

# Alpha-Divergence for Classification, Indexing and Retrieval<sup>0</sup>

Alfred O. Hero, Bing Ma, Olivier Michel, and John Gorman

Communications and Signal Processing Laboratory  
Technical Report CSPL-328

May 2001

<http://www.eecs.umich.edu/~hero>, [hero@eecs.umich.edu](mailto:hero@eecs.umich.edu)

## Abstract

Motivated by Chernoff's bound on asymptotic probability of error we propose the alpha-divergence measure and a surrogate, the alpha-Jensen difference, for feature classification, indexing and retrieval in image and other databases. The alpha-divergence, also known as Renyi divergence, is a generalization of the Kullback-Liebler divergence and the Hellinger affinity between the probability density characterizing image features of the query and the density characterizing features of candidates in the database. As in any divergence-based classification problem, the alpha-divergence must be estimated from the query or reference object and the objects in the database. The surrogate for the alpha-divergence, called the alpha-Jensen difference, can be simply estimated using non-parametric estimation of the joint alpha-entropy of the merged pairs of feature vectors. Two methods of alpha-entropy estimation are investigated: (1) indirect methods based on parametric or non-parametric density estimation over feature space; and (2) direct methods based on combinatorial optimization of minimal spanning trees or other continuous quasi-additive graphs over feature space. On the basis of mean square error convergence rate comparisons the minimal graph entropy estimator can have better performance than an indirect entropy estimator implemented with plug-in density estimates. We illustrate these results for estimation of dependency in the plane and geo-registration of images.

## 1 Introduction

A database of images  $\mathcal{X} = \{X_i\}_{i=1}^K$  is queried for content which is closely related to a reference image  $X_0$ . The answer to the query is a partial re-indexing of the database in decreasing order of similarity to the reference image using an index function. This content-based retrieval problem arises in geographical information systems, digital libraries, medical information processing, video indexing, multi-sensor fusion, and multimedia information retrieval [38, 42, 41, 40]. Common methods for image indexing and retrieval are color histogram matching and texture matching using cross correlation. While these methods are computationally simple they often lack accuracy and discriminatory power.

There are three key ingredients to image retrieval and indexing which impact the accuracy and computation efficiency:

---

<sup>0</sup>Alfred Hero is with the Departments of Electrical Engineering and Computer Science (EECS), Biomedical Engineering, and Statistics at the University of Michigan Ann Arbor, MI 48109-2122. Bing Ma was with the Dept. of EECS at UM and is now with Intervideo, Inc., Fremont, CA. Olivier Michel is with the Department of Astrophysics, University of Nice, France. John Gorman is with ERIM-Veridian, Inc., Ann Arbor, MI.

1. selection of image features which discriminate between different image classes yet possess invariances to unimportant attributes of the images, e.g. rigid translation, rotation and scale;
2. application of an index function that measures feature similarity and is capable of resolving important differences between images;
3. query processing and search optimization which allow fast implementation.

While these ingredients are all closely linked, this paper is primarily concerned with the appropriate choice of the feature similarity measure and its optimization. We consider the class of  $\alpha$ -divergences, also known as Rényi divergences, and a surrogate function called the  $\alpha$ -Jensen difference. The  $\alpha$ -divergences can be roughly viewed as distances between the probability models underlying the query and the database of images. A special case of  $\alpha$ -divergence is the Kullback-Liebler (KL) divergence which has been applied to indexing and image retrieval by Stoica, Zerubia and Francos [40] and Do and Vetterli [10]. A related quantity is the  $\alpha$ -Jensen difference which is a function of the joint  $\alpha$ -entropy of pairs of feature vectors derived from the query and images in the database. The  $\alpha$ -Jensen difference was proposed independently by Ma [31, 30] for registering pairs of image modalities and by He, Ben-Hamza and Krim [19] for registering an arbitrary number of image modalities. Another special case of the  $\alpha$  divergence is the  $\alpha$ -information which is a generalization of the Shannon mutual information. Although we do not explore this extension here, the  $\alpha$ -information can be further generalized to the “ $f$ -information” which has been treated in a recent paper [35] for medical image registration and generalizes the mutual-information method of Viola and Wells [45]

Here we motivate the  $\alpha$ -divergence for indexing by decision theoretic considerations and large deviation theory of detection and classification. A result of this paper is that use of the KL divergence ( $\alpha = 1$ ) can be suboptimal relative to the more general  $\alpha$ -divergence. In particular, we establish that when the feature densities are difficult to discriminate (close together in a weighted sup-norm metric) the theoretically optimal choice of  $\alpha$  is  $\alpha = 1/2$  which corresponds to the Hellinger affinity, related monotonically to the Hellinger-Battacharya distance, as contrasted to the KL divergence. We compare the local discrimination capabilities of the  $\alpha$ -divergence and the  $\alpha$ -Jensen difference. In particular we show that for discrimination between pairs of close feature densities the  $\alpha$ -divergence admits a value  $\alpha = 1/4$  which is universally optimal while for the  $\alpha$ -Jensen difference the optimal value of  $\alpha$  depends on the feature density pair.

When either the  $\alpha$ -divergence or the  $\alpha$ -Jensen difference are used to perform indexing, they must be estimated from the query and the database. In this paper we focus on estimation of the  $\alpha$ -Jensen difference. When a smooth parametric model for the feature densities exists this entropy metric is a smooth non-linear function of these parameters and parametric estimation techniques such as maximum likelihood can be applied [41, 40]. For the parametric case, the entropy estimation error decreases at rate  $1/\sqrt{n}$  where  $n$  is the size of the training sample of feature vectors used for density estimation. When there exists no suitable parametric model for the feature densities non-parametric density estimation methods can be applied to estimate the  $\alpha$ -divergence or the  $\alpha$ -Jensen difference. This technique is called a “density plug-in” method. On the other hand, minimal graph entropy estimation techniques, introduced in Hero and Michel [23], can be applied to directly estimate the  $\alpha$ -Jensen difference. Using recent random graph convergence rates derived by Hero and Ma [20] for densities of bounded variation, we show that the convergence rate of non-parametric plug-in methods based on density plug-in estimation are slower than the rate of direct minimal-graph estimation methods implemented with the minimal spanning tree (MST) or other “continuous quasi-additive” graphs.

Finally, we show how minimal graphs can be applied to estimation of  $\alpha$ -divergence, when a consistent estimator of the reference density is available; estimation of the  $\alpha$ -mutual information and  $\alpha$ -Jensen difference. For purposes of illustration, we apply these results to a geo-registration problem.

## 2 Statistical Framework

Let  $X_0$  be a reference image, called the query, and consider a database  $X_i$ ,  $i = 1, \dots, K$  of images to be indexed relative to the query. Let  $Z_i$  be a feature vectors extracted from  $X_i$ . We assume that image  $X_i$ 's feature vector  $Z_i$  is realization  $Z$  generated by a j.p.d.f.  $f(Z|\underline{\theta})$  which depends on a vector of unknown parameters  $\underline{\theta}$  lying in a specified parameter space  $\Theta$ . Under this probabilistic model the  $k$ -th observed image feature vector  $Z_k$  is assumed to have been generated from model  $f(Z|\underline{\theta}_k)$ , where  $\underline{\theta}_k$  is called the ‘‘true parameter’’ underlying  $Z_k$ ,  $k = 1, \dots, K$ . Under this statistical framework the similarity between images  $X_0, X_i$  is reduced to similarity between feature probability models  $f(Z|\underline{\theta}_0), f(Z|\underline{\theta}_i)$ .

### 2.1 Divergence Measures of Dissimilarity

Define the densities  $f_i = f(Z|\theta_i)$ ,  $i = 0, \dots, K$ . The  $\alpha$ -divergence between  $f_i$  and  $f_0$  of fractional order  $\alpha \in [0, 1]$  is defined as [36, 7, 2]

$$\begin{aligned} D_\alpha(f_i||f_0) &= \frac{1}{\alpha-1} \ln \int f_0 \left( \frac{f_i(z)}{f_0(z)} \right)^\alpha dz \\ &= \frac{1}{\alpha-1} \ln \int f_i^\alpha(z) f_0^{1-\alpha}(z) dz \end{aligned} \quad (1)$$

Note that  $D_\alpha(f_i||f_0) = D_\alpha(\theta_i||\theta_0)$  is indexed by  $\theta_i$  and  $\theta_0$ .

*$\alpha$ -Divergence: Special cases*

When specialized to various values of  $\alpha$  the  $\alpha$ -divergence can be related to other well known divergence measures. Two of the most important examples are the Hellinger affinity  $2 \ln \int \sqrt{f_i(z)f_0(z)} dz$  obtained when  $\alpha = 1/2$ , and is related to the Hellinger-Battacharya distance squared,

$$\begin{aligned} D_{\text{Hellinger}}(f_i||f_0) &= \int \left( \sqrt{f_i(z)} - \sqrt{f_0(z)} \right)^2 dz \\ &= 2 \left( 1 - \exp \left( \frac{1}{2} D_{\frac{1}{2}}(f_i || f_0) \right) \right), \end{aligned}$$

and the Kullback-Liebler (KL) divergence [29], obtained when  $\alpha \rightarrow 1$ ,

$$\lim_{\alpha \rightarrow 1} D_\alpha(f_i, f_0) = \int f_0(z) \ln \frac{f_0(z)}{f_i(z)} dz.$$

Only when  $\alpha = 1/2$  is the divergence monotonically related to a true distance metric between two densities.

