



**HAL**  
open science

# Sketch \*-metric: Comparing Data Streams via Sketching

Emmanuelle Anceaume, Yann Busnel

► **To cite this version:**

Emmanuelle Anceaume, Yann Busnel. Sketch \*-metric: Comparing Data Streams via Sketching. 2012.  
hal-00721211

**HAL Id: hal-00721211**

**<https://hal.science/hal-00721211v1>**

Submitted on 27 Jul 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Sketch $\star$ -metric: Comparing Data Streams via Sketching

## RESEARCH REPORT

Emmanuelle Anceaume  
IRISA / CNRS, Rennes, France  
Emmanuelle.Anceaume@irisa.fr

Yann Busnel  
LINA / Université de Nantes, Nantes, France  
Yann.Busnel@univ-nantes.fr

**Abstract**—In this paper, we consider the problem of estimating the distance between any two large data streams in small-space constraint. This problem is of utmost importance in data intensive monitoring applications where input streams are generated rapidly. These streams need to be processed on the fly and accurately to quickly determine any deviance from nominal behavior. We present a new metric, the *Sketch  $\star$ -metric*, which allows to define a distance between updatable summaries (or sketches) of large data streams. An important feature of the *Sketch  $\star$ -metric* is that, given a measure on the entire initial data streams, the *Sketch  $\star$ -metric* preserves the axioms of the latter measure on the sketch (such as the non-negativity, the identity, the symmetry, the triangle inequality but also specific properties of the  $f$ -divergence). Extensive experiments conducted on both synthetic traces and real data allow us to validate the robustness and accuracy of the *Sketch  $\star$ -metric*.

**Index Terms**—Data stream; metric; randomized approximation algorithm.

### I. INTRODUCTION

The main objective of this paper is to propose a novel metric that reflects the relationships between any two discrete probability distributions in the context of massive data streams. Specifically, this metric, denoted by *Sketch  $\star$ -metric*, allows us to efficiently estimate a broad class of distances measures between any two large data streams by computing these distances only using compact synopses or sketches of the streams. The *Sketch  $\star$ -metric* is distribution-free and makes no assumption about the underlying data volume. It is thus capable of comparing any two data streams, identifying their correlation if any, and more generally, it allows us to acquire a deep understanding of the structure of the input streams. Formalization of this metric is the first contribution of this paper.

The interest of estimating distances between any two data streams is important in data intensive applications. Many different domains are concerned by such analyses including machine learning, data mining, databases, information retrieval, and network monitoring. In all these applications, it is necessary to quickly and precisely process a huge amount of data. For instance, in IP network management, the analysis of input streams will allow to rapidly detect the presence of outliers or intrusions when changes in the communication patterns occur [1]. In sensors networks, such

an analysis will enable the correlation of geographical and environmental informations [2], [3]. Actually, the problem of detecting changes or outliers in a data stream is similar to the problem of identifying patterns that do not conform to the expected behavior, which has been an active area of research for many decades. For instance, depending on the specificities of the domain considered and the type of outliers considered, different methods have been designed, namely classification-based, clustering-based, nearest neighbor based, statistical, spectral, and information theory. To accurately analyze streams of data, a panel of information-theoretic measures and distances have been proposed to answer the specificities of the analyses. Among them, the most commonly used are the Kullback-Leibler (KL) divergence [4], the  $f$ -divergence introduced by Csiszar, Morimoto and Ali & Silvey [5], [6], [7], the Jensen-Shannon divergence and the Battacharyya distance [8]. More details can be found in the comprehensive survey of Basseville [9].

Unfortunately, computing information theoretic measures of distances in the data stream model is challenging essentially because one needs to process streams on the fly (*i.e.*, in one-pass), on huge amount of data, and by using very little storage with respect to the size of the stream. In addition the analysis must be robust over time to detect any sudden change in the observed streams (which may be the manifestation of routers deny of service attack or worm propagation). We tackle this issue by presenting an approximation algorithm that constructs a sketch of the stream from which the *Sketch  $\star$ -metric* is computed. This algorithm is a one-pass algorithm. It uses very basic computations, and little storage space (*i.e.*,  $\mathcal{O}(t(\log n + k \log m))$  where  $k$  and  $t$  are precision parameters,  $m$  is an upper-bound of stream size and  $n$  the number of distinct items in the stream). It does not need any information on the size of input streams nor on their structure. This consists in the second contribution of the paper.

Finally, the robustness of our approach is validated with a detailed experimentation study based on synthetic traces that range from stable streams to highly skewed ones.

The paper is organized as follows. First, Section II reviews the related work on classical generalized metrics and their applications on the data stream model while Section III describes this model. Section IV presents the necessary background

that makes the paper self-contained. Section V formalizes the *Sketch  $\star$ -metric*. Section VI presents the algorithm that fairly approximates the *Sketch  $\star$ -metric* in one pass and Section VII presents extensive experiments (on both synthetic traces and real data) of our algorithm. Finally, we conclude in Section VIII.

## II. RELATED WORK

Work on data stream analysis mainly focuses on efficient methods (data-structures and algorithms) to answer different kind of queries over massive data streams. Mostly, these methods consist in deriving statistic estimators over the data stream, in creating summary representations of streams (to build histograms, wavelets, and quantiles), and in comparing data streams. Regarding the construction of estimators, a seminal work is due to Alon *et al.* [10]. The authors have proposed estimators of the frequency moments  $F_k$  of a stream, which are important statistical tools that allow to quantify specificities of a data stream. Subsequently, a lot of attention has been paid to the strongly related notion of the entropy of a stream, and all notions based on entropy (*i.e.*, norm and relative entropy) [11]. These notions are essentially related to the quantification of the amount of randomness of a stream (*e.g.*, [12], [13], [14], [15], [16], [17]). The construction of synopses or sketches of the data stream have been proposed for different applications (*e.g.*, [18], [19], [20]).

Distance and divergence measures are key measures in statistical inference and data processing problems [9]. There exists two largely used broad classes of measures, namely the  $f$ -divergences and the Bregman divergences. Among them, there exists two classical distances, namely the Kullback-Leibler (KL) divergence and the Hellinger distance, that are very important to quantify the amount of information that separates two distributions. In [16], the authors have proposed a one pass algorithm for estimating the KL divergence of an observed stream compared to an expected one. Experimental evaluations have shown that the estimation provided by this algorithm is accurate for different adversarial settings for which the quality of other methods dramatically decreases. However, this solution assumes that the expected stream is the uniform one, that is a fully random stream. Actually in [21], the authors propose a characterization of the information divergences that are not sketchable. They have proven that any distance that has not “norm-like” properties is not sketchable. In the present paper, we go one step further by formalizing a metric that allows to efficiently and accurately estimate a broad class of distances measures between any two large data streams by computing these distances uniquely on compact synopses or sketches of streams.

## III. DATA STREAM MODEL

We consider a system in which a node  $P$  receives a large data stream  $\sigma = a_1, a_2, \dots, a_m$  of data items that arrive sequentially. In the following, we describe a single instance of  $P$ , but clearly multiple instances of  $P$  may co-exist in a system (*e.g.*, in case  $P$  represents a router, a base station in a sensor

network). Each data item  $i$  of the stream  $\sigma$  is drawn from the universe  $\Omega = \{1, 2, \dots, n\}$  where  $n$  is very large. Data items can be repeated multiple times in the stream. In the following we suppose that the length  $m$  of the stream is not known. Items in the stream arrive regularly and quickly, and due to memory constraints, need to be processed sequentially and in an online manner. Therefore, node  $P$  can locally store only a small fraction of the items and perform simple operations on them. The algorithms we consider in this work are characterized by the fact that they can approximate some function on  $\sigma$  with a very limited amount of memory. We refer the reader to [22] for a detailed description of data streaming models and algorithms.

## IV. INFORMATION DIVERGENCE OF DATA STREAMS

We first present notations and background that make this paper self-contained.

### A. Preliminaries

A natural approach to study a data stream  $\sigma$  is to model it as an empirical data distribution over the universe  $\Omega$ , given by  $(p_1, p_2, \dots, p_n)$  with  $p_i = x_i/m$ , and  $x_i = |\{j : a_j = i\}|$  representing the number of times data item  $i$  appears in  $\sigma$ . We have  $m = \sum_{i \in \Omega} x_i$ .

1) *Entropy*: Intuitively, the entropy is a measure of the randomness of a data stream  $\sigma$ . The entropy  $H(\sigma)$  is minimum (*i.e.*, equal to zero) when all the items in the stream are the same, and it reaches its maximum (*i.e.*,  $\log_2 m$ ) when all the items in the stream are distinct. Specifically, we have  $H(\sigma) = -\sum_{i \in \Omega} p_i \log_2 p_i$ . The log is to the base 2 and thus entropy is expressed in bits. By convention, we have  $0 \log 0 = 0$ . Note that the number of times  $x_i$  item  $i$  appears in a stream is commonly called the frequency of item  $i$ . The norm of the entropy is defined as  $F_H = \sum_{i \in \Omega} x_i \log x_i$ .

2) *2-universal Hash Functions*: In the following, we intensively use hash functions randomly picked from a 2-universal hash family. A collection  $\mathcal{H}$  of hash functions  $h : \{1, \dots, M\} \rightarrow \{0, \dots, M'\}$  is said to be *2-universal* if for every  $h \in \mathcal{H}$  and for every two different items  $i, j \in [M]$ ,  $\mathbb{P}\{h(i) = h(j)\} \leq \frac{1}{M'}$ , which is exactly the probability of collision obtained if the hash function assigned truly random values to any  $i \in [M]$ . In the following, notation  $[M]$  means  $\{1, \dots, M\}$ .

