$ and $h_c(P, Q)$ are equal to zero if and only if $P = Q$. This permits to extend Property 1 to the set Γ_{tr} yielding:

Property 2. For any two values λ_1 and λ_2 such that $\lambda_2 < \lambda_1$, $\Gamma_{tr}(P_X, \lambda_2, L_{max}) \subseteq \Gamma_{tr}(P_X, \lambda_1, L_{max})$.

In a similar way, Lemma 1 can be extended to the set $\Gamma_{tr}(R, \lambda, L_{max})$ (Appendix A).

We are now ready to prove the counterpart of Theorem 3 for the SI_{tr} game.

Theorem 6. Given two sources $X \sim P_X$ and $Y \sim P_Y$ and a maximum allowable average per-letter distortion L_{max} (defined according to an additive distortion measure), the maximum achievable false negative error exponent for the SI_{tr} game is

$$\lim_{\lambda \rightarrow 0} \varepsilon_{tr}(\lambda) = \min_R [c \cdot \mathcal{D}(R||P_X) + \min_{P \in \Gamma(R, L_{max})} \mathcal{D}(P||P_Y)], \quad (38)$$

where $\Gamma(R, L_{max})$ is defined as in (20) by replacing P_X with R^5 .

Proof: The proof goes along the same line of the proof of Theorem 3. From Property 2, we see immediately that $\varepsilon(\lambda)$ is non-increasing when λ decreases, since the innermost minimization in equation (37) is taken over a smaller set when λ decreases. Then, by the same token, we have:

$$\varepsilon_{tr}(\lambda) \leq \min_R (c\mathcal{D}(R||P_X) + \min_{P \in \Gamma(R, D_{max})} \mathcal{D}(P||P_Y)). \quad (39)$$

This implies that $\lim_{\lambda \rightarrow 0} \varepsilon(\lambda)$ exists and is finite. Given that Lemma 1 still holds for the set $\Gamma_{tr}(R, \lambda, L_{max}) \forall R$, we can reason as in the proof of Theorem 3 to conclude that:

$$\min_{P \in \Gamma_{tr}(R, \lambda, L_{max})} \mathcal{D}(P||P_Y) \geq \min_{P \in \Gamma(R, L_{max})} \mathcal{D}(P||P_Y) - \delta(\tau), \quad (40)$$

where $\delta(\tau)$ can be made arbitrarily small by decreasing λ . By adding the term $c\mathcal{D}(R||P_X)$ to both sides of (40) we obtain:

$$\begin{aligned} c\mathcal{D}(R||P_X) + \min_{P \in \Gamma_{tr}(R, \lambda, L_{max})} \mathcal{D}(P||P_Y) &\geq \\ c\mathcal{D}(R||P_X) + \min_{P \in \Gamma(R, L_{max})} \mathcal{D}(P||P_Y) - \delta(\tau). \end{aligned} \quad (41)$$

⁵Note that when λ tends to 0, we do not need anymore to differentiate between the SI_{ks} and SI_{tr} games in the definition of $\Gamma(R, L_{max})$.

Given that (41) holds for any $R \in \mathcal{P}$, we can write:

$$\begin{aligned} \varepsilon_{tr}(\lambda) &= \min_R [c\mathcal{D}(R||P_X) + \min_{P \in \Gamma_{tr}(R, \lambda, L_{max})} \mathcal{D}(P||P_Y)] \\ &\geq \min_R [c\mathcal{D}(R||P_X) + \min_{P \in \Gamma(R, L_{max})} \mathcal{D}(P||P_Y)] - \delta(\tau), \end{aligned} \quad (42)$$

which concludes the proof due to the arbitrariness of $\delta(\tau)$. ■

A consequence of Theorem 6 is that $\lim_{\lambda \rightarrow 0} \varepsilon(\lambda) = 0$ if and only if $P_Y \in \Gamma(P_X, L_{max})$, which then can be seen as the smallest indistinguishability region for the SI_{tr} game. We conclude that the smallest indistinguishability regions for the two cases are the same thus implying that the security margin for the SI_{tr} setting, say \mathcal{SM}_{tr} , is the same of the SI_{ks} game, that is

$$\mathcal{SM}_{tr}(P_X, P_Y) = EMD(P_X, P_Y). \quad (43)$$

We remark that, for any allowed distortion $L_{max} < EMD(P_X, P_Y)$, the minimum value of the false positive error exponent (λ) which allows the defender to take a reliable decision in the SI_{tr} setting is lower than that in the SI_{ks} setting. However, the difference between the two settings regards the decay rate of the error probabilities, not the ultimate distinguishability of the sources.

We conclude this section with a brief discussion on the SI_{tr} game with different training sequences ($t_D^N \neq t_A^K$). Such a scenario provides a more realistic model in which the attacker is not able to compute exactly the acceptance region adopted by the Defender. It is known from [4] that, as long as the length of both sequences grows linearly with n , the indistinguishability region is equal to that of the game with equal training sequences. By relying on this result, it is not difficult to prove that the security margin remains the same even for such version of the game.

VI. SECURITY MARGIN COMPUTATION

In this section we address the problem of the actual computation of the security margin for two generic sources. By following the analysis given so far, we focus on the case of discrete sources, however at the end of the section we extend the analysis so to cover continuous sources as well.

Given two discrete sources $X \sim P_X$ and $Y \sim P_Y$, the computation of the security margin requires the evaluation of $EMD(P_X, P_Y)$. A closed form solution can be found only in some simple cases (see Section VI-A1 and VI-A2). More generally, the EMD between two sources can be computed by resorting to numerical analysis, and in fact, due to its wide use as a similarity measure in computer vision applications, several efficient algorithms have been proposed (see [21] for example). In the following, we describe

a fast iterative algorithm for the computation of the *EMD* between any two sources assuming that the distortion (or cost) function has the general form:

$$d(i, j) = |i - j|^p, \quad (44)$$

with $p \geq 1$. This is a case of great interest for $p = 1$ and $p = 2$, according to which the distortion between y^n and the attacked sequence z^n corresponds, respectively, to the L_1 and L_2^2 distance.