When the density  $f_0$  dominates  $f_1$  and is uniform over a compact domain  $\mathcal{Z} \supset \text{support}\{f_i\}$  the  $\alpha$ -divergence reduces to the  $\alpha$ -entropy, also known as the Rényi entropy:

$$H_\alpha(f_i) = \frac{1}{1-\alpha} \ln \int_{\mathcal{Z}} f_i^\alpha(z) dz \quad (2)$$

## 3 $\alpha$ -divergence as an Index Function

The ordered sequence of increasing  $\alpha$ -divergence measures  $D_\alpha(f_{(1)}||f_0), \dots, D_\alpha(f_{(K)}||f_0)$ , induces an indexing, which we call the ‘‘true indexing,’’ of the images

$$X_i \prec X_j \Leftrightarrow D_\alpha(f_i||f_0) < D_\alpha(f_j||f_0)$$

This indexing is unimplementable given only the  $Z_i$ 's since it requires the underlying probability models  $f_i$  be known to the query processor. The non-statistical indexing problem can now be stated as: given a sequence of divergences  $\{D_\alpha(\theta_i||\theta_0)\}_{i=1}^K$  find the sequence of indices  $i_1, \dots, i_K$  which minimize  $D_\alpha(\theta_{i_k}||\theta_0)$  over the set  $\{1, \dots, K\} - \{i_1, \dots, i_{k-1}\}$ ,  $k = 1, \dots, K$ .

Special cases of the indexing problem are

1. Content-based retrieval: the query is the density of an image object and the database consists of image densities which may “contain” the object in the sense that the object may only be found as a scaled, rotated or ortho-projected version of the query in the database. An invariant feature set is very important for this application.
2. Image registration: the database consists of  $K$  copies of  $Z_0$  which are rotated, translated and possibly locally deformed. The index  $i_1$  finds the pose/orientation in the database closest to that of the query. An invariant feature set is not desirable in this application. When the feature vector  $Z_i$  is defined as the set of pixel pair gray levels associated with each pair of images  $X_i, X_0$  and the mutual information criterion is applied to the pixel pair histogram one obtains the method of Viola and Wells [45]. The MI criterion is equivalent to the KL divergence between the joint distribution of the pixel-pair gray levels and the product of the marginal feature distributions.
3. Target detection: the query is the distribution of the observations and the database is partitioned into of a family of densities  $f_i = f(Z|\theta_i)$  part of which corresponds to the “target-absent” hypothesis and the rest to “target-present.” Target detection is declared if the closest density in the database is in the latter set.
4. Performing parameter estimation by minimizing the Hellinger-Battacharya distance is known as minimum-Hellinger-distance-estimation (MHDE) introduced by Beran [5]. While there are obvious similarities, relations of MHDE to indexing will not be explored in this paper.

### 3.1 Un-normalized $\alpha$ -Divergence and the Chernoff Error Exponent

Here we argue appropriateness of the  $\alpha$ -divergence on the basis of large deviations theory results on the exponential rate of decay of the Bayes-optimal classifier between two densities. Note that the Bayes classification error probability below is different from that defined by Vasconcelos [43, 42] in that here the decision error is averaged over an ensemble of image models. Define the un-normalized  $\alpha$ -divergence as the -log integral in the definition (2) of the  $\alpha$ -divergence:

$$D_\alpha^u(f_1||f_0) = -\ln \int f^\alpha(Z|\underline{\theta}_1) f^{1-\alpha}(Z|\underline{\theta}_0) dZ = (1-\alpha)D_\alpha(f_1||f_0)$$

Assume that from an i.i.d. sequence of images  $X^{(1)}, \dots, X^{(n)}$  we extract feature vectors  $\underline{Z} = [Z^{(1)}, \dots, Z^{(n)}]$  each having density  $f(Z|\underline{\theta})$  for some  $\underline{\theta} \in \Theta$ . Consider testing the hypotheses

$$\begin{aligned} H_0 & : \underline{\theta} \in \Theta_0 \\ H_1 & : \underline{\theta} \in \Theta_1 \end{aligned}$$

where  $\Theta_0$  and  $\Theta_1$  partition the parameter space  $\Theta$ . In the context of image retrieval the parameter range  $\Theta_1$  could cover the  $K$  densities of the images in the database while parameter range  $\Theta_0$  covers densities outside of the database. In this case testing  $H_0$  vs.  $H_1$  is tantamount to testing whether the query lies in the database ( $H_1$ ) or not ( $H_0$ ). If  $H_1$  is decided then sequential hypothesis testing could subsequently be performed to completely search the database for specific query matches by successive refinement of the parameter space  $\Theta_1$  over a depth  $\log_2(K)$  binary tree.

Let  $f(\underline{\theta})$  be a prior over  $\Theta$  and assume that  $P(H_1) = \int_{\Theta_1} f(\underline{\theta})d\underline{\theta}$  and  $P(H_0) = 1 - P(H_1)$  are both positive. Then for any test of  $H_0$  vs.  $H_1$  define the average probability of error

$$P_e(n) = \beta(n)P(H_1) + \alpha(n)P(H_0)$$

where  $\beta(n)$  and  $\alpha(n)$  are Type II and Type I errors of the test, respectively, which depend on  $\underline{\theta}$  in general. The  $\alpha$ -divergence measure can be related to the minimum attainable probability of error through the Chernoff bound of large deviations theory [8]:

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \ln P_e(n) = - \sup_{\alpha \in [0,1]} D_\alpha^u(\bar{f}_1 \| \bar{f}_0), \quad (3)$$

where  $\bar{f}_1(Z) = \int_{\Theta_1} f(Z|\underline{\theta})f(\underline{\theta})d\underline{\theta}$  and  $\bar{f}_0(Z) = \int_{\Theta_0} f(Z|\underline{\theta})f(\underline{\theta})d\underline{\theta}$ . The quantity  $\sup_{\alpha \in [0,1]} D_\alpha^u(\bar{f}_1 \| \bar{f}_0)$  in (3) is called the *Chernoff exponent* which gives the asymptotically optimal rate of exponential decay of the error probability for testing  $H_0$  vs  $H_1$ . The optimal  $\alpha = \alpha_o$  which attains the maximum in (3) is obtained by finding the value of  $\alpha$  which maximizes  $D_\alpha^u(\bar{f}_1 \| \bar{f}_0)$ .

$$\alpha_o = \operatorname{argmax}_{\alpha \in [0,1]} \int \bar{f}_1^\alpha(Z) \bar{f}_0^{1-\alpha}(Z) dZ \quad (4)$$

### 3.2 Selection of $\alpha$

We have empirically determined that for an image indexing problem arising in georegistration (see Section 5) the value of  $\alpha$  leading to highest resolution seems to cluster around either 1 or 1/2 corresponding to the KL divergence and the Hellinger affinity respectively [30]. The determining factor appears to be the degree of differentiation between the densities  $\{f_i\}_{i=0}^K$ . If the densities are very similar, i.e. difficult to discriminate, then the indexing performance of the Hellinger affinity distance ( $\alpha = 1/2$ ) was observed to be better than the KL divergence ( $\alpha = 1$ ). This is consistent with the asymptotic local analysis below.

A locally optimum  $\alpha$  can be explored by asymptotic analysis of the Chernoff exponent. Specifically, the following is a direct result of Proposition 7 in Appendix A.

#### Proposition 1

$$D_\alpha^u(f_0 \| f_1) = \frac{\alpha(1-\alpha)}{2} \int \frac{(f_0(x) - f_1(x))^2}{f_0(x) + f_1(x)} dx + o(\Delta^2), \quad (5)$$

where  $\Delta \in [0, 1]$  is

$$\Delta = 2 \sup_x \frac{|f_1(x) - f_0(x)|}{f_1(x) + f_0(x)}.$$

Recall that the detection error probability decreases exponentially with Chernoff exponent  $\sup_{\alpha \in [0,1]} D_\alpha^u(f_1 \| f_0)$ . A consequence of (5) is that to order  $\Delta^2$  the optimum value of  $\alpha$  in the Chernoff exponent is 1/4.

As an illustrative example consider the case where  $f_0$  and  $f_1$  are multivariate Gaussian densities. The KL information for such a Gaussian feature model was adopted in [41, 40] for performing image indexing. Let  $f(x; \mu, \Lambda)$  be a real  $d$ -dimensional normal density with mean vector  $\mu$  and non-singular covariance matrix  $\Lambda$ . The un-normalized  $\alpha$ -divergence  $D_\alpha^u(f_1 \| f_0) = D_\alpha^u(f(x; \mu_1, \Lambda_1) \| f(x; \mu_0, \Lambda_0))$  of order  $\alpha$  is given by (see Proposition 8 in Appendix B).

$$D_\alpha^u(f(x; \mu_1, \Lambda_1) \| f(x; \mu_0, \Lambda_0)) = \underbrace{-\frac{1}{2} \ln \frac{|\Lambda_0|^\alpha |\Lambda_1|^{1-\alpha}}{|\alpha \Lambda_0 + (1-\alpha) \Lambda_1|}}_{\text{Term A}} + \underbrace{\frac{\alpha(1-\alpha)}{2} \Delta \mu^T (\alpha \Lambda_0 + (1-\alpha) \Lambda_1)^{-1} \Delta \mu}_{\text{Term B}} \quad (6)$$

where  $\Delta \mu = \mu_1 - \mu_0$ .