### B. Metrics and divergences

1) *Metric definitions*: The classical definition of a metric is based on a set of four axioms.

**Definition 1** (Metric) *Given a set  $X$ , a **metric** is a function  $d : X \times X \rightarrow \mathbb{R}$  such that, for any  $x, y, z \in X$ , we have:*

$$\text{Non-negativity: } d(x, y) \geq 0 \quad (1)$$

$$\text{Identity of indiscernibles: } d(x, y) = 0 \Leftrightarrow x = y \quad (2)$$

$$\text{Symmetry: } d(x, y) = d(y, x) \quad (3)$$

$$\text{Triangle inequality: } d(x, y) \leq d(x, z) + d(z, y) \quad (4)$$

In the context of information divergence, usual distance functions are not precisely metric. Indeed, most of divergence

functions do not verify the 4 axioms, but only a subset of them. We recall hereafter some definitions of generalized metrics.

**Definition 2** (Pseudometric) *Given a set  $X$ , a **pseudometric** is a function that verifies the axioms of a metric with the exception of the identity of indiscernible, which is replaced by*

$$\forall x \in X, d(x, x) = 0.$$

Note that this definition allows that  $d(x, y) = 0$  for some  $x \neq y$  in  $X$ .

**Definition 3** (Quasimetric) *Given a set  $X$ , a **quasimetric** is a function that verifies all the axioms of a metric with the exception of the symmetry (cf. Relation 3).*

**Definition 4** (Semimetric) *Given a set  $X$ , a **semimetric** is a function that verifies all the axioms of a metric with the exception of the triangle inequality (cf. Relation 4).*

**Definition 5** (Premetric) *Given a set  $X$ , a **premetric** is a pseudometric that relax both the symmetry and triangle inequality axioms.*

**Definition 6** (Pseudoquasimetric) *Given a set  $X$ , a **pseudoquasimetric** is a function that relax both the identity of indiscernible and the symmetry axioms.*

Note that the latter definition simply corresponds to a premetric satisfying the triangle inequality. Remark also that all the generalized metrics preserve the *non-negativity* axiom.

2) *Divergences*: We now give the definition of two broad classes of generalized metrics, usually denoted as *divergences*.

a) *f-divergence*: Mostly used in the context of statistics and probability theory, a *f-divergence*  $\mathcal{D}_f$  is a premetric that guarantees monotonicity and convexity.

**Definition 7** (*f-divergence*) *Let  $p$  and  $q$  be two  $\Omega$ -point distributions. Given a convex function  $f : (0, \infty) \rightarrow \mathbb{R}$  such that  $f(1) = 0$ , the **f-divergence** of  $q$  from  $p$  is:*

$$\mathcal{D}_f(p||q) = \sum_{i \in \Omega} q_i f\left(\frac{p_i}{q_i}\right),$$

where by convention  $0f\left(\frac{0}{0}\right) = 0$ ,  $af\left(\frac{0}{a}\right) = a \lim_{u \rightarrow 0} f(u)$ , and  $0f\left(\frac{a}{0}\right) = a \lim_{u \rightarrow \infty} f(u)/u$  if these limits exist.

Following this definition, any *f-divergence* verifies both monotonicity and convexity.

**Property 8** (Monotonicity) *Given  $\kappa$  an arbitrary transition probability that respectively transforms two  $\Omega$ -point distributions  $p$  and  $q$  into  $p_\kappa$  and  $q_\kappa$ , we have:*

$$\mathcal{D}_f(p||q) \geq \mathcal{D}_f(p_\kappa||q_\kappa).$$

**Property 9** (Convexity) *Let  $p_1, p_2, q_1$  and  $q_2$  be four  $\Omega$ -point distributions. Given any  $\lambda \in [0, 1]$ , we have:*

$$\begin{aligned} \mathcal{D}_f(\lambda p_1 + (1 - \lambda)p_2 || \lambda q_1 + (1 - \lambda)q_2) \\ \leq \lambda \mathcal{D}_f(p_1 || q_1) + (1 - \lambda) \mathcal{D}_f(p_2 || q_2). \end{aligned}$$

This class of divergences has been introduced in independent works by chronologically Csiszár, Morimoto and Ali & Silvey [5], [6], [7]. All the distance measures in the so-called *Ali-Silvey distances* are applicable to quantifying statistical differences between data streams.

b) *Bregman divergence*: Initially proposed in [23], this class of generalized metrics encloses quasimetrics and semimetrics, as these divergences do not satisfy the triangle inequality nor symmetry.

**Definition 10** (Bregman divergence) *Given  $F$  a continuously-differentiable and strictly convex function defined on a closed convex set  $C$ , the **Bregman divergence** associated with  $F$  for  $p, q \in C$  is defined as*

$$\mathcal{B}_F(p||q) = F(p) - F(q) - \langle \nabla F(q), (p - q) \rangle.$$

where the operator  $\langle \cdot, \cdot \rangle$  denotes the inner product.

In the context of data stream, it is possible to reformulate this definition according to probability theory. Specifically,

**Definition 11** (Decomposable Bregman divergence) *Let  $p$  and  $q$  be two  $\Omega$ -point distributions. Given a strictly convex function  $F : (0, 1] \rightarrow \mathbb{R}$ , the **Bregman divergence** associated with  $F$  of  $q$  from  $p$  is defined as*

$$\mathcal{B}_F(p||q) = \sum_{i \in \Omega} (F(p_i) - F(q_i) - (p_i - q_i)F'(q_i)).$$

Following these definitions, any Bregman divergence verifies non-negativity and convexity in its first argument, but not necessarily in the second argument. Another interesting property is given by thinking of the Bregman divergences as an operator of the function  $F$ .

**Property 12** (Linearity) *Let  $F_1$  and  $F_2$  be two strictly convex and differentiable functions. Given any  $\lambda \in [0, 1]$ , we have that*

$$\mathcal{B}_{F_1 + \lambda F_2}(p||q) = \mathcal{B}_{F_1}(p||q) + \lambda \mathcal{B}_{F_2}(p||q).$$

3) *Classical metrics*: In this section, we present several commonly used metrics in  $\Omega$ -point distribution context. These specific metrics are used in the evaluation part presented in Section VII.

a) *Kullback-Leibler divergence*: The Kullback-Leibler (KL) divergence [4], also called the relative entropy, is a robust metric for measuring the statistical difference between two data streams. The KL divergence owns the special feature that it is both a *f-divergence* and a Bregman one (with  $f(t) = F(t) = t \log t$ ).

Given  $p$  and  $q$  two  $\Omega$ -point distributions, the Kullback-Leibler divergence is then defined as

$$\mathcal{D}_{KL}(p||q) = \sum_{i \in \Omega} p_i \log \frac{p_i}{q_i} = H(p, q) - H(p), \quad (5)$$

where  $H(p) = -\sum p_i \log p_i$  is the (empirical) entropy of  $p$  and  $H(p, q) = -\sum p_i \log q_i$  is the cross entropy of  $p$  and  $q$ .

b) *Jensen-Shannon divergence*: The Jensen-Shannon divergence (JS) is a symmetrized and smoothed version of the Kullback-Leibler divergence. Also known as information radius (IRad) or total divergence to the average, it is defined as

$$\mathcal{D}_{JS}(p||q) = \frac{1}{2} \mathcal{D}_{KL}(p||\ell) + \frac{1}{2} \mathcal{D}_{KL}(q||\ell), \quad (6)$$

where  $\ell = \frac{1}{2}(p+q)$ . Note that the square root of this divergence is a metric.

c) *Bhattacharyya distance*: The Bhattacharyya distance is derived from his proposed measure of similarity between two multinomial distributions, known as the Bhattacharyya coefficient (BC) [8]. It is defined as

$$\mathcal{D}_B(p||q) = -\log(BC(p, q)) \text{ where } BC(p, q) = \sum_{i \in \Omega} \sqrt{p_i q_i}.$$

This distance is a semimetric as it does not verify the triangle inequality. Note that the famous Hellinger distance [24] equal to  $\sqrt{1 - BC(p, q)}$  verifies it.

## V. SKETCH $\star$ -METRIC

We now present a method to sketch two input data streams  $\sigma_1$  and  $\sigma_2$ , and to compute any generalized metric  $\phi$  between these sketches such that this computation preserves all the properties of  $\phi$  computed on  $\sigma_1$  and  $\sigma_2$ . Proof of correctness of this method is presented in this section.

**Definition 13** (Sketch  $\star$ -metric) *Let  $p$  and  $q$  be any two  $n$ -point distributions. Given a precision parameter  $k$ , and any generalized metric  $\phi$  on the set of all  $\Omega$ -point distributions, there exists a **Sketch  $\star$ -metric**  $\hat{\phi}_k$  defined as follows*

$$\hat{\phi}_k(p||q) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{p}_\rho||\hat{q}_\rho) \text{ with } \forall a \in \rho, \hat{p}_\rho(a) = \sum_{i \in a} p(i),$$

where  $\mathcal{P}_k(\Omega)$  is the set of all partitions of an  $\Omega$ -element set into exactly  $k$  nonempty and mutually exclusive cells.

**Remark 14** *Note that for  $k > \Omega$ , it does not exist a partition of  $\Omega$  into  $k$  nonempty parts. By convention, we consider that  $\hat{\phi}_k(p||q) = \phi(p||q)$  in this specific context.*

In this section, we focus on the preservation of axioms and properties of a generalized metric  $\phi$  by the corresponding Sketch  $\star$ -metric  $\hat{\phi}_k$ .