A. Hoffman's greedy algorithm for computing \mathcal{SM}

Let us assume that X and Y are discrete sources with alphabets \mathcal{X} and \mathcal{Y} . The transportation problem we have to solve for computing $\mathcal{SM}(P_Y, P_X)$, i.e. $EMD(P_Y, P_X)$, is known in modern literature as *Hitchcock transportation problem* [22]⁶, which, in turn, can be formulated as a linear programming problem in the following way:

$$EMD(P_X, P_Y) = \min_{S_{XY}} \sum_{i,j} d(i, j) S_{XY}(i, j), \quad (45)$$

where S_{XY} must satisfy the linear constraints:

$$\begin{aligned} \sum_j S_{XY}(i, j) &= P_X(i) && \forall i \in \mathcal{X} \\ \sum_i S_{XY}(i, j) &= P_Y(j) && \forall j \in \mathcal{Y} \\ S_{XY}(i, j) &\geq 0 && \forall i, j, \end{aligned} \quad (46)$$

and where, by referring to the original Monge formulation⁷, $S_{XY}(i, j)$ denotes the quantity of soil shipped from location (source) i to location (sink) j and $d(i, j)$ is the cost for shipping a unitary amount of soil from i to j .

A Transportation Problem (TP) like the one defined by equations (45) and (46) is a particular minimum cost flow problem [24] which, being linear, can be solved through the simplex method [25]. In general, the solution of TP depends on the cost function $d(\cdot, \cdot)$, however there are some classes of cost functions for which the solution can be found through a simple greedy algorithm. Specifically, the algorithm proposed by A.J. Hoffman in 1963 [26], allows to solve the transportation problem whenever $d(\cdot, \cdot)$ satisfies the so called Monge property [27], that is when:

$$d(i, j) + d(r, s) \leq d(i, s) + d(r, j), \quad (47)$$

⁶This is the discrete version of the Monge-Kantorovich mass transportation problem [15].

⁷Monge is considered the founding father of optimal transport [23].

$\forall(i, j, r, s)$ such that $1 \leq i < r \leq |\mathcal{X}|$ and $1 \leq j < s \leq |\mathcal{Y}|$.

It is easy to verify that Monge property is satisfied by any cost function of the form in (44), and, more in general, by any convex function of the quantity $|i - j|$. The iterative procedure proposed by Hoffman to solve the optimal transport problem is known as *north-west corner (NWC) rule* [26] and can be described as follows. Take the bin of \mathcal{X} with the smallest value and start moving its elements into the bin with the smallest value in \mathcal{Y} . When the smallest bin of \mathcal{Y} is filled, go on with the second smallest bin in \mathcal{Y} . Similarly, when the smallest bin in \mathcal{X} is emptied, go on with the second smallest bin in \mathcal{X} . The procedure is iterated until all the bins in \mathcal{X} have been moved into those of \mathcal{Y} . Let i^{low} (i^{up}) and j^{low} (j^{up}) denote the lower (upper) non-empty bins of \mathcal{X} and \mathcal{Y} respectively. A pseudocode description of the *NWC* rule is given below.

- 1) Initialize: $i := i^{low}$, $j := j^{low}$.
- 2) Set $S_{XY}(i, j) := \min\{P_X(i), P_Y(j)\}$.
- 3) Adjust the ‘supply’ distribution $P_X(i) := P_X(i) - S_{XY}(i, j)$ and the ‘demand’ distribution $P_Y(j) := P_Y(j) - S_{XY}(i, j)$.
If $P_X(i) = 0$ then $i := i + 1$ and if $P_Y(j) = 0$ then $j := j + 1$.
- 4) If $j < j^{up}$ or $P_Y(j^{up}) > 0$ go back to Step 2).

The above procedure is described graphically in Figure 3. In the figure, we chose two distributions with disjoint supports for sake of clarity, however the procedure is valid regardless of how the two distributions are spread along the real line. Interestingly, the *NWC* rule does not depend explicitly on the cost matrix, so the transportation map obtained through it is the same regardless of the Monge cost. According to Hoffman’s greedy algorithm, when the cost function satisfies Monge’s property, the *EMD* can be computed in linear running time: the number of elementary operations, in fact, is at most equal to $|\mathcal{X}| + |\mathcal{Y}|$ ⁸. This represents a dramatic simplification with respect to the complexity required to solve a general Hitchcock transportation problem (see for example [28]).

As detailed below, in some cases, it is possible to derive a closed form expression for the security margin.

⁸For sake of simplicity, the iterative algorithm described by the pseudocode spans all the bins between the minimum and the maximum non-empty bins. However, only the values $i \in \mathcal{X}$ and $j \in \mathcal{Y}$ must be considered given that for all the empty bins i and j we have $S_{XY}(i, j) = 0$.

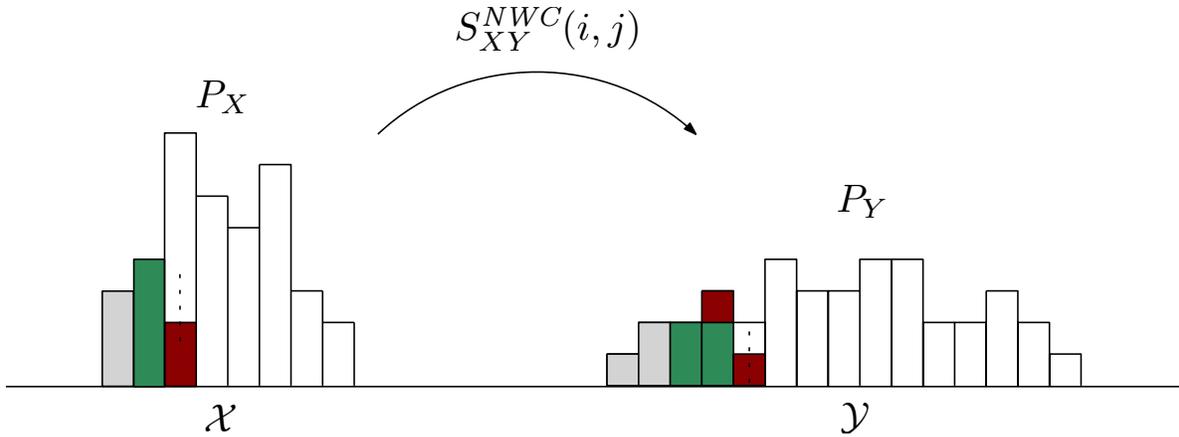


Fig. 3. Graphical representation of the north-west corner rule for the earth mover transportation problem (Monge problem). P_X and P_Y are two generic earth piles (source and sink) \mathcal{X} and \mathcal{Y} , while $S_{XY}^{NWC}(i, j)$ denotes the amount of earth moved from location i to j .