The divergence consists of two terms  $A$  and  $B$ .  $A$  is equal to zero when  $\Lambda_0 = \Lambda_1$  and  $B$  is equal to zero when  $\mu_0 = \mu_1$ . Term  $A$  is the log of the ratio of the determinants of the geometric mean and the arithmetic means of  $\Lambda_1$  and

$\Lambda_0$  with mean weights  $\alpha$  and  $1 - \alpha$ . Term  $B$  is the quadratic difference of mean vectors normalized by the arithmetic mean of  $\Lambda_1$  and  $\Lambda_0$  with mean weights  $\alpha$  and  $1 - \alpha$ .

An asymptotic expansion yields the following expression for the case that  $\Delta\mu = 0$ , i.e. equal means,

$$D_\alpha^u(f_1||f_0) = \frac{\alpha(1-\alpha)}{4} \text{tr}(\Lambda_1 - \Lambda_0)^2 + o(\text{tr}(\Lambda_1 - \Lambda_0)^2).$$

so that locally the Chernoff exponent increases in the trace norm of the differences between the feature covariances and, as expected,  $\alpha = 1/4$  is optimal.

## 4 Divergence and Entropy Estimation

In practice the image model parameters  $\underline{\theta}_k$ 's are unknown so that the actual relative ordering of  $\alpha$ -divergences  $\{D_\alpha(\underline{\theta}_k||\underline{\theta}_0)\}_{k=1}^K$  is also unknown. The statistical problem of indexing can be stated as follows: based on a single realization  $X_k = X_k^{(1)}$  of the  $k$ -th image,  $k = 0, \dots, K$ , estimate the actual rank ordering of  $\alpha$ -divergences  $\{D_\alpha(\underline{\theta}_k||\underline{\theta}_0)\}_{k=1}^K$  between feature distributions. Divergence estimation is closely related to entropy estimation which has a long history in the statistics and information theory communities.

Estimation of entropy is an important problem that arises in statistical pattern recognition, adaptive vector quantization, image registration and indexing, and other areas. Non-parametric estimation of Shannon entropy has been of interest to many in non-parametric statistics, pattern recognition, model identification, image registration and other areas [18, 27, 1, 44, 4, 45, 11]. Estimation of  $\alpha$ -entropy arises as a step towards Shannon entropy estimation, e.g., Mokkadem [32] constructed a non-parametric estimate of the Shannon entropy from a convergent sequence of  $\alpha$ -entropy estimates. However, estimation of the  $\alpha$ -entropy is of interest in its own right. The problem arises in vector quantization where Rényi entropy is related to asymptotic quantizer distortion via the Panter-Dite factor and Bennett's integral [15, 34]. The  $\alpha$ -entropy parametrizes the Chernoff exponent governing the minimum probability of error in binary detection problems [26, 6]. It also has been used for image registration from multiple modalities via the  $\alpha$ -Jensen difference [31, 30, 19]. The most natural estimation method is to substitute a non-parametric density estimator  $\hat{f}$  into the expression for entropy. This method has been widely applied to estimation of the Shannon entropy and is called "plug-in" estimation in [4]. Other methods of Shannon entropy estimation discussed in [4] include sample spacing estimators, restricted to  $d = 1$ , and estimates based on nearest neighbor distances.

Three general classes of methods can be identified: parametric estimators, non-parametric estimators based on density or function estimation, and non-parametric estimators based on direct estimation. The first two methods can be classified as *plug-in* techniques where a parametric or non-parametric density estimate  $\hat{f}$  or function estimate  $\hat{f}^\alpha$  are simply plugged into the divergence formula. When an accurate parametric model and good parameter estimates are available parametric plug-in estimates of divergence are attractive due to their  $1/\sqrt{n}$  convergence properties. An analytical parametric form of the divergence can often be derived over the parametric class of densities considered and maximum likelihood can be used to estimate parameters in the divergence formula. This approach was adopted under a multivariate Gaussian image model by Stoica *et al* [40] for image retrieval. For Gaussian  $f_1$  and  $f_0$  the KL divergence  $D_1(f_1||f_0)$  has a simple closed form expression, which can be derived as the limit of (37) as  $\alpha \rightarrow 1$ , and the authors in [40] proposed using maximum likelihood or least squares estimates of the mean and covariance parameters of each image.

Non-parametric plug-in divergence estimates do not benefit from closed form parametric expressions for divergence but avoid pitfalls of model dependent estimates. For example, when a non-parametric estimate of  $\hat{f}$  or of  $\hat{f}^\alpha$  are available the following plug-in estimates of  $\alpha$ -entropy are natural

$$H_\alpha(\hat{f}) = \frac{1}{1-\alpha} \ln \int \hat{f}^\alpha(z) dz \quad (7)$$

$$H_\alpha(\widehat{f^\alpha}) = \frac{1}{1-\alpha} \ln \int \widehat{f^\alpha}(z) dz. \quad (8)$$

For the special case of estimation of Shannon entropy  $\lim_{\alpha \rightarrow 1} H_\alpha(f) = -\int f(z) \ln f(z) dz$  recent non-parametric estimation proposals have included: histogram estimation plug-in [16]; kernel density estimation plug-in [1]; and sample-spacing density estimator plug-in [17]. The reader is referred to [4] for a comprehensive overview of work in non-parametric estimation of Shannon entropy. The main difficulties with non-parametric methods are due to the infinite dimension of the spaces in which the unconstrained densities lie. Specifically: density estimator performance is poor without stringent smoothness conditions; no unbiased density estimators generally exist; density estimators have high variance and are sensitive to outliers; the high dimensional integration in (7) might be difficult.

The problems with the above methods can be summarized by the basic observation: on the one hand parameterizing the divergence and entropy functionals with infinite dimensional density function models is a costly over-parameterization, while on the other hand artificially enforcing lower dimensional density parametrizations can produce significant bias in the estimates. This observation has motivated us to develop direct methods which accurately estimate the entropy without the need for performing artificial low dimensional parameterizations or non-parametric density estimation [21, 23, 22]. These methods are based on constructing minimal graphs spanning the feature vectors in the feature space. The overall length of these minimal graphs can be used to specify a strongly consistent estimator of entropy for Lebesgue continuous densities. In particular, let  $\mathcal{Z}^{(n)} = \{Z^{(1)}, \dots, Z^{(n)}\}$  and define the Euclidean functional of order  $\gamma$ :  $L_\gamma = L_\gamma(\mathcal{Z}^{(n)}) = \min_{e \in \mathcal{E}} \sum_e |e|^\gamma$  the overall length of a graph spanning  $n$  i.i.d. vectors  $Z^{(i)}$  in  $\mathbf{R}^d$  each with density  $f$ . Here  $\gamma \in (0, d)$  is real,  $e$  are edges in a graph connecting pairs of  $Z^{(i)}$ 's and the minimization is over some suitable subsets  $\mathcal{E}$  of the  $\binom{n}{2}$  edges of the complete graph. Examples include the minimal spanning tree (MST), Steiner tree (ST), minimal matching bipartite graph, and traveling salesman tour. The asymptotic behavior of  $L_\gamma$  over random points  $\mathcal{Z}^{(n)}$  as  $n \rightarrow \infty$  has been studied for over half a decade [3, 46, 39] and, based on these studies, in [23] we gave conditions under which

$$\hat{H}_\alpha(\mathcal{Z}^{(n)}) = \ln L_\gamma(\mathcal{Z}^{(n)})/n^\alpha - \ln \beta_{L_\gamma, d} \quad (9)$$

is an asymptotically unbiased and almost surely consistent estimator of the un-normalized  $\alpha$ -entropy of  $f$  where  $\alpha = (d - \gamma)/d$  and  $\beta_{L_\gamma, d}$  is a constant bias correction depending on the graph minimality criterion over  $\mathcal{E}$  but independent of  $f$ .

As shown in [23], optimal pruning of greedy implementations of the minimal graph can robustify the entropy estimator against outliers from contaminating distributions. This procedure consists of constructing the  $k$ -minimal graph which is defined as the minimum weight minimal graph spanning any  $k$  out of the  $n$  points in the realization of  $f$ . Divergence  $D_\alpha(f_1 \| f_0)$  between the observed feature density  $f$  and a reference feature density  $f_0$  can be estimated similarly via performing a preprocessing step before implementing the minimal-graph entropy estimator. In this preprocessing step one applies a measure transformation on the feature space which converts the reference density to a uniform density over the unit cube [22].