### A. Axioms preserving

**Theorem 15** *Given any generalized metric  $\phi$  then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\hat{\phi}_k$  preserves all the axioms of  $\phi$ .*

*Proof*: The proof is directly derived from Lemmata 16, 17, 18 and 19. ■

**Lemma 16** (Non-negativity) *Given any generalized metric  $\phi$  verifying the Non-negativity axiom then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\hat{\phi}_k$  preserves the Non-negativity axiom.*

*Proof*: Let  $p$  and  $q$  be any two  $\Omega$ -point distributions. By definition,

$$\hat{\phi}_k(p||q) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{p}_\rho||\hat{q}_\rho)$$

As for any two  $k$ -point distributions,  $\phi$  is positive we have  $\hat{\phi}_k(p||q) \geq 0$  that concludes the proof. ■

**Lemma 17** (Identity of indiscernible) *Given any generalized metric  $\phi$  verifying the Identity of indiscernible axiom then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\hat{\phi}_k$  preserves the Identity of indiscernible axiom.*

*Proof*: Let  $p$  be any  $\Omega$ -point distribution. We have

$$\hat{\phi}_k(p||p) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{p}_\rho||\hat{p}_\rho) = 0,$$

due to  $\phi$  Identity of indiscernible axiom.

Consider now two  $\Omega$ -point distributions  $p$  and  $q$  such that  $\hat{\phi}_k(p||q) = 0$ . Metric  $\phi$  verifies both the non-negativity axiom (by construction) and the Identity of indiscernible axiom (by assumption). Thus we have  $\forall \rho \in \mathcal{P}_k(\Omega), \hat{p}_\rho = \hat{q}_\rho$ , leading to

$$\forall \rho \in \mathcal{P}_k(\Omega), \forall a \in \rho, \sum_{i \in a} p(i) = \sum_{i \in a} q(i). \quad (7)$$

Moreover, for any  $i \in \Omega$ , there exists a partition  $\rho \in \mathcal{P}_k(\Omega)$  such that  $\{i\} \in \rho$ . By Equation 7,  $\forall i \in \Omega, p(i) = q(i)$ , and so  $p = q$ .

Combining the two parts of the proof leads to  $\hat{\phi}_k(p||q) = 0 \iff p = q$ , which concludes the proof of the Lemma. ■

**Lemma 18** (Symmetry) *Given any generalized metric  $\phi$  verifying the Symmetry axiom then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\hat{\phi}_k$  preserves the Symmetry axiom.*

*Proof*: Let  $p$  and  $q$  be any two  $\Omega$ -point distributions. We have

$$\hat{\phi}_k(p||q) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{p}_\rho||\hat{q}_\rho).$$

Let  $\bar{\rho} \in \mathcal{P}_k(\Omega)$  be a  $k$ -cell partition such that  $\phi(\hat{p}_{\bar{\rho}}||\hat{q}_{\bar{\rho}}) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{p}_\rho||\hat{q}_\rho)$ . We get

$$\hat{\phi}_k(p||q) = \phi(\hat{p}_{\bar{\rho}}||\hat{q}_{\bar{\rho}}) = \phi(\hat{q}_{\bar{\rho}}||\hat{p}_{\bar{\rho}}) \leq \hat{\phi}_k(q||p).$$

By symmetry, considering  $\underline{\rho} \in \mathcal{P}_k(\Omega)$  such that  $\phi(\hat{q}_{\underline{\rho}}||\hat{p}_{\underline{\rho}}) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{q}_\rho||\hat{p}_\rho)$ , we also have  $\hat{\phi}_k(q||p) \leq \hat{\phi}_k(p||q)$ , which concludes the proof. ■

**Lemma 19** (Triangle inequality) *Given any generalized metric  $\phi$  verifying the Triangle inequality axiom then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\hat{\phi}_k$  preserves the Triangle inequality axiom.*

*Proof*: Let  $p, q$  and  $r$  be any three  $\Omega$ -point distributions. Let  $\bar{\rho} \in \mathcal{P}_k(\Omega)$  be a  $k$ -cell partition such that  $\phi(\hat{p}_{\bar{\rho}}||\hat{q}_{\bar{\rho}}) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{p}_\rho||\hat{q}_\rho)$ .

We have

$$\begin{aligned} \hat{\phi}_k(p||q) &= \phi(\hat{p}_{\bar{\rho}}||\hat{q}_{\bar{\rho}}) \\ &\leq \phi(\hat{p}_{\bar{\rho}}||\hat{r}_{\bar{\rho}}) + \phi(\hat{r}_{\bar{\rho}}||\hat{q}_{\bar{\rho}}) \\ &\leq \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{p}_\rho||\hat{r}_\rho) + \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\hat{r}_\rho||\hat{q}_\rho) \\ &= \hat{\phi}_k(p||r) + \hat{\phi}_k(r||q) \end{aligned}$$

that concludes the proof. ■

### B. Properties preserving

**Theorem 20** *Given a  $f$ -divergence  $\phi$  then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\hat{\phi}_k$  is also a  $f$ -divergence.*

*Proof*: From Theorem 15,  $\hat{\phi}_k$  preserves the axioms of the generalized metric. Thus,  $\hat{\phi}_k$  and  $\phi$  are in the same equivalence class. Moreover, from Lemma 22,  $\hat{\phi}_k$  verifies the monotonicity

property. Thus, as the  $f$ -divergence is the only class of decomposable information *monotonic* divergences (cf. [25]),  $\widehat{\phi}_k$  is also a  $f$ -divergence. ■

**Theorem 21** *Given a Bregman divergence  $\phi$  then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\widehat{\phi}_k$  is also a Bregman divergence.*

*Proof:* From Theorem 15,  $\widehat{\phi}_k$  preserves the axioms of the generalized metric. Thus,  $\widehat{\phi}_k$  and  $\phi$  are in the same equivalence class. Moreover, the Bregman divergence is characterized by the property of transitivity (cf. [26]) defined as follows. Given  $p, q$  and  $r$  three  $\Omega$ -point distributions such that  $q = \Pi(L|r)$  and  $p \in L$ , with  $\Pi$  is a selection rule according to the definition of Csiszár in [26] and  $L$  is a subset of the  $\Omega$ -point distributions, we have the Generalized Pythagorean Theorem:

$$\phi(p||q) + \phi(q||r) = \phi(p||r).$$

Moreover the authors in [27] show that the set  $S_n$  of all discrete probability distributions over  $n$  elements ( $X = \{x_1, \dots, x_n\}$ ) is a Riemannian manifold, and it owns another different dually flat affine structure. They also show that these dual structures give rise to the generalized Pythagorean theorem. This is verified for the coordinates in  $S_n$  and for the dual coordinates [27]. Combining these results with the projection theorem [26], [27], we obtain that

$$\begin{aligned} \widehat{\phi}_k(p||r) &= \max_{\rho \in \mathcal{P}_k(n)} \phi(\widehat{p}_\rho||\widehat{r}_\rho) \\ &= \max_{\rho \in \mathcal{P}_k(n)} (\phi(\widehat{p}_\rho||\widehat{q}_\rho) + \phi(\widehat{q}_\rho||\widehat{r}_\rho)) \\ &= \max_{\rho \in \mathcal{P}_k(n)} \phi(\widehat{p}_\rho||\widehat{q}_\rho) + \max_{\rho \in \mathcal{P}_k(n)} \phi(\widehat{q}_\rho||\widehat{r}_\rho) \\ &= \widehat{\phi}_k(p||q) + \widehat{\phi}_k(q||r) \end{aligned}$$

Finally, by the characterization of Bregman divergence through transitivity [26], and reinforced with Lemma 24 statement,  $\widehat{\phi}_k$  is also a Bregman divergence. ■

In the following, we show that the *Sketch  $\star$ -metric* preserves the properties of divergences.

**Lemma 22** (Monotonicity) *Given any generalized metric  $\phi$  verifying the Monotonicity property then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\widehat{\phi}_k$  preserves the Monotonicity property.*

*Proof:* Let  $p$  and  $q$  be any two  $\Omega$ -point distributions. Given  $c < N$ , consider a partition  $\mu \in \mathcal{P}_c(\Omega)$ . As  $\phi$  is monotonic, we have  $\phi(p||q) \geq \phi(\widehat{p}_\mu||\widehat{q}_\mu)$  [28]. We split the proof into two cases:

Case (1). Suppose that  $c \geq k$ . Computing  $\widehat{\phi}_k(\widehat{p}_\mu||\widehat{q}_\mu)$  amounts in considering only the  $k$ -cell partitions  $\rho \in \mathcal{P}_k(\Omega)$  that verify

$$\forall b \in \mu, \exists a \in \rho, b \subseteq a.$$

These partitions form a subset of  $\mathcal{P}_k(\Omega)$ . The maximal value of  $\phi(\widehat{p}_\rho||\widehat{q}_\rho)$  over this subset cannot be greater than the maximal value over the whole  $\mathcal{P}_k(\Omega)$ . Thus we have

$$\widehat{\phi}_k(p||q) = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\widehat{p}_\rho||\widehat{q}_\rho) \geq \widehat{\phi}_k(\widehat{p}_\mu||\widehat{q}_\mu).$$