1) *Uniform sources with different cardinalities:* Let X and Y be two uniform pmf's with alphabets \mathcal{X} and \mathcal{Y} such that $|\mathcal{X}| = \alpha|\mathcal{Y}|$, with $\alpha \in \mathbb{N}$. In this case, thanks to Hoffman's algorithm we can express $\mathcal{SM}(P_X, P_Y)$ in closed form:

$$\mathcal{SM}_{L_p^p}(P_X, P_Y) = \frac{1}{|\mathcal{Y}|} \sum_{i=0}^{|\mathcal{X}|-1} \sum_{j=0}^{\alpha-1} (|i^{low} - j^{low}| - j - (\alpha - 1)i)^p, \quad (48)$$

The formula implicitly assumes that $j^{low} > i^{low}$, the extension to the case in which such a relationship does not hold being immediate.

2) *Security Margin under the L_1 distance:* If the distortion function corresponds to the L_1 distance, the *EMD* (and hence the security margin) assumes a particularly simple form. Specifically, by applying the flow decomposition principle [29], it is easy to see that the security margin between P and Q can be calculated as follows:

$$\mathcal{SM}_{L_1}(P, Q) = \sum_{i=\min\{i^{low}, j^{low}\}}^{\max\{i^{up}, j^{up}\}} \left| \sum_{s=1}^i (P(s) - Q(s)) \right|. \quad (49)$$

B. Continuous sources

The analysis carried out in the previous sections is limited to discrete sources. When continuous sources are considered, we can quantize the probability density functions (pdf's) of the sources and apply the analysis for discrete sources. By letting the quantization step tend to zero, the *EMD* between P_X and P_Y can still be regarded as the security margin between the two sources. In this case, a general expression

for the \mathcal{SM} can be derived by considering the *continuous transportation problem* (CTP), known as Monge-Kantorovic formulation of the mass transportation problem:

$$\mathcal{SM}(P_X, P_Y) = \min_{S_{XY}(x,y)} \iint c(x,y) S_{XY}(x,y) dx dy, \quad (50)$$

subject to the constraints

$$\begin{aligned} \int S_{XY}(x,y) dx &= P_Y(y) \\ \int S_{XY}(x,y) dy &= P_X(x) \\ S_{XY}(x,y) &\geq 0 \quad \forall x,y, \end{aligned} \quad (51)$$

where c is a continuous cost function $c(x,y) : X \times Y \rightarrow \mathbb{R}$. If $c(x,y)$ satisfies the continuous Monge property [27], that is if:

$$c(x,y) + c(x',y') \leq c(x',y) + c(x,y'), \quad (52)$$

for all $x \leq x', y \leq y'$, the optimum transportation map corresponds to the Hoeffding distribution [17] defined as follows. Let $C_X(x)$ and $C_Y(y)$ be the cumulative distributions of X and Y respectively, and let $C_{XY}(x,y)$ be the cumulative transportation map, that is:

$$C_{XY}(x,y) = \int_{-\infty}^x \int_{-\infty}^y S_{XY}(u,v) du dv. \quad (53)$$

The optimum transportation map is obtained by letting:

$$C_{XY}^*(x,y) = \min\{C_X(x), C_Y(y)\}, \quad \forall (x,y) \in \mathbb{R}^2, \quad (54)$$

which generalizes the *NWC* rule. Given the optimum transportation map, one can compute $\mathcal{SM}(P_Y, P_X)$ by evaluating the integral in (50). In general, however, finding a closed form expression is not an easy task.

A particularly simple and insightful formula can be obtained when the cost function corresponds to the squared Euclidean distance. Let us assume, then, that $c(x,y) = (x-y)^2$ (in this case $\mathcal{SM}(P_X, P_Y)$ corresponds to the squared Mallows distance - see equation (18) - and let X and Y be two continuous sources with means μ_X and μ_Y , variances σ_X and σ_Y and covariance $covXY$. As shown in [30] (decomposition theorem), the expectation in (18) can be rewritten as follows:

$$\begin{aligned} E_{XY}[(X-Y)^2] &= (\mu_X - \mu_Y)^2 + (\sigma_X - \sigma_Y)^2 \\ &\quad + 2[\sigma_X \sigma_Y - covXY], \end{aligned} \quad (55)$$

where the three terms express, respectively, the difference in location, spread and shape between the variables X and Y [31]. Interestingly, the covariance $covXY$ is the only term in (55) which depends on the joint pdf of X and Y . Then, in order to find the security margin, we only have to compute the maximum covariance over all the possible joint pdf's:

$$\begin{aligned} \mathcal{SM}_{L_2^2}(P_X, P_Y) &= (\mu_X - \mu_Y)^2 + (\sigma_X - \sigma_Y)^2 \\ &+ 2[\sigma_X \sigma_Y - \max_{P_{XY}: \substack{\sum_y P_{XY} = P_X \\ \sum_x P_{XY} = P_Y}} covXY]. \end{aligned} \quad (56)$$

By assuming that X and Y are independent, i.e. $P_{XY} = P_X P_Y$, we have $covXY = 0$, hence permitting us to derive a general upper bound for the security margin:⁹

$$\mathcal{SM}_{L_2^2}(P_X, P_Y) \leq (\mu_X - \mu_Y)^2 + \sigma_X^2 + \sigma_Y^2. \quad (57)$$

When P_X and P_Y have the same form, for instance when the random variables X and Y are both distributed according to a Gaussian or a Laplacian distribution, the security margin assumes a particularly simple expression. In this case, in fact, it is possible to turn P_X into P_Y by imposing a deterministic relationship between X and Y , namely $Y = \frac{\sigma_Y}{\sigma_X} X + (\mu_Y - \frac{\sigma_Y}{\sigma_X} \mu_X)$. In this way the covariance term is maximum and equal to $\sigma_X \sigma_Y$, and hence the contribution of the shape term in the security margin vanishes, yielding:

$$\mathcal{SM}_{L_2^2}(P_X, P_Y) = (\mu_X - \mu_Y)^2 + (\sigma_X - \sigma_Y)^2. \quad (58)$$

This is a remarkable, and somewhat surprising, result stating that the distinguishability of two sources belonging to the same class depends only on their means and variances, regardless of their particular pdf.

VII. THE SECURITY MARGIN WITH L_∞ DISTANCE

We conclude the paper by extending the definition of the Security Margin to the case in which the distortion measure constraining the attacker is expressed in terms of the maximum absolute distance between the samples of y^n and z^n , that is to the case in which the distortion is measured by relying on the L_∞ distance.