As contrasted with density-based estimates of entropy, minimal graph entropy estimators enjoy the following properties: they can have faster asymptotic convergence rates (see next sub-section), especially for non-smooth densities and for low dimensional feature spaces; they completely bypass the complication of choosing and fine tuning parameters such as histogram bin size, density kernel width, complexity, and adaptation speed; the  $\alpha$  parameter in the  $\alpha$ -entropy function is varied by varying the interpoint distance measure used to compute the weight of the minimal graph. On the other hand, the need for combinatorial optimization may be a bottleneck for a large number of feature samples, for example the MST or the  $k$ -NNG are ‘‘almost linear’’ algorithms of complexity  $O(n \log n)$ .

When  $f_0$  is known the  $\alpha$ -divergence can be estimated by minimal graph methods using the measure transformation method outlined [22]. For unknown  $f_0$  and unknown  $f_1$  the existence of consistent minimal-graph estimators of  $D_\alpha(f_1 \| f_0)$  is an open problem. The sequel of this paper will be concerned with an alternative index function, called the  $\alpha$ -Jensen difference, which is a function of the joint entropy of the query and candidate image feature sets. This function can be estimated using the entropy estimation techniques discussed above.

## 4.1 Robust Entropy Estimation: the $k$ -MST

In Hero and Michel [23] we established strong convergence results for a greedy approximation to the following minimal  $k$ -point Euclidean graph problem. Assume that one is given a set  $\mathcal{X}_n = \{x_1, \dots, x_n\}$  of  $n$  random points in the unit cube  $[0, 1]^d$  of dimension  $d \geq 2$  following a Lebesgue density  $f$ . Fix  $k$  and denote by  $\mathcal{X}_{n,k}$  a  $k$ -point subset of  $\mathcal{X}_n$ ,  $0 < k \leq n$ . The minimal  $k$ -point Euclidean graph problem is to find the subset of points  $\mathcal{X}_{n,k}$  and the set of edges connecting these points such that the resultant graph has minimum total weight  $L(\mathcal{X}_{n,k})$ . This problem arises in many combinatorial optimization problems, see references in [23] for a partial list. In addition to these problems, in [23] we noted that the weight function  $k$ -minimal graph could be useful for robust estimation of the Rényi entropy of order  $\nu$ , where  $\nu = (d - \gamma)/d \in (0, 1)$  is specified by the dimension  $d \geq 2$  and a weight exponent  $\gamma \in (0, d)$  on the Euclidean distance between vertices of the graph. The intuition is that the  $k$ -MST tends to only accept those points that are most clustered near each other, thereby the rejected points tend to be stragglers outside of the  $k$  cluster. In Hero and Michel [23] we established almost sure (a.s.) convergence of the normalized weight function of a greedy approximation to the a class of  $k$ -point minimal graphs including the  $k$ -MST. This normalized weight function converges to a limit which, for  $k \ll n$ , is a close approximation to the entropy integral in (9). The influence function was investigated in Hero and Michel [23] and quantitatively establishes the  $k$ -MST as a robust estimator of entropy.

The greedy approximation was introduced in [23] to reduce the exponential complexity of the exact  $k$ -MST algorithm yet retain its outlier resistant properties. This greedy approximation involves the following partitioning heuristic: dissect the support of the density  $f$ , assumed to be  $[0, 1]^d$ , into a set of  $m^d$  cells of equal volumes  $1/m^d$ ; rank the cells in increasing order of numbers of points contained; starting with the highest ranked cell and continuing down the list compute the minimal spanning graphs in each cell until at least  $k$  points are covered. Stitching together these small graphs gives a graph which is an approximation to the  $k$ -minimal graph and, when  $k = \alpha n$ , for which the log of the normalized weight function  $\hat{L}_{n,k}/k^\nu$  converges to an approximation of the Rényi entropy of order  $\nu$ . The computational advantage of the greedy algorithm comes from its divide-and-conquer multi-resolution structure: it only requires solving the difficult non-linear minimal graph construction on cells containing smaller numbers of points. When  $\alpha = 1$  this greedy approximation reduces to a partitioning approximation to the full minimal graph spanning all of the  $n$  points. By selecting the “progressive-resolution parameter”  $m$  as a function  $m(n)$  of  $n$  we obtain an adaptive multi-resolution approximation to the  $k$ -MST.

## 4.2 Entropy Estimator Convergence Comparisons

It can be shown that when  $f$  is a density supported on the unit cube  $[0, 1]^d$  the bias and variance of direct minimal-graph entropy estimators (9) and indirect density plug-in entropy estimators (7) converge to zero as a function of the number  $n$  of i.i.d. observations [23, 1]. Here we attack a harder problem: comparing the asymptotic convergence rates of the mean square error.

The estimators discussed below will be of the form  $\hat{H}_\alpha = \frac{1}{1-\alpha} \ln \hat{I}_\alpha$ , where  $\hat{I}_\alpha$  will be a consistent estimator of the integral

$$I_\alpha(f) = \int f^\alpha(z) dz.$$

Given estimates of  $f$  and  $f^\alpha$  define the density plug-in estimator  $I_\alpha(\hat{f})$  and function plug-in estimator  $I_\alpha(\hat{f}^\alpha)$  as

$$I_\alpha(\hat{f}) = \int \hat{f}^\alpha(x) dx \tag{10}$$

and

$$I_\alpha(\hat{f}^\alpha) = \int \hat{f}^\alpha(x) dx. \tag{11}$$



Define the direct graph-based estimator  $\hat{I}_\alpha$  as

$$\hat{I}_\alpha = L_\gamma(Z_1, \dots, Z_n) / (\beta_{L_\gamma, d} n^{(d-\gamma)/d}) \quad (12)$$

where  $(d - \gamma)/d = \alpha$ . A standard perturbation analysis of  $\ln(z)$  establishes that

$$|\hat{H}_\alpha - H_\alpha(f)| = \frac{1}{1 - \alpha} \frac{|\hat{I}_\alpha - I_\alpha(f)|}{I_\alpha(f)} + o(|\hat{I}_\alpha - I_\alpha(f)|).$$

Thus as a function of  $n$  the asymptotic rate of convergence of  $\hat{H}_\alpha - H_\alpha(f)$  will be identical to that of  $\hat{I}_\alpha - I_\alpha(f)$ . In the sequel we will therefore focus on the convergence of  $\hat{I}_\alpha$ .

### 4.3 Plug-In Estimators

Modern methods of non-parametric density estimation attempt to minimize the estimation error as the density varies over a function space. Common spaces of variation which are considered are the Hölder spaces  $\Sigma_d(\kappa, c)$ , Besov spaces  $B_{p,q}^\sigma$ , and the space of functions of bounded variation  $BV(c, d)$ . We will restrict the density  $f$  on  $[0, 1]^d$ , and associated functions, to one of these  $d$ -dimensional function spaces in this report.

The class  $\Sigma_d(\kappa, c)$  of order- $\kappa$  Hölder continuous functions over  $[0, 1]^d$  are defined as

$$\Sigma_d(\kappa, c) = \left\{ f(x) : |f(x) - p_x^{\lfloor \kappa \rfloor}(z)| \leq c \|x - z\|^\kappa \right\}$$

where  $p_x^k(z)$  is the Taylor polynomial (multinomial) of  $f$  of order  $k$  expanded about the point  $x$ ,  $\|x\|$  denotes the  $L_2$  norm and  $\lfloor \kappa \rfloor$  is defined as the greatest integer strictly less than  $\kappa$ . As  $\kappa$  becomes large the class  $\Sigma_d(\kappa, c)$  contains functions which are increasingly smooth. For example,  $\Sigma_d(1, c)$  is the space of Lipschitz functions and  $\Sigma_d(\infty, c)$  contains all infinitely differentiable functions.

The class  $B_{p,q}^\sigma(c)$  of Besov functions over  $[0, 1]^d$  is defined

$$B_{p,q}^\sigma(c) = \left\{ f : \|f\|_p + \|f - T(f)\|_{l_{p,q}^\sigma} \leq c \right\}$$

where  $\|f\|_p$  is the standard  $L_p$  norm and  $\|f - T(f)\|_{l_{p,q}^\sigma}$  is a norm of the approximation error of the best dyadic spline approximation to  $f$  of order  $r$ , where  $r$  is determined by  $\sigma, p$  and  $q$  (see Devore and Popov [9]).

The class  $BV(c, d)$  of functions over  $[0, 1]^d$  is defined as [37]

$$BV(c, d) = \left\{ f : \sup_{\{z_i\}} \sum_i |f(z_i) - f(z_{i-1})| \leq c \right\}, \quad (13)$$

where the maximum is taken over all countable subsets  $\{z_1, z_2, \dots\}$  of points in  $[0, 1]^d$ .

We have the following simple results which are given without proof.