Case (2). Suppose now that  $c < k$ . By definition, we have  $\widehat{\phi}_k(\widehat{p}_\mu||\widehat{q}_\mu) = \phi(\widehat{p}_\mu||\widehat{q}_\mu)$ . Consider  $\rho' \in \mathcal{P}_k(\Omega)$  such that  $\forall b \in \mu, \exists a \in \rho', b \subseteq a$ . It then exists a transition probability that respectively transforms  $\widehat{p}_{\rho'}$  and  $\widehat{q}_{\rho'}$  into  $\widehat{p}_\mu$  and  $\widehat{q}_\mu$ . As  $\phi$  is monotonic, we have

$$\begin{aligned} \widehat{\phi}_k(p||q) &= \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\widehat{p}_\rho||\widehat{q}_\rho) \\ &\geq \phi(\widehat{p}_{\rho'}||\widehat{q}_{\rho'}) \\ &\geq \phi(\widehat{p}_\mu||\widehat{q}_\mu) = \widehat{\phi}_k(\widehat{p}_\mu||\widehat{q}_\mu). \end{aligned}$$

Finally for any value of  $c$ ,  $\widehat{\phi}_k$  guarantees the monotonicity property. This concludes the proof. ■

**Lemma 23** (Convexity) *Given any generalized metric  $\phi$  verifying the Convexity property then, for any  $k \in \Omega$ , the corresponding Sketch  $\star$ -metric  $\widehat{\phi}_k$  preserves the Convexity property.*

*Proof:* Let  $p_1, p_2, q_1$  and  $q_2$  be any four  $\Omega$ -point distributions. Given any  $\lambda \in [0, 1]$ , we have:

$$\begin{aligned} \widehat{\phi}_k(\lambda p_1 + (1-\lambda)p_2||\lambda q_1 + (1-\lambda)q_2) \\ = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\lambda \widehat{p}_{1\rho} + (1-\lambda)\widehat{p}_{2\rho}||\lambda \widehat{q}_{1\rho} + (1-\lambda)\widehat{q}_{2\rho}) \end{aligned}$$

Let  $\rho' \in \mathcal{P}_k(\Omega)$  such that

$$\begin{aligned} \phi(\lambda \widehat{p}_{1\rho'} + (1-\lambda)\widehat{p}_{2\rho'}||\lambda \widehat{q}_{1\rho'} + (1-\lambda)\widehat{q}_{2\rho'}) \\ = \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\lambda \widehat{p}_{1\rho} + (1-\lambda)\widehat{p}_{2\rho}||\lambda \widehat{q}_{1\rho} + (1-\lambda)\widehat{q}_{2\rho}). \end{aligned}$$

As  $\phi$  verifies the Convex property, we have:

$$\begin{aligned} \widehat{\phi}_k(\lambda p_1 + (1-\lambda)p_2||\lambda q_1 + (1-\lambda)q_2) \\ = \phi(\lambda \widehat{p}_{1\rho'} + (1-\lambda)\widehat{p}_{2\rho'}||\lambda \widehat{q}_{1\rho'} + (1-\lambda)\widehat{q}_{2\rho'}) \\ \leq \lambda \phi(\widehat{p}_{1\rho'}||\widehat{q}_{1\rho'}) + (1-\lambda)\phi(\widehat{p}_{2\rho'}||\widehat{q}_{2\rho'}) \\ \leq \lambda \left( \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\widehat{p}_{1\rho}||\widehat{q}_{1\rho}) \right) + (1-\lambda) \left( \max_{\rho \in \mathcal{P}_k(\Omega)} \phi(\widehat{p}_{2\rho}||\widehat{q}_{2\rho}) \right) \\ = \lambda \widehat{\phi}_k(p_1||q_1) + (1-\lambda)\widehat{\phi}_k(p_2||q_2) \end{aligned}$$

that concludes the proof. ■

**Lemma 24** (Linearity) *The Sketch  $\star$ -metric definition preserves the Linearity property.*

*Proof:* Let  $F_1$  and  $F_2$  be two strictly convex and differentiable functions, and any  $\lambda \in [0, 1]$ . Consider the three Bregman divergences generated respectively from  $F_1, F_2$  and  $F_1 + \lambda F_2$ .

Let  $p$  and  $q$  be two  $\Omega$ -point distributions. We have:

$$\begin{aligned} \widehat{\mathcal{B}}_{F_1 + \lambda F_2}(p||q) &= \max_{\rho \in \mathcal{P}_k(\Omega)} \mathcal{B}_{F_1 + \lambda F_2}(\widehat{p}_\rho||\widehat{q}_\rho) \\ &= \max_{\rho \in \mathcal{P}_k(n)} (\mathcal{B}_{F_1}(\widehat{p}_\rho||\widehat{q}_\rho) + \lambda \mathcal{B}_{F_2}(\widehat{p}_\rho||\widehat{q}_\rho)) \\ &\leq \widehat{\mathcal{B}}_{F_1}(p||q) + \lambda \widehat{\mathcal{B}}_{F_2}(p||q) \end{aligned}$$

As  $F_1$  and  $F_2$  two strictly convex functions, and taken a leaf out of the Jensen's inequality, we have:

$$\begin{aligned} \widehat{\mathcal{B}}_{F_1}(p||q) + \lambda \widehat{\mathcal{B}}_{F_2}(p||q) \\ \leq \max_{\rho \in \mathcal{P}_k(\Omega)} (\mathcal{B}_{F_1}(\widehat{p}_\rho||\widehat{q}_\rho) + \lambda \mathcal{B}_{F_2}(\widehat{p}_\rho||\widehat{q}_\rho)) \\ = \widehat{\mathcal{B}}_{F_1 + \lambda F_2}(p||q) \end{aligned}$$

that conclude the proof.  $\blacksquare$

This concludes the proof that the *Sketch  $\star$ -metric* preserves all the axioms of a metric as well as the properties of  $f$ -divergences and Bregman divergences. We now show how to efficiently implement such a metric.

## VI. APPROXIMATION ALGORITHM

In this section, we propose an algorithm that computes the *Sketch  $\star$ -metric* in one pass on the stream. By definition of the metric (cf., Definition 13), we need to generate all the possible  $k$ -cell partitions. The number of these partitions follows the Stirling numbers of the second kind, which is equal to  $S(n, k) = \frac{1}{k!} \sum_{j=0}^k (-1)^{k-j} \binom{k}{j} j^n$ , where  $n$  is the size of the items universe. Therefore,  $S(n, k)$  grows exponentially with  $n$ . Unfortunately,  $n$  is very large. As the generating function of  $S(n, k)$  is equivalent to  $x^n$ , it is unreasonable in term of space complexity. We show in the following that generating  $t = \lceil \log(1/\delta) \rceil$  random  $k$ -cell partitions, where  $\delta$  is the probability of error of our randomized algorithm, is sufficient to guarantee good overall performance of our metric.

Our algorithm is inspired from the Count-Min Sketch algorithm proposed in [19] by Cormode and Muthukrishnan. Specifically, the Count-Min algorithm is an  $(\varepsilon, \delta)$ -approximation algorithm that solves the *frequency-estimation* problem. For any items in the input stream  $\sigma$ , the algorithm outputs an estimation  $\hat{f}_v$  of the frequency of item  $v$  such that  $\mathbb{P}\{|\hat{f}_v - f_v| > \varepsilon f_v\} < \delta$ , where  $\varepsilon, \delta > 0$  are given as parameters of the algorithm. The estimation is computed by maintaining a two-dimensional array  $C$  of  $t \times k$  counters, and by using  $t$  2-universal hash functions  $h_i$  ( $1 \leq i \leq t$ ), where  $k = 2/\varepsilon$  and  $t = \lceil \log(1/\delta) \rceil$ . Each time an item  $v$  is read from the input stream, this causes one counter of each line to be incremented, i.e.,  $C[h_i(v)]$  is incremented by one for each  $i \in [1..t]$ .

To compute the *Sketch  $\star$ -metric* of two streams  $\sigma_1$  and  $\sigma_2$ , two sketches  $\hat{\sigma}_1$  and  $\hat{\sigma}_2$  of these streams are constructed according to the above description. Note that there is no particular assumption on the length of both streams  $\sigma_1$  and  $\sigma_2$ . That is their respective length is finite but unknown. By construction of the 2-universal hash functions  $h_i$  ( $1 \leq i \leq t$ ), each line of  $C_{\sigma_1}$  and  $C_{\sigma_2}$  corresponds to one partition  $\rho_i$  of the  $N$ -point empirical distributions of both  $\sigma_1$  and  $\sigma_2$ . Thus when a query is issued to compute the given distance  $\phi$  between these two streams, the maximal value over all the  $t$  partitions  $\rho_i$  of the distance  $\phi$  between  $\hat{\sigma}_{1_{\rho_i}}$  and  $\hat{\sigma}_{2_{\rho_i}}$  is returned, i.e., the distance  $\phi$  applied to the  $i^{\text{th}}$  lines of  $C_{\sigma_1}$  and  $C_{\sigma_2}$  for  $1 \leq i \leq t$ . Figure 1 presents the pseudo-code of our algorithm.