The interest in this case is motivated by the importance that the L_∞ distance has in applications where the perceptual distortion between the sequence y^n and the attacked sequence z^n must be taken into account. This is the case, for instance, in image forensics applications [32], [33], [34], [3], [7], wherein the attacker is interested in hiding the true source of an image. In this case, the use of a distortion measure

⁹We point out that relation (55), as well as the upper bound in (57), holds for the discrete case too.

based on the L_∞ distance ensures that the attacked image is perceptually similar to the original one. In our analysis, we will refer to the case of known sources, the extension to the SI_{tr} case being immediate.

A. The SI_{ks} game with L_∞ distance

We start by observing that the adoption of the L_∞ distance requires that the SI_{ks} game is, partly, redefined due to the non-additive nature of the distortion constraint. In this case, in fact, it does not make any sense to define the distortion constraint in terms of average per-letter distortion and let the overall allowed distortion to increase with n .

Similarly to the previous cases, it is possible to express the distortion constraint by limiting the set of transportation maps the attacker can choose from. More specifically, we observe that the maximum distance between the two sequences y^n and z^n can be rewritten as follows:

$$d_{L_\infty}(y^n, z^n) = \max_j |z_j - y_j| = \max_{(i,j): S_{YZ}^n(i,j) \neq 0} |i - j|. \quad (59)$$

By using the above formula in the definition of the set of admissible maps (i.e. in the second line of (7)), we can still define the set of strategies of the attacker as the set of rules associating an admissible map to the to-be-attacked sequence, as in (29). In the following, we will refer to the set of admissible maps resulting from the use of the d_{L_∞} distance as $\mathcal{A}_{L_\infty}^n(L_{max}, P_{y^n})$.

Passing to the analysis of the indistinguishability region, it is easy to see that relation (12) continues to hold by replacing $\mathcal{A}^n(L_{max}, P_{y^n})$ with $\mathcal{A}_{L_\infty}^n(L_{max}, P_{y^n})$. In fact, the dominant strategy for the defender does not depend on the set of strategies available to the attacker. The asymptotic version of $\Gamma_{L_\infty}^n(P_X, \lambda, L_{max})$ can also be defined as in (13), namely:

$$\begin{aligned} \Gamma_{L_\infty}(P_X, \lambda, L_{max}) = & \quad (60) \\ \{P \in \mathcal{P} : \exists S_{YZ} \in \mathcal{A}_{L_\infty}(L_{max}, P) \text{ s.t. } S_Z \in \Lambda^*(P_X, \lambda)\}, & \end{aligned}$$

where $\mathcal{A}_{L_\infty}(L_{max}, P)$ is the asymptotic counterpart of $\mathcal{A}_{L_\infty}^n(L_{max}, P)$. The next step requires the extension of Theorem 2 to the SI_{ks} game with L_∞ distance, that is we need to prove that the set in (60) contains all the sources that can not be distinguished from X because of the attack, even when the length of the observed sequence tends to infinity. This is a critical step since such theorem was proved in [5] by assuming an additive distortion measure, which clearly is not the case when the L_∞ distance is adopted. Roughly speaking, we need to prove that when $n \rightarrow \infty$ the elements of $\Gamma_{L_\infty}^n(P_X, \lambda, L_{max})$ are dense in $\Gamma_{L_\infty}(P_X, \lambda, L_{max})$ (in which case Theorem 2 can be proven in a way similar to Sanov's Theorem [5]). More formally, we need to prove that for any $P_Y \in \Gamma_{L_\infty}(P_X, \lambda, L_{max})$ and any $\delta > 0$, a

pmf $Q^n \in \Gamma_{L_\infty}^n(P_X, \lambda, L_{max})$ exists such that the distance between P_Y and Q^n is smaller than δ . The proof requires only some minor modifications with respect to the proof given in [2] (Lemma 2 in the Appendix) and is skipped for sake of brevity.

B. Security Margin for the SI_{ks} game with L_∞ distance

As a next step, we must study the behavior of the indistinguishability region of the test when $\lambda \rightarrow 0$ (to determine the smallest indistinguishability region). As we will see, even if the adoption of the d_{L_∞} distance prevents a direct formulation of the problem in terms of *EMD*, the distinguishability between two sources X and Y is still closely related to the optimal transportation map between P_X and P_Y . The basis for such a connection is rooted in the following property.

Property 3. *Given two distributions P and Q , the transportation map S_{PQ}^{NWC} obtained by applying the NWC rule to P and Q is a solution of the problem*

$$\min_{S_{YZ}: S_Y=P, S_Z=Q} \left(\max_{(i,j) \in S_{YZ}(i,j) \neq 0} |i - j| \right). \quad (61)$$

Proof: Let $S^* \neq S_{PQ}^{NWC}$ be a generic transformation mapping P into Q . Given that $S^* \neq S_{PQ}^{NWC}$ there exists at least one quadruple of bins (t, r, v, s) , with $t < r$ and $v < s$, for which, $S^*(t, s) > 0$ and $S^*(r, v) > 0$. Let us assume, without loss of generality, that $S^*(t, s) \leq S^*(r, v)$. We now define a new map S' which is obtained from S^* by letting:

$$S'(t, v) = S^*(t, v) + S^*(t, s) \quad (62)$$

$$S'(t, s) = 0$$

$$S'(r, v) = S^*(r, v) - S^*(t, s)$$

$$S'(r, s) = S^*(r, s) + S^*(t, s).$$

Since $\max\{|t - s|, |r - v|\} > \max\{|t - v|, |r - s|\}$, the maximum distortion introduced by S' is lower than or equal to that introduced by S^* , that is:

$$\max_{(i,j) \in S^*(i,j) \neq 0} |i - j| \geq \max_{(i,j) \in S'(i,j) \neq 0} |i - j|. \quad (63)$$

We now inspect S' , if there is another quadruple of bins (t', r', v', s') satisfying the same properties of (t, r, v, s) , we let $S^* = S'$ and iterate the above procedure. The process ends when no quadruple of bins with the required properties exists and hence when $S' = S_{PQ}^{NWC}$. Since at each step the distortion introduced by the new map does not increase, the above procedure proves that S_{PQ}^{NWC} introduces a distortion lower

than or equal to that introduced by any other S^* mapping P into Q , thus proving that S_{PQ}^{NWC} achieves the minimum in (61). ■

Thanks to Property 3, the set $\Gamma_{L_\infty}(P_X, \lambda, L_{max})$ in (60) can be rewritten as follows:

$$\Gamma_{L_\infty}(P_X, \lambda, L_{max}) = \{P \in \mathcal{P} : \exists Q \in \Lambda^*(P_X, \lambda) \text{ s.t.} \quad (64)$$

$$\max_{(i,j): S_{PQ}^{NWC}(i,j) \neq 0} |i - j| \leq L_{max}\}.$$

By letting λ tend to 0, we obtain the smallest indistinguishability region, thus extending Theorem 3 to the SI_{ks} game with d_{L_∞} distance.