**Proposition 2** *Assume that the Lebesgue density  $f$  is in a function space  $\mathcal{C}$  and that  $\hat{f}$  is a plug-in estimator with uniform rms convergence rate of order  $O(n^{-r})$  over  $\mathcal{C}$ . If  $\int f^{\alpha-1}(z) dz < \infty$  then*

$$E \left[ \left| I_\alpha(\hat{f}) - I_\alpha(f) \right|^2 \right]^{1/2} = O(n^{-r}). \quad (14)$$

When  $f^\alpha \in \mathcal{C} = \text{BV}(C, d)$  Proposition 2 can be tightened

**Proposition 3** Assume the Lebesgue density  $f$  is such that  $f^\alpha \in \text{BV}(C, d)$ . Let  $\widehat{f}^\alpha$  be a plug-in estimator with uniform rms convergence rate of order  $O(n^{-r})$  over  $\text{BV}$ .

$$E \left[ \left| I_\alpha(\widehat{f}^\alpha) - I_\alpha(f) \right|^2 \right]^{1/2} = O(n^{-r}). \quad (15)$$

#### 4.4 Minimal Graph Estimators

For a direct minimal-graph estimator (12) exact convergence rates have been elusive except in some special cases. Even for a uniform density  $f$ , exact rates are known only for  $d = 2$  [39, 46]. Only bounds on convergence rates are available for this case when  $d > 2$  and these bounds form the basis for proving the convergence results discussed below. For example it has been shown [46] that when  $f$  is a uniform density the normalized MST length functional  $L_\gamma/n^{(d-\gamma)/d}$  converges to the integral  $\beta_{L_\gamma, d} I_\alpha(f)$  with rate upper bounded by  $O(n^{-1/d})$ . This bound is tight for  $d = 2$ . In [20] we establish that the rms convergence rate of the normalized MST functional is upper bounded by  $n^{-1/(d+1)}$  for arbitrary density  $f$  such that  $f^\alpha$  is of bounded variation. This result holds under the assumptions that the minimal graph is continuous quasi-additive. This rate is better than the best possible rate attainable by an entropy plug-in estimator. Specifically, in [20] we establish the following

**Proposition 4** Assume that the Lebesgue density  $f$  on  $[0, 1]^d$  is such that  $f^\alpha \in \text{BV}(C, d)$  where  $\alpha \in [1/2, (d-1)/d]$ ,  $d \geq 2$ . Then, for  $p = 1, 2, \dots$ , and any plug-in estimator  $I_\alpha(\widehat{f}^\alpha)$

$$\sup_{f^\alpha \in \text{BV}(C, d)} E^{1/p} \left[ \left| I_\alpha(\widehat{f}^\alpha) - I_\alpha(f) \right|^p \right] \geq O(n^{-1/(d+2)}), \quad (16)$$

while for the MST-based entropy estimator  $\widehat{I}_\alpha$

$$\sup_{f^\alpha \in \text{BV}(C, d)} E^{1/p} \left[ \left| \widehat{I}_\alpha - I_\alpha(f) \right|^p \right] \leq O(n^{-1/(d+1)}). \quad (17)$$

##### 4.4.1 Achievability of Minimal Graph Estimator Bound

Some of the comments below are explored in more detail in Hero and Ma [20]. It is unknown whether the bound (17) is tight except in the case  $d = 2$ , for which the bound only holds for  $\alpha = 1/2$ , i.e. estimation of the Hellinger affinity. We point out that the general bound (17) also holds for other continuous quasi-additive graphs such as the  $k$  nearest neighbor graph. When  $\alpha < 1/2$  a potentially slower bound of order  $O(\max\{n^{-1/(d+1)}, n^{-\alpha/2}\})$  is available. When  $f^\alpha \in \text{BV}(C, d)$  we conclude from Proposition 4 that the worst case convergence rate  $n^{-1/(d+1)}$  of the minimal graph estimator is faster than the best rate  $n^{-1/(d+2)}$  of a plug-in estimator using a non-parametric function estimate of  $f^\alpha$ . In particular, this implies that the rate of convergence of the MST estimator of the Hellinger affinity ( $\alpha = 1/2$ ) is faster than any estimator based on minimax density estimation. However, the assumption  $0 < \alpha \leq (d-1)/d$  prevents the application of the rms convergence rate bound (17) to estimates of the Shannon entropy ( $\alpha \rightarrow 1$ ). In particular, we cannot use it to bound the rms of a minimal-graph based analog to the method of Makkadem [32] in which one estimates Shannon entropy by a sequence  $\widehat{I}_{\alpha_n}(\widehat{f}_n)$  of  $\alpha$ -entropy estimators where  $\alpha_n < 1$  and  $\lim_{n \rightarrow \infty} \alpha_n = 1$ .

If it is known *a priori* that  $f$  is piecewise constant with known regions of support a faster rate of convergence bound for the minimal graph estimator than (17) is available:  $O(\max\{n^{-1/d}\})$ . Thus for piecewise constant  $f$  the histogram plug in estimator has  $1/\sqrt{n}$  rms convergence rate and we conclude that the minimal graph and plug-in

estimators have identical convergence rate for  $d = 2$ . On the other hand, for  $d > 2$  the histogram plug-in estimator has faster rms convergence rate for piecewise constant  $f$ . Finally we mention a shortcoming of the minimal graph estimator: it is not consistent for estimating  $\alpha$ -entropy of any singular components of  $f$ , e.g. dirac delta functions.

Finally we point out that the extension of the rate of convergence bound (17) to the greedy  $k$ -MST algorithm is an open problem.

#### 4.4.2 Achievability of Plug-In Estimator Bound

The issue of achievability of the minimax bound (16) in Proposition 4 by a specific estimator is of course of interest but appears to be an open question. Thus the bound may be optimistic and the gap between convergence rates of plug-in and minimal graph estimators of entropy may be even larger than indicated in Proposition 4. Two classes of estimators have been introduced for non-parametric function estimation: linear Parzen-Rosenblatt kernel density estimators and non-linear wavelet shrinkage estimators.

Parzen-Rosenblatt (PR) kernel density estimators are defined as

$$\hat{f}(z) = \frac{1}{nh_n^d} \sum_{i=1}^n V\left(\frac{z - Z_i}{h_n}\right) \quad (18)$$

where  $V(z)$  is a kernel function satisfying  $\int V(z)dz = 1$ ,  $\int z^j V(z)dz < \infty$ ,  $j = 1, 2, \dots, [\kappa]$ , and  $h_n$  is a positive sequence satisfying  $h_n \rightarrow 0$  and  $nh_n \rightarrow \infty$ . For fixed and known  $\kappa$  the estimator (18) has minimax rate of convergence over the Hölder class  $\Sigma_d(\kappa, c)$  when the kernel-width  $h_n$  is chosen as:  $h_n = an^{-1/(2\kappa+d)}$  for some positive constant  $a$  [25, 28]. The PR estimator has root MISE (rms) asymptotic convergence rate which is constant over  $\Sigma_d(\kappa, c)$  and given by

$$\sup_{f \in \Sigma_d(\kappa, c)} E \left[ \int (\hat{f}(x) - f(x))^2 dx \right]^{1/2} = O\left(n^{-\kappa/(2\kappa+d)}\right) \quad (19)$$

for a positive constant  $C$ . Such linear estimators can achieve rate exponent  $r = 1/d + 2$  uniformly over  $\Sigma_d(1, c)$ . However, they cannot uniformly achieve this rate over the larger space  $BV(C, d)$ .

Minimax wavelet shrinkage methods of function estimation were introduced by Donoho and are discussed in detail in a discussion paper by Donoho *etal* [14]. Therein the authors show that wavelet shrinkage function estimators can be made to have rms convergence rates of order  $O\left(n^{(\ln n)^e \sigma / (2\sigma + d)}\right)$  uniformly over functions in  $B_{p,q}^\sigma$ , where  $e$  is a positive constant given by Donoho *etal* [12]. For  $\sigma > d/p$  this rate exponent is nearly equal to the minimax MISE rate exponent  $\alpha/(2\alpha + d)$  over the Besov space of functions  $B_{p,q}^\sigma$  [14, 13]. As pointed out by these authors,  $BV$  can be sandwiched between  $B_{1,1}^1$  and  $B_{1,\infty}^1$  so therefore if this “nearly minimax” result could be extended to the case  $\sigma = 1$  for some generalized class of multidimensional wavelet shrinkage estimators we would have nearly optimal rms convergence rate exponent  $(\ln n)^e / (d + 2)$  of the plug-in entropy estimate over the space of  $BV$  functions. As far as we know, existence of such a generalization is an open question.

## 4.5 $\alpha$ -Jensen Difference Index Function

Here we study an alternative index function based on the Jensen entropy difference. This index function was independently proposed by Ma [30] and He *etal* [19] for image registration problems. Let  $f_0$  and  $f_1$  be two densities and  $\beta \in [0, 1]$  be a mixture parameter. The  $\alpha$ -Jensen difference is the difference between the  $\alpha$ -entropies of the mixture  $f = \beta f_0 + (1 - \beta)f_1$  and the mixture of the  $\alpha$ -entropies of  $f_0$  and  $f_1$  [2]:

$$\Delta H_\alpha(\beta, f_0, f_1) \triangleq H_\alpha(\beta f_0 + (1 - \beta)f_1) - [\beta H_\alpha(f_0) + (1 - \beta)H_\alpha(f_1)], \quad \alpha \in (0, 1). \quad (20)$$

The  $\alpha$ -Jensen difference is measure of dissimilarity between  $f_0$  and  $f_1$ : as the  $\alpha$ -entropy  $H_\alpha(f)$  is concave in  $f$  it is clear from Jensen's inequality that  $\Delta H_\alpha(\beta, f_0, f_1) = 0$  iff  $f_0 = f_1$  a.e.