**Lemma 25** *Given parameters  $k$  and  $t$ , Algorithm 1 gives an approximation of the Sketch  $\star$ -metric, using  $\mathcal{O}(t(\log n + k \log m))$  bits of space.*

*Proof:* The matrices  $C_{\sigma_i}$ , for any  $i \in \{1, 2\}$ , are composed of  $t \times k$  counters, which uses  $\mathcal{O}(\log m)$ . On the other hand, with a suitable choice of hash family, we can store the hash functions above in  $\mathcal{O}(t \log n)$  space.  $\blacksquare$

Figure 1. Sketch  $\star$ -metric algorithm

---

**Input:** Two input streams  $\sigma_1$  and  $\sigma_2$ ; the distance  $\phi$ ,  $k$  and  $t$  settings;

**Output:** The distance  $\bar{\phi}$  between  $\sigma_1$  and  $\sigma_2$

- 1 Choose  $t$  functions  $h : [n] \rightarrow [k]$ , each from a 2-universal hash function family;
- 2  $C_{\sigma_1}[1..t][1..k] \leftarrow 0$ ;
- 3  $C_{\sigma_2}[1..t][1..k] \leftarrow 0$ ;
- 4 **for**  $a_j \in \sigma_1$  **do**
- 5      $v = a_j$ ;
- 6     **for**  $i = 1$  **to**  $t$  **do**
- 7          $C_{\sigma_1}[i][h_i(v)] \leftarrow C_{\sigma_1}[i][h_i(v)] + 1$ ;
- 8 **for**  $a_j \in \sigma_2$  **do**
- 9      $w = a_j$ ;
- 10     **for**  $i = 1$  **to**  $t$  **do**
- 11          $C_{\sigma_2}[i][h_i(w)] \leftarrow C_{\sigma_2}[i][h_i(w)] + 1$ ;
- 12 On query  $\bar{\phi}_k(\sigma_1 || \sigma_2)$  **return**  
 $\bar{\phi} = \max_{1 \leq i \leq t} \phi(C_{\sigma_1}[i][-], C_{\sigma_2}[i][-])$ ;

---

## VII. PERFORMANCE EVALUATION

We have implemented our *Sketch  $\star$ -metric* and have conducted a series of experiments on different types of streams and for different parameters settings. We have fed our algorithm with both real-world data and synthetic traces. Real data give a realistic representation of some real systems, while the latter ones allow to capture phenomenon which may be difficult to obtain from real-world traces, and thus allow to check the robustness of our metric. We have varied all the significant parameters of our algorithm, that is, the maximal number of distinct data items  $n$  in each stream, the number of cells  $k$  of each generated partition, and the number of generated partitions  $t$ . For each parameters setting, we have conducted and averaged 100 trials of the same experiment, leading to a total of more than 300,000 experiments for the evaluation of our metric. Real data have been downloaded from the repository of Internet network traffic [29]. We have used five traces among the available ones. Two of them represent two weeks logs of HTTP requests to the Internet service provider ClarkNet WWW server – ClarkNet is a full Internet access provider for the Metro Baltimore-Washington DC area – the other two ones contain two months of HTTP requests to the NASA Kennedy Space Center WWW server, and the last one represents seven months of HTTP requests to the WWW server of the University of Saskatchewan, Canada. In the following these data sets will be respectively referred to as ClarkNet, NASA, and Saskatchewan traces. Table I presents the statistics of these data traces, in term of stream size (cf. “# items” in the table), number of distinct items in each stream (cf. “# distinct items”) and the number of occurrences of the most frequent item (cf. “max. freq.”). For more information on these data traces, an extensive analysis is available in [30]. We

Data trace	# items	# distinct items	max. freq.
NASA (July)	1,891,715	81,983	17,572
NASA (August)	1,569,898	75,058	6,530
ClarkNet (August)	1,654,929	90,516	6,075
ClarkNet (September)	1,673,794	94,787	7,239
Saskatchewan	2,408,625	162,523	52,695

Table I  
STATISTICS OF REAL DATA TRACES.

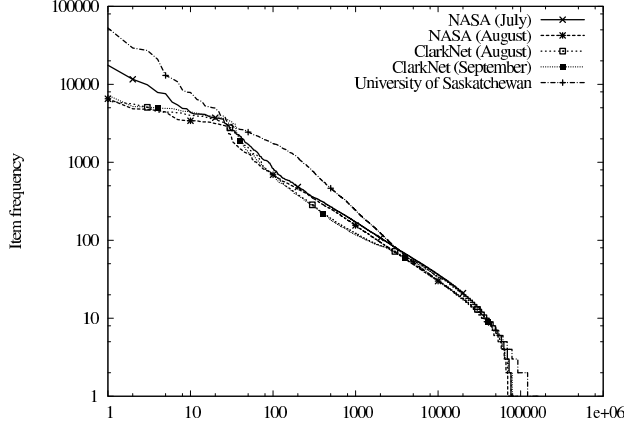


Figure 2. Logscale distribution of frequencies for each real data trace.

have evaluated the accuracy of our metric by comparing for each data set (real and synthetic), the results obtained with our algorithm on the stream sketches (referred to as *Sketch* in the legend) and the ones obtained on full streams (referred to as *Ref* distance in the legend). That is, for each couple of input streams, and for each generalized metric  $\phi$ , we have computed both the exact distance between the two streams and the one as generated by our metric. We now present the main lessons drawn from these experiments.

Figure 3 and 4 shows the accuracy of our metric as a function of the different input streams and the different generalized metrics applied on these streams. All the histograms shown in Figures 3(a)–4(e) share the same legend, but for readability reasons, this legend is only indicated on histogram 3(c). Three generalized metrics have been used, namely the Bhattacharyya distance, the Kullback-Leibler and the Jensen-Shannon divergences, and five distribution families denoted by  $p$  and  $q$  have been compared with these metrics.

Let us focus on synthetic traces. The first noticeable remark is that our metric behaves perfectly well when the two compared streams follow the same distribution, whatever the generalized metric  $\phi$  used (*cf.*, Figure 3(a) with the uniform distribution, Figures 3(b), 3(d) and 3(f) with the Zipf distribution, Figure 3(c) with the Pascal distribution, Figure 3(e) with the Binomial distribution, and Figure 3(g) with the Poisson one). This result is interesting as it allows the sketch  $\star$ -metric to be a very good candidate as a parametric method for making inference about the parameters of the distribution that follow an input stream. The general tendency is that when the distributions of input streams are close (*e.g.*, Zipf distribution with different parameter, Pascal and the Zipf with  $\alpha = 4$ ), then applying the generalized metrics  $\phi$  on sketches give a good

estimation of the distance as computed on the full streams.

Now, when the two input distributions exhibit a totally different shape, this may have an impact on the precision of our metric. Specifically, let us consider as input distributions the Uniform and the Pascal distributions (see Figure 3(a) and 3(c)). Sketching the Uniform distribution leads to  $k$ -cell partitions whose value is well distributed, that is, for a given partition all the  $k$  cell values have with high probability the same value. Now, when sketching the Pascal distribution, the repartition of the data items in the cells of any given partitions is such that a few number of data items (those with high frequency) populate a very few number of cells. However, the values of these cells is very large compared to the other cells, which are populated by a large number of data items whose frequency is small. Thus the contribution of data items exhibiting a small frequency and sharing the cells of highly frequent items will be biased compared to the contribution of the other items. This explains why the accuracy of the sketch  $\star$ -metric is slightly lowered in these cases.

We can also observe the strong impact of the non-symmetry of the Kullback-Leibler divergence on the computation of the distance (computed on full streams or on sketches) with a clear influence when the input streams follow a Pascal and Zipf with  $\alpha = 1$  distributions (see Figures 3(c) and 3(b)).

Finally, Figure 3(h) summarizes the good properties of our method whatever the input streams to be compared and the generalized metric  $\phi$  used to do this comparison.

The same general remarks hold when considering real data sets. Indeed, Figure 4 shows that when the input streams are close to each other, which is the case for both (July and August) NASA and (August and September) ClarkNet traces (*cf.* Figure 2), then applying the generalized metrics  $\phi$  on sketches gives good results w.r.t. full streams. When the shapes of the input streams are different (which is the case for Saskatchewan w.r.t. the 4 other input streams), the accuracy of the sketch  $\star$ -metric decreases a little bit but in a small proportion. Notice that the scales on the y-axis differ significantly in Figure 3 and in Figure 4

Figure 5 presents the impact of the number of cells per generated partition on the accuracy of our metric on both synthetic traces and real data. It clearly shows that, by increasing  $k$ , the number of data items per cell in the generated partition shrinks and thus the absolute error on the computation of the distance decreases. The same feature appears when the number  $n$  of distinct data items in the stream increases. Indeed, when  $n$  increases (for a given  $k$ ), the number data items per cell augments and thus the precision of our metric decreases. This gives rise to a shift of the inflection point, as illustrated in Figure 5(b), due to the fact that data sets have almost twenty times more distinct data items than the synthetic ones. As aforementioned, the input streams exhibit very different shapes which explain the strong impact of  $k$ . Note also that  $k$  has the same influence on the *Sketch*  $\star$ -metric for all the generalized distances  $\phi$ .