Theorem 7. *Given two sources $X \sim P_X$ and $Y \sim P_Y$ and a maximum allowable per-letter distortion L_{max} , and given:*

$$\Gamma(P_X, L_{max}) = \{P \in \mathcal{P} : \max_{(i,j) \in S_{PP_X}^{NWC}} |i - j| \leq L_{max}\}, \quad (65)$$

the maximum achievable false negative error exponent ε for the SI_{ks} game with L_∞ distance is

$$\lim_{\lambda \rightarrow 0} \lim_{n \rightarrow \infty} -\frac{1}{n} \log P_{fn} = \min_{P \in \Gamma_{L_\infty}(P_X, L_{max})} \mathcal{D}(P || P_Y). \quad (66)$$

Proof: The proof relies on the extension of Property 1 and Lemma 1 to the L_∞ case. The extension of Property 1 is immediate since, once again, the indistinguishability region depends on λ only through $\Lambda^*(P_X, \lambda)$, whose form does not depend on the particular norm adopted to express the distortion constraint. The extension of Lemma 1 requires some more care and is proven in Appendix B. For the rest, the theorem can be proven by reasoning as in the proof of Theorem 3. ■

As a consequence of Theorem 7, the distinguishability of two sources depends again on the optimum transportation map between the pmf's of the two sources. Specifically, given the sources X and Y , the defender is able to distinguish between them in this adversarial setting, only if

$$\max_{(i,j) \in S_{P_Y P_X}^{NWC}} |i - j| > L_{max}. \quad (67)$$

Condition (67) can be used to determine the maximum attacking distortion for which D is able to distinguish the sources X and Y , i.e. $\mathcal{SM}(P_X, P_Y)$.

Definition 2 (Security Margin for the L_∞ case). *Let $X \sim P_X$ and $Y \sim P_Y$ be two discrete memoryless sources. The maximum distortion for which the two sources can be reliably distinguished in the SI_{ks} setting with L_∞ distance is given by*

$$\mathcal{SM}_{L_\infty}(P_Y, P_X) = \max_{(i,j): S_{P_Y P_X}^{NWC}(i,j) \neq 0} |i - j|, \quad (68)$$

where $S_{P_Y P_X}^{\text{NWC}}$ is obtained by applying the NWC rule to map P_Y into P_X .

Even if we proved Theorem 7 for the case of known sources, it is possible to extend it to the SI_{tr} game. The proof goes along the same lines used for the SI_{ks} case and is omitted for sake of brevity.

VIII. CONCLUSIONS

By interpreting the attacker's optimum strategy in the SI_{ks} (and SI_{tr}) game as the solution of an optimum transport problem, we introduced the concept of security margin, a single measure summarizing the distinguishability of two sources under adversarial conditions. We also described an efficient algorithm to compute the security margin between several classes of sources. By relying on the security margin concept, we can understand who between the attacker and the defender is going to win the source identification game under asymptotic conditions. Among the practical applications of our analysis we mention image forensics, wherein the defender is interested in distinguishing images produced by different devices, and intrusion detection, in which the defender is willing to distinguish normal and anomalous behaviors. In the first case, knowing the SM between the statistics ruling the emission of images from different sources permits to compute the amount of distortion required to make the images produced by the two sources indistinguishable. In the latter case, the SM determines how much an intruder must deviate from the intended, anomalous, behavior to make its presence undetectable by the analyst.

ACKNOWLEDGMENT

We thank Alessandro Agnetis for the useful discussions on the optimization problems underlying the computation of the EMD .

This work has been partially supported by the European Office of Aerospace Research and Development under Grant FA8655-12-1- 2138: AMULET - A multi-clue approach to image forensics, and the the REWIND Project, funded by the Future and Emerging Technologies (FET) programme within the 7FP of the EC, under grant 268478.

REFERENCES

- [1] M. Barni and F. Pérez-González, "Coping with the enemy: advances in adversary-aware signal processing," in *ICASSP 2013, IEEE Int. Conf. Acoustics, Speech and Signal Processing*, Vancouver, Canada, 26-31 May 2013, pp. 8682–8686.
- [2] M. Barni and B. Tondi, "The source identification game: an information-theoretic perspective," *IEEE Transactions on Information Forensics and Security*, vol. 8, no. 3, pp. 450–463, March 2013.
- [3] —, "Optimum forensic and counter-forensic strategies for source identification with training data," in *Proc. of WIFS'12, IEEE International Workshop on Information Forensics and Security*, Tenerife, Spain, 2-5 December 2012, pp. 199–204.