The  $\alpha$ -Jensen difference can be motivated as an index function as follows. Assume that two sets of labeled feature vectors  $\mathcal{Z}_0 = \{Z_0^{(i)}\}_{i=1, \dots, n_0}$  and  $\mathcal{Z}_1 = \{Z_1^{(i)}\}_{i=1, \dots, n_1}$  are extracted from images  $X_0$  and  $X_1$ , respectively. Assume that each of these sets consist of independent realizations from densities  $f_0$  and  $f_1$ , respectively. Define the union  $\mathcal{Z} = \mathcal{Z}_0 \cup \mathcal{Z}_1$  containing  $n = n_0 + n_1$  unlabeled feature vectors. Any consistent entropy estimator constructed on the unlabeled  $Z^{(i)}$ 's will converge to  $H_\alpha(\beta f_0 + (1 - \beta)f_1)$  as  $n \rightarrow \infty$  where  $\beta = \lim_{n \rightarrow \infty} n_0/n$ . This motivates the following consistent minimal-graph estimator of Jensen difference (20) for  $\beta = n_0/n$ :

$$\widehat{\Delta H}_\alpha(\beta, f_0, f_1) \triangleq \hat{H}_\alpha(\mathcal{Z}_0 \cup \mathcal{Z}_1) - \left[ \beta \hat{H}_\alpha(\mathcal{Z}_0) + (1 - \beta) \hat{H}_\alpha(\mathcal{Z}_1) \right], \quad \alpha \in (0, 1).$$

where  $\hat{H}_\alpha(\mathcal{Z}_0 \cup \mathcal{Z}_1)$  is the minimal-graph entropy estimator (9) constructed on the  $n$  point union of both sets of feature vectors and  $\hat{H}_\alpha(\mathcal{Z}_0)$ ,  $\hat{H}_\alpha(\mathcal{Z}_1)$  are constructed on the individual sets of  $n_0$  and  $n_1$  feature vectors, respectively. We can similarly define the density-based estimator of Jensen difference based on entropy estimates of the form (7) constructed on  $\mathcal{Z}_0 \cup \mathcal{Z}_1$ ,  $\mathcal{Z}_0$  and  $\mathcal{Z}_1$ .

For some indexing problems the marginal entropies  $\{H_\alpha(f_i)\}_{i=1}^K$  over the database are all identical so that the indexing function  $\{H_\alpha(\beta f_0 + (1 - \beta)f_i)\}_{i=1}^K$  is equivalent to  $\{\Delta H_\alpha(\beta, f_0, f_i)\}_{i=1}^K$ . This fact was used in Ma *et al* [31] for registering a query image to a database of images which are generated by entropy-preserving rigid transformations of a reference image.

## 4.6 Comparisons of $\alpha$ -Jensen Difference and $\alpha$ -Divergence

The local discrimination capabilities of the  $\alpha$ -Jensen difference and the  $\alpha$ -divergence can easily be compared using the results (Propositions 6 and 7) obtained in Appendix A:

$$\ln D_\alpha(f_0 \| f_1) = \ln \left( E_{f_{\frac{1}{2}}} \left[ \left( \frac{f_0 - f_1}{f_{\frac{1}{2}}} \right)^2 \right] \right) + C_1 + o(\Delta^2) \quad (21)$$

$$\begin{aligned} \ln \Delta H_\alpha(\beta, f_0, f_1) &= \ln \left( E_{\tilde{f}_{\frac{1}{2}}^\alpha} \left[ \left( \frac{f_0 - f_1}{f_{\frac{1}{2}}} \right)^2 \right] + \frac{\alpha}{1 - \alpha} \left( E_{\tilde{f}_{\frac{1}{2}}^\alpha} \left[ \frac{f_0 - f_1}{f_{\frac{1}{2}}} \right] \right)^2 \right) + C_2 \\ &\quad + o(\Delta^2) \end{aligned} \quad (22)$$

where  $E_f[g(x)] = \int f(x)g(x) dx$ ,  $\tilde{f}_{\frac{1}{2}}^\alpha(x) \triangleq \frac{f_{\frac{1}{2}}^\alpha(x)}{\int f_{\frac{1}{2}}^\alpha(x) dx}$  is a "tilted" pdf,  $\Delta$  is a term that decreases in the difference  $f_0 - f_1$ , and  $C_1, C_2$  are constants independent of  $f_0, f_1$ .

There are a number of interesting properties of  $\ln D_\alpha(f_0 \| f_1)$  and  $\ln \Delta H_\alpha(\beta, f_0, f_1)$ :

- Up to order  $\Delta^2$  the leading terms in (21) and (22) are the curvatures of the log  $\alpha$ -divergence and the log  $\alpha$ -Jensen difference, respectively. These curvatures are a measure of the sensitivity of these index functions for discriminating between density pairs  $f_0, f_1$ . The discrimination capability of the  $\alpha$ -divergence is locally independent of  $\alpha$  while that of the  $\alpha$ -Jensen difference depends on  $\alpha$ .
- When  $\alpha$  approaches 0, tail differences between the two densities  $f_0$  and  $f_1$  are influential on  $\Delta H_\alpha(\beta, f_0, f_1)$ .

- When  $\alpha$  approaches 1, central differences between the two densities become highly pronounced in  $\Delta H_\alpha(\beta, f_0, f_1)$ . Therefore, if the feature densities differ in regions where there is a lot of mass one should choose  $\alpha$  close to 1 to ensure locally optimum discrimination with  $\Delta H_\alpha(\beta, f_0, f_1)$ .
- $\Delta H_\alpha(\beta, f_0, f_1)$  has maximal discriminative capability when  $\beta = \frac{1}{2}$ , i.e., when two images yield the same number of feature vectors.

## 4.7 Estimation of $\alpha$ -Divergence

Here we describe extensions of the entropy estimation procedures described above to information divergence estimation. Let  $g(x)$  be a reference density on  $\mathbf{R}^d$  which dominates the density  $f(x)$  of a sample point  $x = [x^1, \dots, x^d]^T$  in the sense that for all  $x$  such that  $g(x) = 0$  we have  $f(x) = 0$ . The plug-in estimator of the  $\alpha$ -divergence based on independent estimation of  $g$  and  $f$  will have rms convergence rate  $O(n^{-1/(d+2)})$  when  $f^\alpha$  and  $g^{1-\alpha}$  are of bounded variation. As described below, when  $g(x)$  is known and  $f^\alpha$  is of bounded variation the minimal graph estimator can be applied, achieving faster rms convergence rate of at worst  $O(n^{-1/(d+1)})$ .

As introduced in Hero and Michel [22] minimal graph divergence estimation is performed by constructing a minimal graph on a transformed sample where the transformation corresponds to a change of measure which flattens the reference distribution  $g$ . For any  $x$  such that  $g(x) > 0$  let  $g(x)$  have the product representation  $g(x) = g(x^1)g(x^2|x^1) \dots g(x^d|x^{d-1}, \dots, x^1)$  where  $g(x^k|x^{k-1}, \dots, x^1)$  denotes the conditional density associated with  $g(x)$  of the  $k$ -th component. In what follows we will ignore the set  $\{x : g(x) = 0\}$  since, as  $f(x) = 0$  over this set, it has probability zero. Now consider generating the vector  $y = [y^1, \dots, y^d]^T \in \mathbf{R}^d$  by the following vector transformation

$$\begin{aligned} y^1 &= G(x^1) \\ y^2 &= G(x^2|x^1) \\ &\vdots \\ y^d &= G(x^d|x^{d-1}, \dots, x^1) \end{aligned} \quad (23)$$

where  $G(x^j|x^{j-1}, \dots, x^1) = \int_{-\infty}^{x^j} g(\tilde{x}^j|x^{j-1}, \dots, x^1)d\tilde{x}^j$  is the cumulative conditional distribution of the  $j$ -th component, which is monotone increasing except on the zero probability set  $\{x : g(x) = 0\}$ . Thus, except for this probability zero set, the conditional distribution has an inverse  $x^j = G^{-1}(y^j|x^{j-1}, \dots, x^1) = G^{-1}(y^j|y^{j-1}, \dots, y^1)$  and it can be shown (via the standard Jacobian formula for transformation of variables) that the induced joint density,  $h(y)$ , of the vector  $y$  takes the form:

$$h(y) = \frac{f(G^{-1}(y^1), \dots, G^{-1}(y^d|y^{d-1}, \dots, y^1))}{g(G^{-1}(y^1), \dots, G^{-1}(y^d|y^{d-1}, \dots, y^1))} \quad (24)$$

Now let  $L(\mathcal{Y}_n)$  denote the length of the greedy approximation to the MST constructed on the transformed random variables  $\mathcal{Y}_n = \{Y_1, \dots, Y_n\}$ . Then, by the consistency property of the MST estimator, we know that

$$\hat{H}_\alpha(\mathcal{Y}_n) \rightarrow \frac{1}{1-\alpha} \ln \int h^\alpha(y) dy \quad (a.s.) \quad (25)$$

and, from Section 4.4, the r.m.s. convergence rate will be  $O(n^{-1/(d+1)})$ . Making the inverse transformation  $y \rightarrow x$  specified by (23) in the above integral, noting that, by the Jacobian formula,  $dy = g(x)dx$ , and using the expression (24) for  $h$ , it is easy to see that the integral in the right hand side of (25) is equivalent to the Rényi information divergence of  $f(x)$  with respect to  $g(x)$

$$\frac{1}{1-\alpha} \ln \int h^\alpha(y) dy = \frac{1}{1-\alpha} \ln \int \left( \frac{f(x)}{g(x)} \right)^\alpha g(x) dx.$$

Hence we have established that  $\hat{H}_\alpha(\mathcal{Y}_n)$  is a strongly consistent estimator of the Rényi information divergence above. The results of Hero and Ma [20] and Hero and Michel [23] can thus be easily be extended to classification against any *arbitrary* distribution  $f_o$ , and not just the uniform distribution initially studied by Hoffman and Jain [24] and Hero and Michel [21]. This extension also holds for the  $k$ -MST described in Section 4.1.