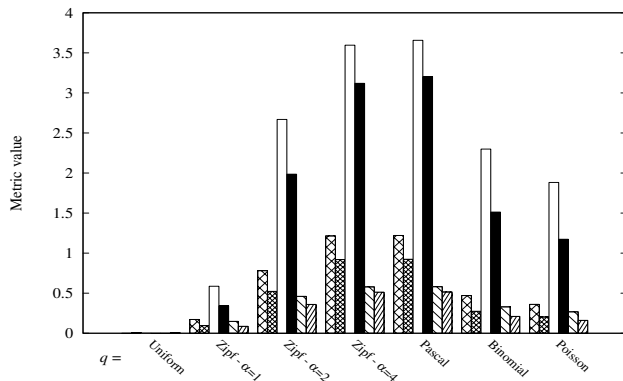
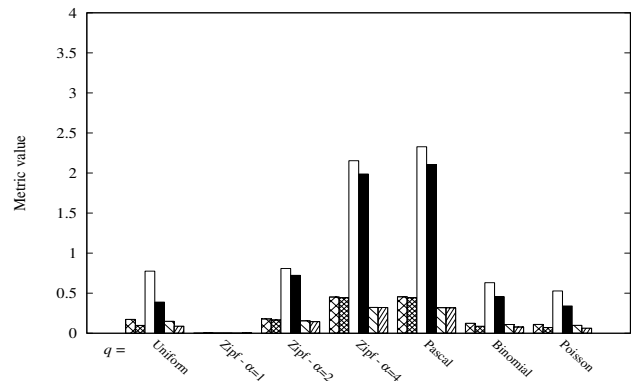
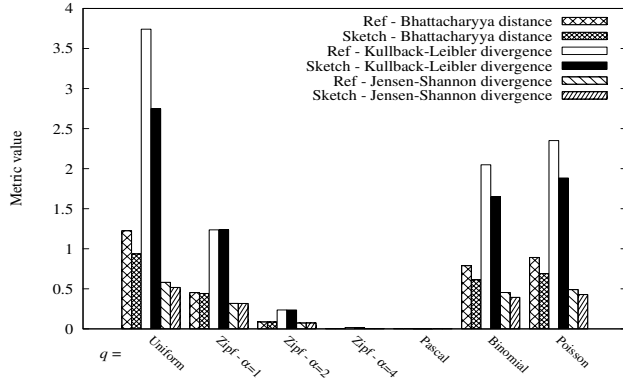
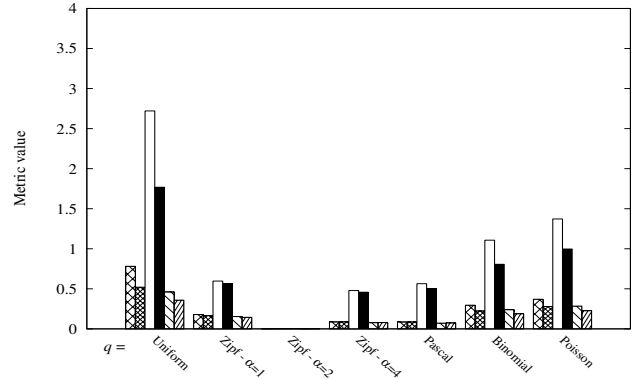
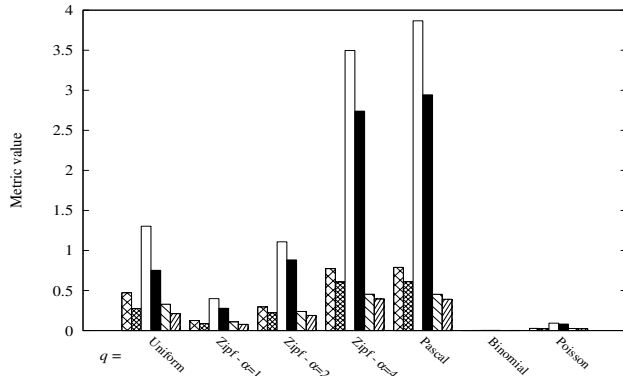
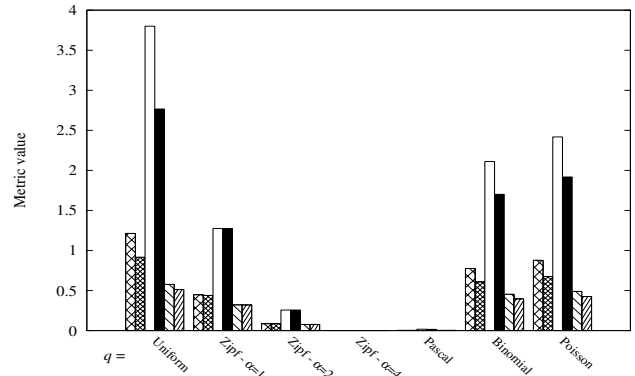
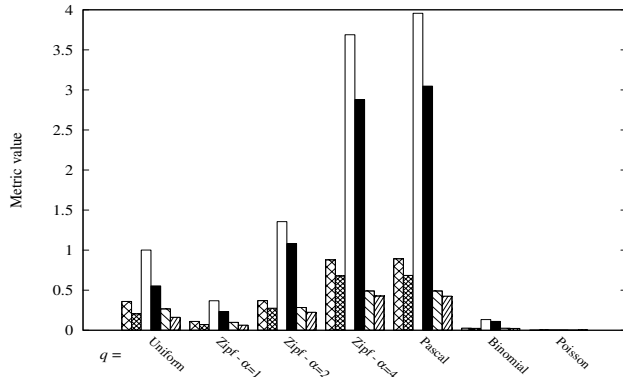
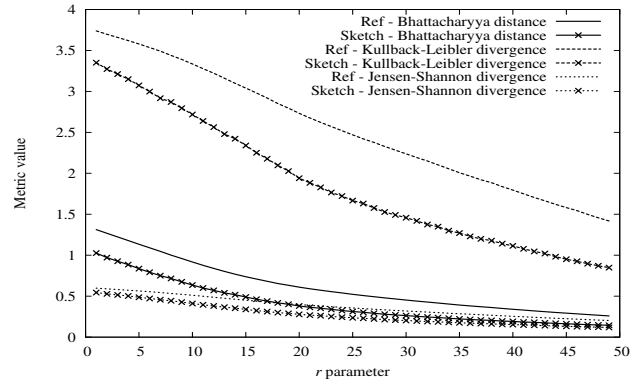
(a)  $p = \text{Uniform distribution}$ (b)  $p = \text{Zipf distribution with } \alpha = 1$ (c)  $p = \text{Pascal distribution with } r = 3 \text{ and } p = \frac{n}{2r+n}$ (d)  $p = \text{Zipf distribution with } \alpha = 2$ (e)  $p = \text{Binomial distribution with } p = 0.5$ (f)  $p = \text{Zipf distribution with } \alpha = 4$ (g)  $p = \text{Poisson distribution with } p = \frac{n}{2}$ (h)  $p = \text{Uniform distribution and } q = \text{Pascal distribution, as a function of its parameter } r (p = \frac{n}{2r+n})$ .

Figure 3. Sketch  $\star$ -metric accuracy as a function of  $p$  and  $q$  (or  $r$  for 3(h)). Parameters setting is as follows:  $m = 200,000$ ;  $n = 4,000$ ;  $k = 200$ ;  $t = 4$  where  $m$  represents the size of the stream,  $n$  the number of distinct data items in the stream,  $t$  the number of generated partitions and  $k$  the number of cells per generated partition.

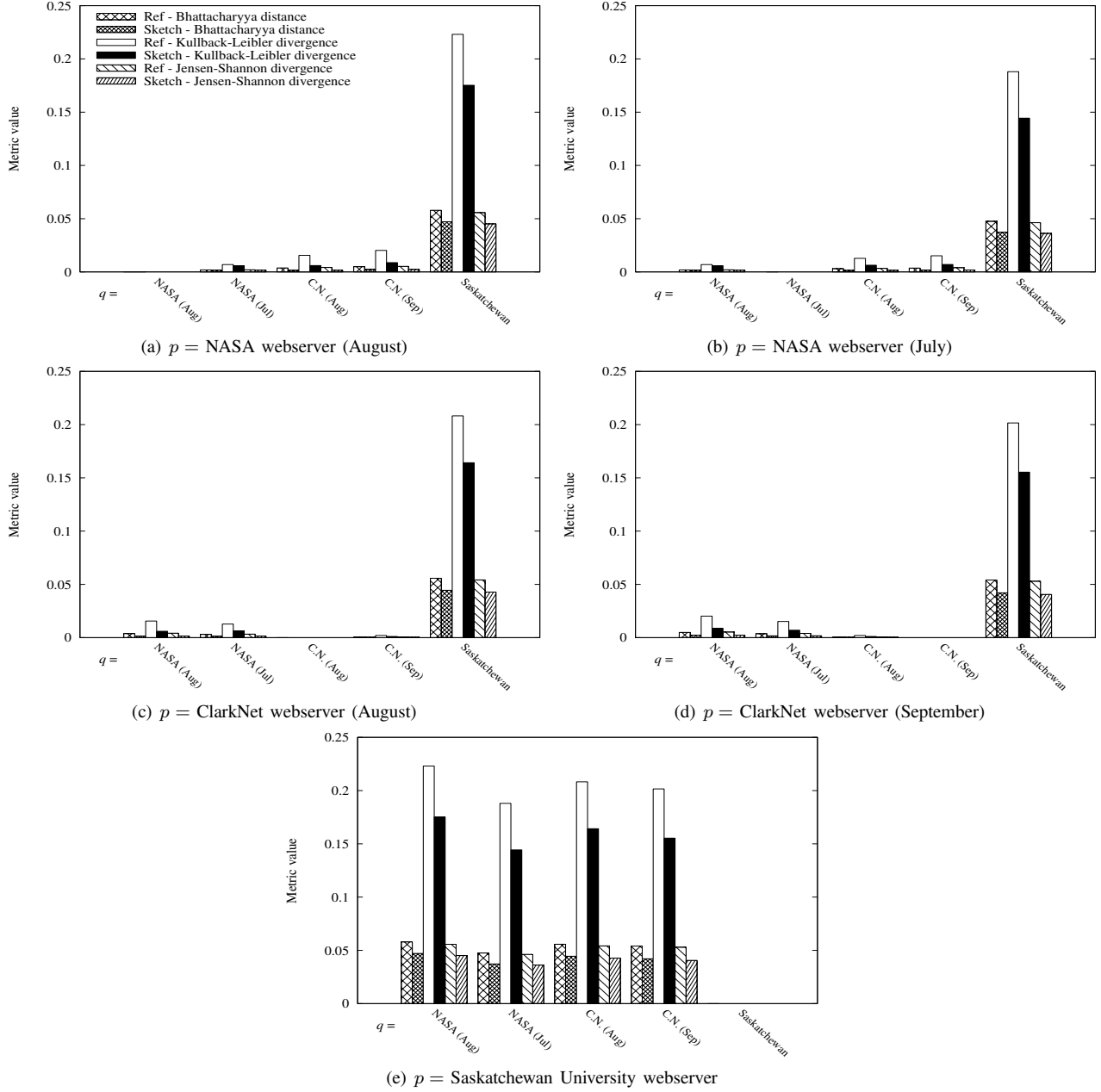
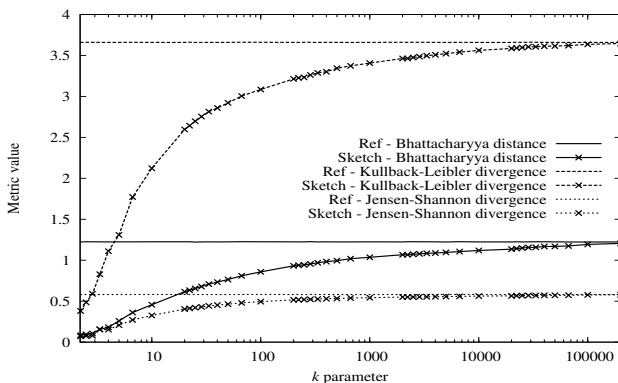
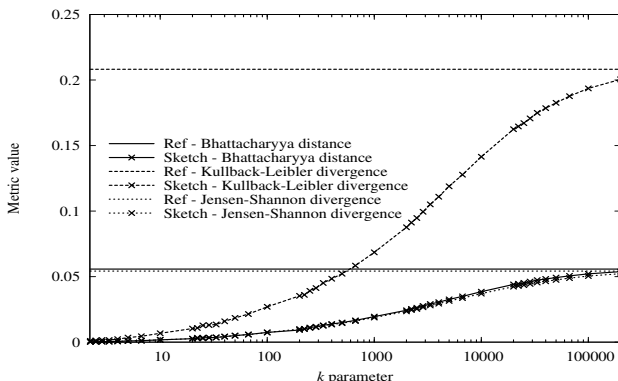