- [4] —, “Binary hypothesis testing game with training data,” *Information Theory, IEEE Transactions on*, vol. PP, no. 99, pp. 1–1, 2014.
- [5] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley Interscience, 1991.
- [6] C. Villani, *Optimal Transport: Old and New*. Berlin: Springer-Verlag, 2009.
- [7] P. Comesana-Alfaro and F. Pérez-González, “Optimal counterforensics for histogram-based forensics,” in *Proc. IEEE Int. Conf. Acoust., Speech, and Signal Process*, 2013.
- [8] F. Balado, “The role of permutation coding in minimum-distortion perfect counterforensics,” in *ICASSP 2013, IEEE Int. Conf. Acoustics, Speech and Signal Processing*, Florence, Italy, 4-9 May 2013.
- [9] F. Balado and D. Haughton, “Permutation codes and steganography,” in *Proc. of ICASSP 2013, IEEE International Conference on Acoustics, Speech and Signal Processing*. Vancouver, Canada: IEEE, 26-31 May 2013, pp. 2954–2958.
- [10] M. Barni and B. Tondi, “The security margin: a measure of source distinguishability under adversarial conditions,” in *Proc. of GlobalSip’13, IEEE Global Conference on Signal and Information Processing*, Austin, Texas, 3-5 December 2013.
- [11] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems. 2nd edition*. Cambridge University Press, 2011.
- [12] M. J. Osborne and A. Rubinstein, *A Course in Game Theory*. MIT Press, 1994.
- [13] Y. C. Chen, N. Van Long, and X. Luo, “Iterated strict dominance in general games,” *Games and Economic Behavior*, vol. 61, no. 2, pp. 299–315, November 2007.
- [14] Y. Rubner, C. Tomasi, and L. J. Guibas, “The earth mover’s distance as a metric for image retrieval,” *Int. J. Comput. Vision*, vol. 40, no. 2, pp. 99–121, November 2000.
- [15] S. T. Rachev, *Mass Transportation Problems: Volume I: Theory*. Springer, 1998, vol. 1.
- [16] C. Villani, *Topics in Optimal Transportation*. Graduate Studies in Mathematics Series: American Mathematical Society, 2003, vol. 58.
- [17] S. T. Rachev, “The monge-kantorovich mass transference problem and its stochastic applications,” *Theory of Probability & Its Applications*, vol. 29, no. 4, pp. 647–676, 1985.
- [18] E. Levina and P. Bickel, “The Earth Mover’s distance is the Mallows distance: some insights from statistics,” in *Proc. of ICCV 2001, Eighth IEEE International Conference on Computer Vision*, vol. 2, 2001, pp. 251–256 vol.2.
- [19] M. Gutman, “Asymptotically optimal classification for multiple tests with empirically observed statistics,” *IEEE Transactions on Information Theory*, vol. 35, no. 2, pp. 401–408, March 1989.
- [20] M. Kendall and S. Stuart, *The Advanced Theory of Statistics, vol. 2, 4th edition*. New York: MacMillan, 1979.
- [21] O. Pele and M. Werman, “Fast and robust Earth Mover’s distances,” in *Proc. ICCV’09, 12th IEEE International Conference on Computer Vision*, 2009, pp. 460–467.
- [22] F. L. Hitchcock, “The distribution of a product from several sources to numerous localities,” *Journal of Mathematical Physics*, vol. 20, pp. 224–230.
- [23] G. Monge, *Mémoire sur la théorie des déblais et des remblais*. De l’Imprimerie Royale, 1781.
- [24] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin, *Network Flows: Theory, Algorithms, and Applications*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1993.
- [25] V. Chvatal, “Linear programming,” *A Series of Books in the Mathematical Sciences, New York: Freeman, 1983*, vol. 1, 1983.
- [26] A. Hoffman, “On simple linear programming problems,” in *Proceedings of Symposia in Pure Mathematics*, vol. 7. World Scientific, 1963, pp. 317–327.

- [27] R. E. Burkard, B. Klinz, and R. Rudolf, “Perspectives of monge properties in optimization,” *Discrete Applied Mathematics*, vol. 70, no. 2, pp. 95–161, 1996.
- [28] J. B. Orlin, “A faster strongly polynomial minimum cost flow algorithm,” *Operations research*, vol. 41, no. 2, pp. 338–350, 1993.
- [29] A. C. Williams, “A treatment of transportation problems by decomposition,” *Journal of the Society for Industrial and Applied Mathematics*, vol. 10, no. 1, pp. pp. 35–48.
- [30] A. Irpino and E. Romano, “Optimal histogram representation of large data sets: Fisher vs piecewise linear approximation.” in *EGC*, ser. Revue des Nouvelles Technologies de l’Information, M. Noirhomme-Fraiture and G. Venturini, Eds., vol. RNTI-E-9. Cepadues-Editions, 2007, pp. 99–110.
- [31] K. Košmelj and L. Billard, “Mallows’ L_2 distance in some multivariate methods and its application to histogram-type data,” *Metodoloski Zvezki*, vol. 9, no. 2, pp. 107–118, 2012.
- [32] R. Böhme and M. Kirchner, “Counter-forensics: Attacking image forensics,” in *Digital Image Forensics*, H. T. Sencar and N. Memon, Eds. Springer Berlin / Heidelberg, 2012.
- [33] M. Barni, M. Fontani, and B. Tondi, “A universal technique to hide traces of histogram-based image manipulations,” in *Proc. of the ACM Multimedia and Security Workshop*, Coventry, UK, 6-7 September 2012, pp. 97–104.
- [34] M. Barni, M. Fontani, and B. Tondi, “A universal attack against histogram-based image forensics,” *International Journal of Digital Crime and Forensics (IJDCF)*, vol. 5, no. 3, 2013.
- [35] D. Bertsimas and J. Tsitsiklis, *Introduction to Linear Optimization*, 1st ed. Athena Scientific, 1997.

APPENDIX

A. Behavior of $\Gamma(P_X, \lambda, L_{max})$ and $\Gamma_{tr}(R, \lambda, L_{max})$ for $\lambda \rightarrow 0$.

We start by studying the behavior of $\Gamma(P_X, \lambda, L_{max})$ when $\lambda \rightarrow 0$. More specifically, we show that for small values of λ the set $\Gamma(P_X, \lambda, L_{max})$ approaches $\Gamma(P_X, L_{max})$ smoothly.

As a first step, we highlight the following property.

Property 4. $EMD(P, Q)$ is a continuous and convex function of P and Q .

Proof: Property 4 follows immediately if we look at the EMD as the solution of a Linear Programming (LP) problem (see Section VI-A), wherein P and Q are the known terms of the linear constraints. In fact, it is a known result in operations research that the minimum of the objective function of an LP problem is a continuous and convex function of the known terms of the linear constraints [35]. ■

By exploiting the continuity of the divergence and the continuity and convexity of the EMD , we now show that when λ tends to 0, the set $\Gamma(P_X, \lambda, L_{max})$ tends to $\Gamma(P_X, L_{max})$ regularly. More precisely, the following lemma holds.