#### 4.7.1 Application to robust divergence estimation

Here we applied the  $k$ -MST to robustly classify a triangular vs. uniform density on the plane. 256 samples were simulated from a triangle-uniform mixture density  $f = (1 - \epsilon)f_1 + \epsilon f_0$  where  $f_1(x) = (\frac{1}{2} - |x^1 - \frac{1}{2}|)(\frac{1}{2} - |x^2 - \frac{1}{2}|)$  is a (separable) triangular shaped product density and  $f_0 = 1$  is a uniform density, both supported on the unit square  $x = (x^1, x^2) \in [0, 1]^2$ . The Rényi information divergences  $I(f, f_0)$  and  $I(f, f_1)$  were estimated by  $\hat{H}_\alpha(\mathcal{X}_n)$  and  $\hat{H}_\alpha(\mathcal{Y}_n)$ , respectively, for  $\alpha = \frac{1}{2}$  ( $\gamma = 1$  in the  $k$ -MST construction).  $\mathcal{Y}_n$  was obtained by applying the mapping  $y = (y^1, y^2) = (F_1(x^1), F_1(x^2))$  to the data sample  $\mathcal{X}_n$ , where  $F_1(u)$  is the marginal cumulative distribution function associated with the triangular density.

In a first sequence of experiments the estimates  $\hat{H}_\alpha(\mathcal{X}_n)$  and  $\hat{H}_\alpha(\mathcal{Y}_n)$  of the respective quantities  $I(f, f_0)$  and  $I(f, f_1)$  were thresholded to decide between the hypotheses  $H_0 : \epsilon = 0$  vs.  $H_1 : \epsilon \neq 0$  and  $H_0 : \epsilon = 1$  vs.  $H_1 : \epsilon \neq 1$ , respectively. The receiver operating characteristic (ROC) curves are indicated in Figures 1 and 2. Note that, as expected, in each case the detection performance improves as the difference between the assumed  $H_0$  and  $H_1$  densities increases.

In a second sequence of experiments we selected two realizations of the triangle-uniform mixture model for the values  $\epsilon = 0.1$  and  $\epsilon = 0.9$ . For the former case the triangular is the dominating density and for the latter case the uniform is the dominating density. In each case the  $k$ -MST was implemented ( $k = 90$ ) as a robust clustering algorithm to identify data points from the dominating densities - in the former case the  $k$ -MST was applied directly to  $\mathcal{X}_n$  while in the latter case it was applied to  $\mathcal{Y}_n$ . The resulting  $k$ -MST quantities  $\hat{H}_\alpha(\mathcal{X}_{n,k})$  and  $\hat{H}_\alpha(\mathcal{Y}_{n,k})$  can be interpreted as robust estimates of the uncontaminated Rényi information divergences  $I(f_1, f_0)$  and  $I(f_0, f_1)$ , respectively. Figure 3-5 illustrate the effectiveness of these estimates as “outlier rejection” algorithms.

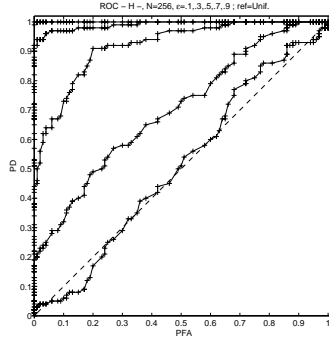


Figure 1: ROC curves for the Rényi information divergence test for detecting triangle-uniform mixture density  $f = (1 - \epsilon)f_1 + \epsilon f_0$  ( $H_1$ ) against the uniform hypothesis  $f = f_0$  ( $H_0$ ). Curves are decreasing in  $\epsilon$  over the range  $\epsilon \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$ .

## 4.8 Estimation of Dependency in the Plane

One indexing application is to rank order images according to the degree of feature dependence. For example, if two features  $X$  and  $Y$  are horizontal and vertical changes over local neighborhoods of pixels one can search for evidence

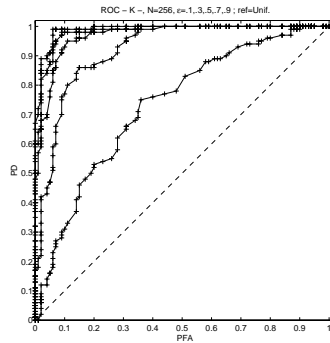


Figure 2: Same as Figure 1 except test is against triangular hypothesis  $f = f_1 (H_0)$ . Curves are increasing in  $\epsilon$ .

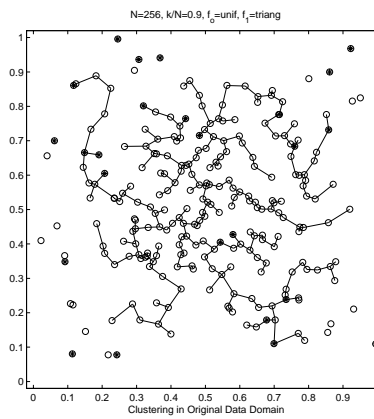


Figure 3: A scatterplot of a 256 point sample from triangle-uniform mixture density with  $\epsilon = 0.1$ . Labels 'o' and '\*' mark those realizations from the uniform and triangular densities, respectively. Superimposed is the  $k$ -MST implemented directly on the scatterplot  $\mathcal{X}_n$  with  $k = 230$ .

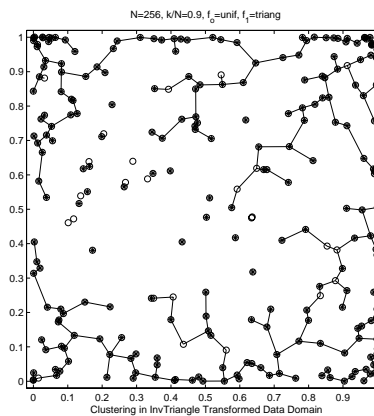


Figure 4: A scatterplot of a 256 point sample from triangle-uniform mixture density with  $\epsilon = 0.9$  in the transformed domain  $\mathcal{Y}_n$ . Labels 'o' and '\*' mark those realizations from the uniform and triangular densities, respectively. Superimposed is the  $k$ -MST implemented on the transformed scatterplot  $\mathcal{Y}_n$  with  $k = 230$

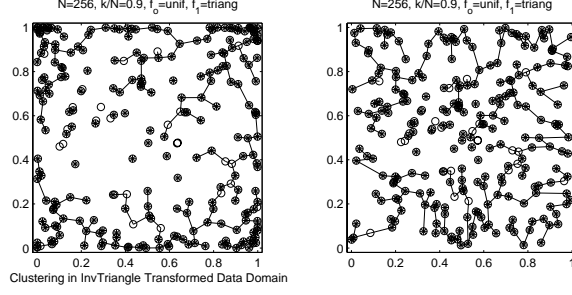


Figure 5: Same as Figure 4 except displayed in the original data domain.

of anisotropy by evaluating a measure of statistical dependence of  $X$  and  $Y$ . One possible measure is the mutual  $\alpha$ -information

$$MI_\alpha(X, Y) = \frac{1}{\alpha - 1} \ln \int f^\alpha(X, Y) (f(X)f(Y))^{1-\alpha} dX dY.$$

This quantity converges to the standard Shannon mutual information in the limit as  $\alpha \rightarrow 1$  and is equal to zero if and only if  $X, Y$  are independent.

Assume that the “two-dimensional function”  $f(X, Y)^\alpha$  and “one-dimensional functions”  $f^{1-\alpha}(X)$  and  $f^{1-\alpha}(Y)$  are all of bounded variation. If one could find minimax plug-in estimates based on independent estimates of  $f(X, Y)^\alpha$ ,  $f^{1-\alpha}(X)$  and  $f^{1-\alpha}(Y)$ , e.g. each based on separate segments of the data sample, this would result in overall rms convergence rate no less than  $O(n^{-1/2})$  ( $O(n^{-1/2})$  contributed by minimax estimation of the two-dimensional function and  $O(n^{-1})$  contributed by minimax estimation of the one dimensional functions). Note that this plug-in estimation procedure requires a.e. positive estimates of  $f^{1-\alpha}(X)$  and  $f^{1-\alpha}(Y)$  or at least these marginal density estimates should dominate the joint estimate of  $f(X, Y)^\alpha$ .