Figure 4. Sketch  $\star$ -metric accuracy as a function of real data traces. Parameters setting is as follows:  $k = 2,000$ ;  $t = 4$ .



(a) Sketch  $\star$ -metric accuracy as a function of  $k$ . We have  $m = 200,000$ ;  $n = 4,000$ ;  $t = 4$ ;  $r = 3$



(b) Sketch  $\star$ -metric accuracy between data trace extracted from ClarkNetwork (August) and Saskatchewan University, as a function of  $k$

Figure 5. Sketch  $\star$ -metric between the Uniform distribution and Pascal with parameter  $p = \frac{n}{2r+n}$  (Figures 5(a) and 6(a)), and between data trace extracted from ClarkNetwork (August) and Saskatchewan University (Figures 5(b) and 6(b)).

Figure 6 shows the slight influence of the number  $t$  of generated partitions on the accuracy of our metric. The reason comes from the use of 2-universal hash functions, which guarantee for each of them and with high probability that data items are uniformly distributed over the cells of any partition. As a consequence, augmenting the number of such hash functions has a weak influence on the accuracy of the metric. Figure 7 focuses on the error made on the Sketch  $\star$ -metric for five different values of  $t$  as a function of parameter  $r$  of the Pascal distribution (recall that increasing values of  $r$  – while maintaining the mean value – makes the shape of the Pascal distribution flatter). Figures 7(b), 7(d), and 7(f) respectively depict for each value of  $t$  the difference between the reference and the sketch values which makes more visible the impact of  $t$ . The same main lesson drawn from these figures is the moderate impact of  $t$  on the precision of our algorithm.

## VIII. CONCLUSION AND OPEN ISSUES

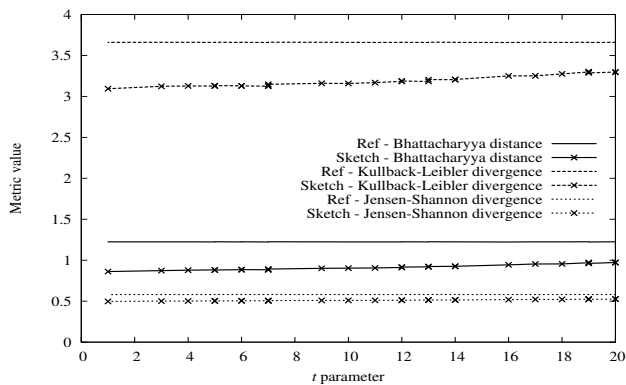
In this paper, we have introduced a new metric, the Sketch  $\star$ -metric, that allows to compute any generalized metric  $\phi$  on the summaries of two large input streams. We have presented a simple and efficient algorithm to sketch streams and compute

this metric, and we have shown that it behaves pretty well whatever the considered input streams. We are convinced of the undisputable interest of such a metric in various domains including machine learning, data mining, databases, information retrieval and network monitoring.

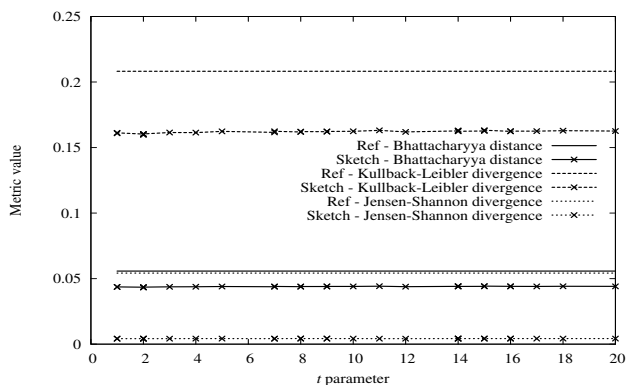
Regarding future works, we plan to characterize our metric among Rényi divergences [31], also known as  $\alpha$ -divergences, which generalize different divergence classes. We also plan to consider a distributed setting, where each site would be in charge of analyzing its own streams and then would propagate its results to the other sites of the system for comparison or merging. An immediate application of such a “tool” would be to detect massive attacks in a decentralized manner (e.g., by identifying specific connection profiles as with worms propagation, and massive port scan attacks or by detecting sudden variations in the volume of received data).

## REFERENCES

- [1] B. K. Subhabrata, E. Krishnamurthy, S. Sen, Y. Zhang, and Y. Chen, “Sketch-based change detection: Methods, evaluation, and applications,” in *Internet Measurement Conference*, 2003, pp. 234–247.
- [2] Y. Busnel, M. Bertier, and A.-M. Kermarec, “SOLIST or How To Look For a Needle in a Haystack?” in the *4th IEEE International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob’2008)*, Avignon, France, October 2008.
- [3] M. Chu, H. Haussecker, and F. Zhao, “Scalable information-driven sensor querying and routing for ad hoc heterogeneous sensor networks,” *International Journal of High Performance Computing Applications*, vol. 16, no. 3, pp. 293–313, 2002.
- [4] S. Kullback and R. A. Leibler, “On information and sufficiency,” *The Annals of Mathematical Statistics*, vol. 22, no. 1, pp. 79–86, 1951. [Online]. Available: <http://dx.doi.org/10.2307/2236703>
- [5] I. Csiszár, “Eine informationstheoretische ungleichung und ihre anwendung auf den beweis der ergodizitat von markoffschen ketten,” *Magyar. Tud. Akad. Mat. Kutató Int. Közl.*, vol. 8, pp. 85–108, 1963.
- [6] T. Morimoto, “Markov processes and the  $h$ -theorem,” *Journal of the Physical Society of Japan*, vol. 18, no. 3, pp. 328–331, 1963.
- [7] S. M. Ali and S. D. Silvey, “General Class of Coefficients of Divergence of One Distribution from Another,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 28, no. 1, pp. 131–142, 1966.
- [8] A. Bhattacharyya, “On a measure of divergence between two statistical populations defined by their probability distributions,” *Bulletin of the Calcutta Mathematical Society*, vol. 35, pp. 99–109, 1943.
- [9] M. Basseville and J.-F. Cardoso, “On entropies, divergences, and mean values,” in *Proceedings of the IEEE International Symposium on Information Theory*, 1995.
- [10] N. Alon, Y. Matias, and M. Szegedy, “The space complexity of approximating the frequency moments,” in *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing (STOC)*, 1996, pp. 20–29.
- [11] T. Cover and J. Thomas, “Elements of information theory,” Wiley New York, 1991.
- [12] A. Chakrabarti, G. Cormode, and A. McGregor, “A near-optimal algorithm for computing the entropy of a stream,” in *In ACM-SIAM Symposium on Discrete Algorithms*, 2007, pp. 328–335.
- [13] S. Guha, A. McGregor, and S. Venkatasubramanian, “Streaming and sublinear approximation of entropy and information distances,” in *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2006, pp. 733–742.
- [14] A. Chakrabarti, K. D. Ba, and S. Muthukrishnan, “Estimating entropy and entropy norm on data streams,” in *In Proceedings of the 23rd International Symposium on Theoretical Aspects of Computer Science (STACS)*. Springer, 2006.
- [15] A. Lall, V. Sekar, M. Ogihara, J. Xu, and H. Zhang, “Data streaming algorithms for estimating entropy of network traffic,” in *Proceedings of the joint international conference on Measurement and modeling of computer systems (SIGMETRICS)*. ACM, 2006.