Lemma 1. Let $X \sim P_X$ be an information source and L_{max} the maximum allowable average per-letter

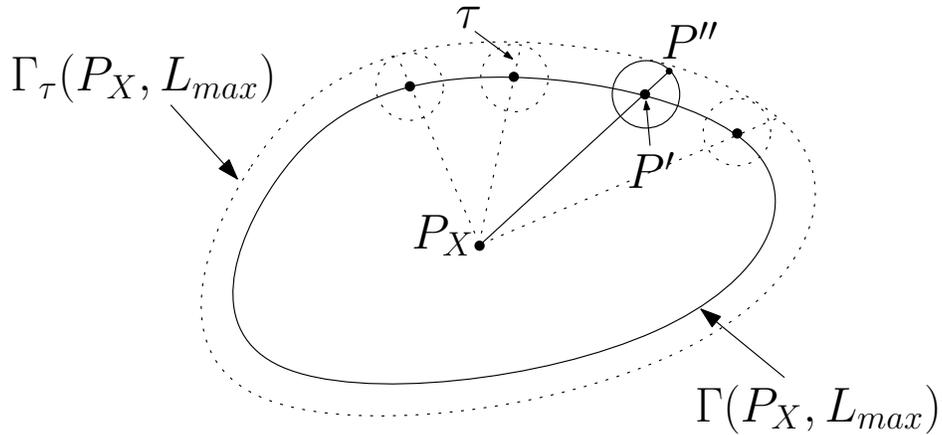


Fig. 4. Graphical representation of the set $\Gamma_\tau(P_X, L_{max})$.

distortion in the SI_{ks} game. The set $\Gamma(P_X, \lambda, L_{max})$, defined in (19), satisfies the following property:

$$\begin{aligned} \forall \tau > 0, \exists \lambda > 0 \text{ s.t. } \forall P \in \Gamma(P_X, \lambda, L_{max}) & \quad (\text{A1}) \\ \exists P' \in \Gamma(P_X, L_{max}) \text{ s.t. } P \in B(P', \tau), & \end{aligned}$$

where $\Gamma(P_X, L_{max})$ is defined as in (20) and $B(P', \tau)$ is a ball centered in P' with radius τ .

Proof: Throughout the proof we will refer to Figure 4 where all the sets and quantities involved in the proof are sketched. For any $\tau > 0$, we consider the set:

$$\begin{aligned} \Gamma_\tau(P_X, L_{max}) = & \quad (\text{A2}) \\ \{P : \exists P' \in \Gamma(P_X, L_{max}) \text{ s.t. } P \in B(P', \tau)\}. & \end{aligned}$$

With such a definition, we can rephrase (A2) as follows:

$$\forall \tau > 0, \exists \lambda > 0 \text{ s.t. } \Gamma(P_X, \lambda, L_{max}) \subseteq \Gamma_\tau(P_X, L_{max}). \quad (\text{A3})$$

For sake of simplicity, we will prove a slightly stronger version of the lemma by means of the following two-step proof. First, we will show that a subset of $\Gamma_\tau(P_X, L_{max})$ exists having the following form:

$$\Gamma_\tau^{sub}(P_X, L_{max}) = \{P : EMD(P, P_X) \leq L_{max} + \delta(\tau)\}, \quad (\text{A4})$$

for some $\delta(\tau) > 0$. Then, we will prove that for small enough λ , any $P \in \Gamma(P_X, \lambda, L_{max})$ belongs to $\Gamma_\tau^{sub}(P_X, L_{max})$.

To start with, let P' be any point on $\mathcal{B}(\Gamma(P_X, L_{max}))$, the boundary of $\Gamma(P_X, L_{max})$. Among all the points on the boundary of the ball of radius τ and centered in P' , consider the one, name it P'' ,

lying along the direction given by the line joining P_X and P' and falling outside $\Gamma(P_X, D_{max})$ (see Figure 4). By the convexity of the EMD (Property 4) and since $EMD = 0$ if and only if $P = P_X$, we conclude that $EMD(P'', P_X) > EMD(P', P_X)$. Since P' lies on the boundary of $\Gamma(P_X, L_{max})$ we know that $EMD(P'', P_X) = L_{max} + \mu$, where $\mu = \mu(P', \tau)$ is a strictly positive quantity. We now show that the first part the proof holds by letting $\delta(\tau) = \min_{P' \in B(\Gamma(P_X, L_{max}))} \mu(P', \tau)$. To this purpose, let P be any point in set $\Gamma_\tau^{sub}(P_X, L_{max})$ for the above choice of $\delta(\tau)$. If $P \in \Gamma(P_X, L_{max})$, then, by definition, P also belongs to $\Gamma_\tau(P_X, L_{max})$. On the other side, if P lies outside $\Gamma(P_X, L_{max})$, let us denote by P^* the point lying on the boundary of the set $\Gamma(P_X, L_{max})$ along the line joining P and P_X , and let P^{**} be the point where the same line crosses the ball $B(P^*, \tau)$ outside $\Gamma(P_X, L_{max})$. Now, $EMD(P, P_X) \leq L_{max} + \delta(\tau) \leq EMD(P^{**}, P_X)$ by construction. Because of the convexity of EMD , then $P \in B(P^*, \tau)$ as required.

Let us now pass to the second part of the proof. First, we notice that set $\Gamma(P_X, \lambda, L_{max})$ depends on λ only through the acceptance region $\Lambda^*(P_X, \lambda)$. If λ is small, due to the continuity of the divergence, for any $Q \in \Lambda^*(P_X, \lambda)$ we will have $Q \in B(P_X, \kappa(\lambda))$ for some $\kappa(\lambda)$ such that $\kappa(\lambda) \rightarrow 0$ when $\lambda \rightarrow 0$. Let, then, P be a pmf in $\Gamma(P_X, \lambda, L_{max})$. By definition, a $Q \in \Lambda^*(P_X, \lambda)$ exists s.t. $EMD(P, Q) \leq L_{max}$. If λ is small, due to the proximity of Q to P_X and the continuity of the EMD we have that $EMD(P, P_X) < EMD(P, Q) + \eta(\lambda) \leq L_{max} + \eta(\lambda)$ with $\eta(\lambda)$ approaching 0 when $\lambda \rightarrow 0$. In particular, if λ is small enough $\eta(\lambda) < \delta(\tau)$ and hence $P \in \Gamma_\tau^{sub}(P_X, L_{max})$ which in turn is entirely contained in $\Gamma_\tau(P_X, L_{max})$ thus completing the proof. ■

In the same way, we can prove that Lemma 1 holds also when $\Gamma(P_X, \lambda, L_{max})$ is replaced by $\Gamma_{tr}(R, \lambda, L_{max})$ and $\Gamma(P_X, L_{max})$ by $\Gamma(R, L_{max})$ with a generic R instead of P_X . To be convinced about that, it is sufficient to note that the only difference between Γ and Γ_{tr} relies on the test function which defines the acceptance region, respectively the divergence and the h_c function. Since the h_c function is still a continuous and convex function and, likewise \mathcal{D} , is equal to zero if and only if its arguments are identical, the proof that we used for Lemma 1 still holds.

B. Behavior of $\Gamma_{L_\infty}(P_X, \lambda, L_{max})$ for $\lambda \rightarrow 0$.

We prove that when $\lambda \rightarrow 0$, $\Gamma_{L_\infty}(P_X, \lambda, L_{max})$ approaches $\Gamma_{L_\infty}(P_X, L_{max})$ regularly, in the sense stated by the following lemma.

Lemma 2 (Extension of Lemma 1 to the L_∞ case). *Let $X \sim P_X$ be an information source and L_{max} the maximum per-sample distortion allowed to the attacker. The set $\Gamma_{L_\infty}(P_X, \lambda, L_{max})$, defined in Section*

VII, satisfies the following property:

$$\begin{aligned} \forall \tau > 0, \exists \lambda > 0 \text{ s.t.}, \forall P \in \Gamma_{L_\infty}(P_X, \lambda, L_{max}) \\ \exists P' \in \Gamma_{L_\infty}(P_X, L_{max}) \text{ s.t. } P \in B(P', \tau), \end{aligned} \quad (\text{A5})$$

where $B(P', \tau)$ is a ball centered in P' with radius τ .

Proof: We will prove the lemma by assuming that the distance defining the ball $B(P', \tau)$ is the L_1 distance, extending the proof to other distances being straightforward.

For a fixed $\tau > 0$, let P be a pmf in $\Gamma_{L_\infty}(P_X, \lambda, L_{max})$ for some λ . This means that at least one pmf $Q \in \Lambda^*(P_X, \lambda)$ exists, such that P can be mapped into Q with maximum shipment distance lower than or equal to L_{max} . From equation (14) and by exploiting the continuity of the divergence function, we argue that $Q \in \mathcal{B}(P_X, \gamma(\lambda))$ for some positive $\gamma(\lambda)$, and where $\gamma(\lambda) \rightarrow 0$ as $\lambda \rightarrow 0$. Accordingly, P_X can be written as $P_X(j) = Q(j) + \gamma(j)$, $\forall j$, where $\sum_{j \in \mathcal{X}} |\gamma(j)| < \gamma(\lambda)$. Note that, by construction, $\sum_j \gamma(j) = 0$ and $\gamma(j) \rightarrow 0$ when $\lambda \rightarrow 0$. Let S_{PQ} be an admissible map bringing P into Q (such a map surely exists by construction). We prove the lemma by explicitly building a pmf P' and a new admissible transportation map S' , such that, P' is arbitrarily close to P (for a small enough λ) and S' maps P' into P_X . We start by introducing two new quantities, namely $\gamma^+(j)$, defined as follows:

$$\begin{aligned} \gamma^+(j) &= \gamma(j) && \text{if } P_X(j) - Q(j) \geq 0 \\ \gamma^+(j) &= 0 && \text{if } P_X(j) - Q(j) < 0, \end{aligned} \quad (\text{A6})$$

and $\gamma^-(j)$ defined as

$$\begin{aligned} \gamma^-(j) &= -\gamma(j) && \text{if } P_X(j) - Q(j) < 0 \\ \gamma^-(j) &= 0 && \text{if } P_X(j) - Q(j) \geq 0. \end{aligned} \quad (\text{A7})$$

A graphical interpretation of γ^+ and γ^- is given in Figure 5. Clearly, $\sum_j \gamma^-(j) = \sum_j \gamma^+(j)$. With the above definitions, we can look at the demand distribution Q as consisting of two amounts: the mass distribution D , with $D(j) = \min\{P_X(j), Q(j)\}$, and γ^- . According to the superposition principle, the map S_{PQ} can then be split into two sub-maps: one which satisfies the demand of D (let us call it S_{PQ}^D), and one that satisfies the demand of γ^- (let us call it $S_{PQ}^{\gamma^-}$). The same distinction can be made in the

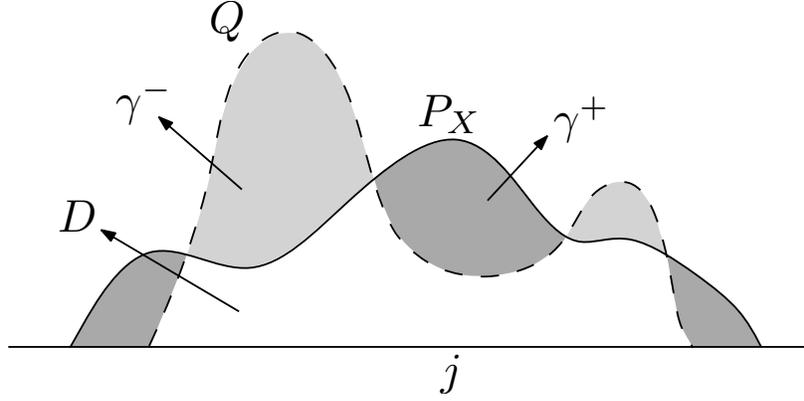


Fig. 5. Geometric interpretation of γ^+ , γ^- and $D(j)$.

source distribution, as follows:

$$\begin{aligned} P(i) &= \sum_j S_{PQ}(i, j) \\ &= \sum_j S_{PQ}^D(i, j) + \sum_j S_{PQ}^\gamma(i, j) = P_D(i) + P_\gamma(i), \end{aligned} \quad (\text{A8})$$

where P_D and P_γ are the masses in the source distribution which are used to satisfy the mass demand pertaining to D and γ^- according to mapping S_{PQ} . Then, $\sum_i P_D(i) = D$ and $\sum_i P_\gamma(i) = \gamma^-$. In order to construct the pmf P' we are looking for, we simply remove from P the amount of mass P_γ used to fill γ^- and redistribute it according to γ^+ . Specifically, we have

$$P'(i) = P_D(i) + \gamma^+(i) \quad (\text{A9})$$

$$S'(i, j) = S_{PQ}^D(i, j) + \gamma^+(j)\delta(i, j), \quad (\text{A10})$$

where $\delta(i, j)$ is equal to 1 if $i = j$ and 0 otherwise. It is easy to see that applying the transportation map $S'(i, j)$ to P' yields P_X . Besides, from the procedure adopted to build S' , it is evident that

$$\max_{(i,j):S'(i,j)\neq 0} |i - j| \leq \max_{(i,j):S_{PQ}(i,j)\neq 0} |i - j| \leq L_{max}, \quad (\text{A11})$$

(the only new shipments introduced are from a bin to itself). In addition, the distance between P' and P is, by construction, lower than $\gamma(\lambda)$, which can be made arbitrarily small by decreasing λ , thus completing the proof of the lemma. ■