(a) Sketch  $\star$ -metric accuracy as a function of parameter  $t$ . We have  $m = 200,000$ ;  $n = 4,000$ ;  $k = 200$  and  $r = 3$

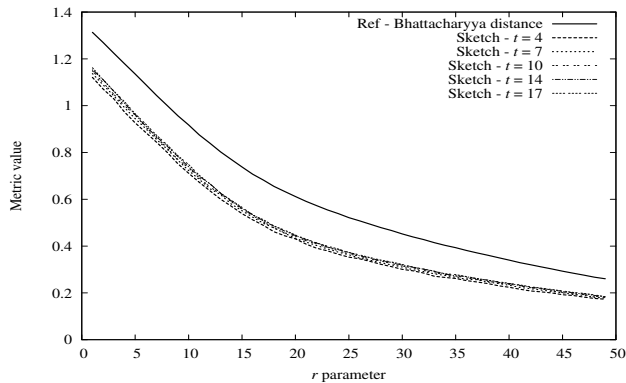


(b) Sketch  $\star$ -metric accuracy between data trace extracted from ClarkNetwork (August) and Saskatchewan University, as a function of  $t$

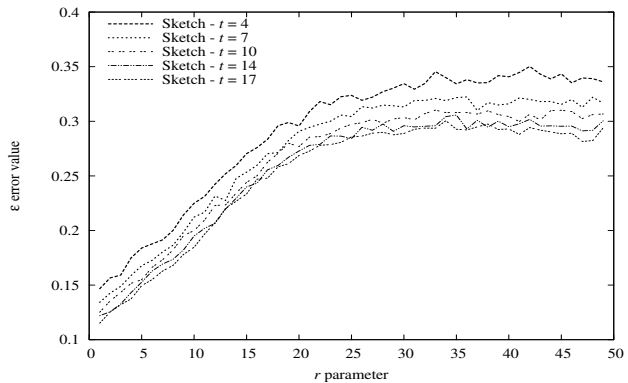
Figure 6. Sketch  $\star$ -metric between the Uniform distribution and Pascal with parameter  $p = \frac{n}{2r+n}$  (Figures 5(a) and 6(a)), and between data trace extracted from ClarkNetwork (August) and Saskatchewan University (Figures 5(b) and 6(b)).

[16] E. Anceaume, Y. Busnel, and S. Gamba, “AnKLe: detecting attacks in large scale systems via information divergence,” in *Proceedings of the*

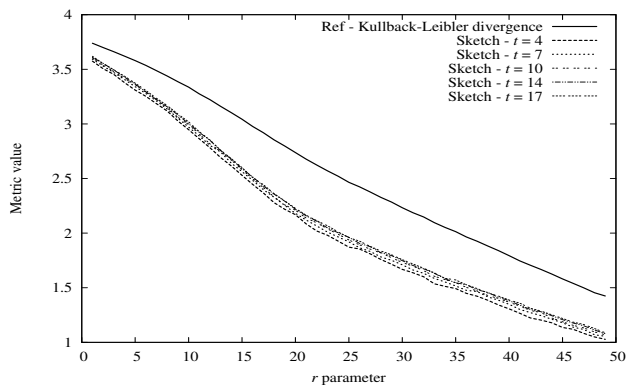
- 9th European Dependable Computing Conference (EDCC)*, 2012.
- [17] E. Anceaume and Y. Busnel, “An information divergence estimation over data streams,” in *Proceedings of the 11th IEEE International Symposium on Network Computing and Applications (NCA)*, 2012.
- [18] M. Charikar, K. Chen, and M. Farach-Colton, “Finding frequent items in data streams,” *Theoretical Computer Science*, vol. 312, no. 1, pp. 3–15, 2004.
- [19] G. Cormode and S. Muthukrishnan, “An improved data stream summary: the count-min sketch and its applications,” *J. Algorithms*, vol. 55, no. 1, pp. 58–75, 2005.
- [20] G. Cormode and M. Garofalakis, “Sketching probabilistic data streams,” in *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*, 2007, pp. 281–292.
- [21] S. Guha, P. Indyk, and A. McGregor, “Sketching information divergences,” *Machine Learning*, vol. 72, no. 1-2, pp. 5–19, 2008.
- [22] Muthukrishnan, *Data Streams: Algorithms and Applications*. Now Publishers Inc., 2005.
- [23] L. M. Bregman, “The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming,” *USSR Computational Mathematics and Mathematical Physics*, vol. 7, no. 3, pp. 200–217, 1967.
- [24] E. Hellinger, “Neue begründung der theorie quadratischer formen von unendlichvielen veränderlichen,” *J. Reine Angew. Math.*, vol. 136, pp. 210–271, 1909.
- [25] I. Csiszár, “Information Measures: A Critical Survey,” in *Transactions of the Seventh Prague Conference on Information Theory, Statistical Decision Functions, Random Processes*. Dordrecht: D. Riedel, 1978, pp. 73–86.
- [26] I. Csiszár, “Why least squares and maximum entropy? an axiomatic approach to inference for linear inverse problems,” *The Annals of Statistics*, vol. 19, no. 4, pp. 2032–2066, 1991.
- [27] S.-I. Amari and A. Cichocki, “Information geometry of divergence functions,” *Bulletin of the Polish Academy of Sciences: Technical Sciences*, vol. 58, no. 1, pp. 183–195, 2010.
- [28] S.-I. Amari, “ $\alpha$ -Divergence Is Unique, Belonging to Both  $f$ -Divergence and Bregman Divergence Classes,” *IEEE Transactions on Information Theory*, vol. 55, no. 11, pp. 4925–4931, nov 2009.
- [29] the Internet Traffic Archive, “<http://ita.ee.lbl.gov/html/traces.html>,” Lawrence Berkeley National Laboratory, Apr. 2008.
- [30] M. F. Arlitt and C. L. Williamson, “Web server workload characterization: the search for invariants,” *SIGMETRICS Performance Evaluation Review*, vol. 24, no. 1, pp. 126–137, 1996.
- [31] A. Renyi, “On measures of information and entropy,” in *Proceedings of the 4th Berkeley Symposium on Mathematics, Statistics and Probability*, 1960, pp. 547–561.



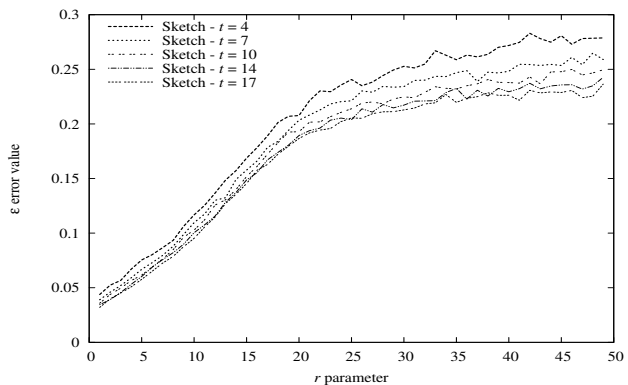
(a) Value of Bhattacharyya distance



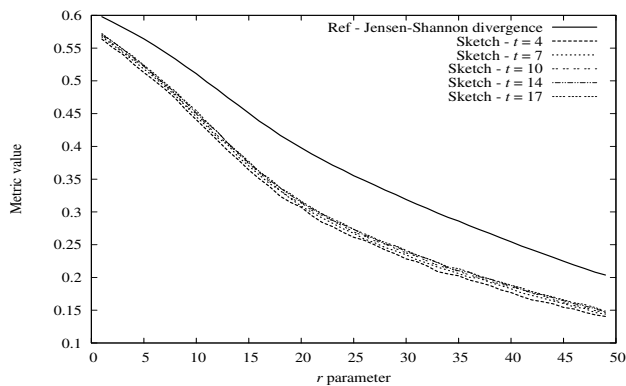
(b) Difference with Bhattacharyya distance



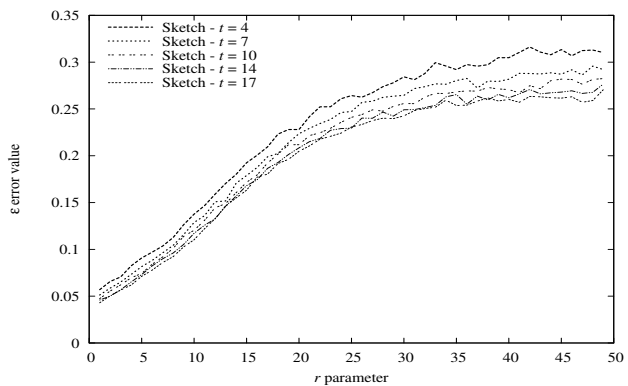
(c) Value of Kullback-Leibler divergence



(d) Difference with Kullback-Leibler divergence



(e) Value of Jensen-Shannon divergence



(f) Difference with Jensen-Shannon divergence

Figure 7.  $Sketch \star$ -metric estimation between Uniform distribution and Pascal with parameter  $p = \frac{n}{2r+n}$ , as a function of  $k$ ,  $t$  and  $r$ .