A hybrid method of estimation of the mutual  $\alpha$ -information which has faster  $O(n^{-2/3})$  rms convergence rate is the following. Step 1: generate estimates  $\hat{f}(X)$  and  $\hat{f}(Y)$  of the one-dimensional functions using minimax density estimation applied to  $X$  and  $Y$  components. Such estimates will converge a.s. with rms rates greater than  $n^{-1}$ , when the marginals densities are of bounded variation. Step 2: apply the separable measure transformation  $dx dy \rightarrow \hat{f}(x)\hat{f}(y)$  to the plane, as described in the previous section, and implement the MST estimator on the the  $(X, Y)$  realization imbedded into the transformed coordinates. The resultant estimator will converge a.s. to the mutual  $\alpha$ -information with rms convergence rate bounded above by  $O(n^{-2/3})$ . This procedure easily generalizes to estimating the  $\alpha$ -mutual information of a  $d$  dimensional sample  $(X, Y, \dots, Z)$  for which the rms convergence rates of the minimax plug in estimator is no less than  $O(n^{-1/(d+2)})$  and that of the hybrid estimator is  $O(n^{-1/(d+1)})$ . For an application of multi-dimensional  $\alpha$ -mutual information estimation to image registration see Neemuchwalla and *etal* [33].

The one dimensional density estimation step in the hybrid procedure can be circumvented by considering a related measure to the mutual  $\alpha$ -information: the mutual  $\alpha$ -entropy difference

$$\begin{aligned} \Delta_0 H_\alpha(X, Y) &= H_\alpha(X, Y) - H_\alpha(X) - H_\alpha(Y) \\ &= \frac{1}{1 - \alpha} \ln \frac{\int f^\alpha(X, Y) dX dY}{\int f^\alpha(X) dX \int f^\alpha(Y) dY} \end{aligned} \quad (26)$$

which also converges to the standard Shannon mutual information in the limit as  $\alpha \rightarrow 1$ . Given an i.i.d. sample  $\{(X_i, Y_i)\}_{i=1}^n$  the length of the MST  $L_\gamma(\{(X_i, Y_i)\}_{i=1}^n)/n^\alpha$  converges w.p.1 to the numerator of (26) times the scale factor  $\beta_{L_\gamma, d}$ ,  $\alpha = (d - \gamma)/d$ . Furthermore, let  $\{\pi(i)\}_{i=1}^n$  be a permutation function, selected at random. Then  $L_\gamma(\{(X_{\pi(i)}, Y_{\pi(i)})\}_{i=1}^n)/n^\alpha$  converges w.p.1 to the denominator of (26) times the same scale factor. It can be



concluded that a consistent estimator of  $\Delta_0 H_\alpha(X, Y)$  is given by the ratio

$$\widehat{\Delta_0 H_\alpha}(X, Y) = \frac{1}{1 - \alpha} \ln \frac{L_\gamma(\{(X_i, Y_i)\}_{i=1}^n)}{L_\gamma(\{X_{\pi(i)}, Y_{\pi(i)}\}_{i=1}^n)}$$

which does not depend on the factor  $\beta_{L_\gamma, d}$ . By comparing this statistic to a threshold we obtain a simple test for dependence of two random variables  $X, Y$  based on  $n$  i.i.d. observations. To reduce bias for finite  $n$  it is suggested that the denominator  $L_\gamma \stackrel{\text{def}}{=} L_\gamma(\{X_{\pi(i)}, Y_{\pi(i)}\}_{i=1}^n)$  be replaced by a sample average  $\overline{L_\gamma^\pi} = 1/|\Pi| \sum_{\pi \in \Pi} L_\gamma^\pi$  where  $\Pi$  is a set of randomly selected permutation functions  $\pi$ . When  $f(X, Y)^\alpha$ ,  $f^{1-\alpha}(X)$  and  $f^{1-\alpha}(Y)$  are of bounded variation the minimax rms convergence rate of the mutual  $\alpha$ -entropy difference will be  $O(n^{-2/3})$ .

## 5 Application to Geo-Registration

It is desired to register two images taken on different sensor planes by potentially different sensor modalities for geo-registration applications. Our objective is to register two types of images — a set of electro-optical(EO) images and a terrain height map. For this multisensor image registration problem, there usually exists distortions between the two types of images. The distortions are due to difference acquisition conditions of the images such as shadowing, diffraction, terrain changes over time, clouds blocking the illumination sources, seasonal variations, etc. Existence of such differences between the images to be registered requires that the registration algorithms to be robust to noise and other small perturbations in intensity values. Here we describe an application of minimal graph entropy estimation to a feature set which is the set of gray level pixels.

For this image registration problem the set of EO images are generated from the *a priori* digital elevation model (DEM)<sup>1</sup> of a terrain patch (the terrain height map) at different look angles (determined by the sensor's location) and with different lighting positions. With different sensor and light locations, we can simulate the distortions mentioned above. For example, shadows are generated by taking into account both the sensor location and the lighting location as follows. The scene is first rendered using the lighting source as the viewing location. Depth values (distance from the light source) are generated for all pixels in the scene and stored in a depth buffer. Next, the scene is rendered using the sensor's location as the viewpoint. Before drawing each pixel, its depth value as measured from the sensor is compared to the transformed depth value as measured from the light source. This comparison determines if a particular pixel is illuminated by the source. Shadows are placed on those pixels that fail this comparison.

Geo-registration of a EO reference image to DEM's in an image database is accomplished by selecting a candidate DEM image from the database and projecting it into the EO image plane of the reference image. The objective is to find the correct viewing angle such that the corresponding EO image is the best match to the EO reference image. Figure 6 shows an DEM projected into the EO image plane with viewing angles (290, -20, 130) and the reference EO image. Clearly they are not aligned.

For matching criterion we use the  $\alpha$ -Jensen difference, with  $\alpha$  chosen arbitrarily as 0.5, applied to grey level features extracted from the reference images and candidate EO images derived from the DEM database. For illustration purposes we selected a very simple set of features via stratified sampling of the grey levels with centroid refinements. This sampling method produces a set of  $n$  three dimensional feature vectors  $Z_i = (x_i, y_i, F(x_i, y_i))$  where  $F(x, y)$  is a sample of the grey level at planar position  $x, y$ . The points  $\{(x_i, y_i)\}_{i=1}^n$  approximate the centroids of Voronoi cells and  $\{F(x_i, y_i)\}_{i=1}^n$  correspond to the set of  $n$  samples of the image from which we could reconstruct the original image with minimum mean square error. For more details see [30]. When the union of features from reference and target images are rendered as points in three dimensions we obtain a point cloud of features over which the MST can be constructed and the Jensen difference estimated.

<sup>1</sup>DEM stores the terrain height information in a three dimensional array where each element of the array consists of the locations (x and y coordinates) and the height of the terrain at that location.

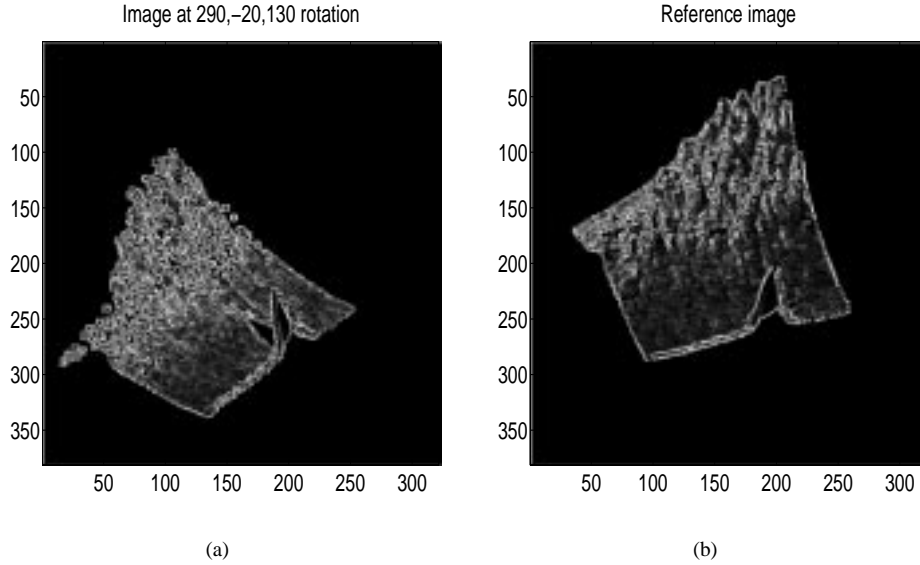


Figure 6: Misaligned EO and reference images

Figure 7 illustrates the MST-based registration procedure over the union of the reference and candidate image features for misaligned images, while Figure 8 shows the same for aligned images. In both Figures 7(a) and 8(a), circle points denote the pixels from Image  $X_1$  and cross points denote the pixels from Image  $X_0$ . From Figures 7(a) and 8(a) we see that for misaligned images, the representation points have larger distances than those for aligned images. Therefore the corresponding MST for the misaligned images has a longer length than that for the aligned images (Figures 7(b) and 8(b)).

We repeat this MST construction process over the union of reference features and features derived from each of the images in the DEM database. The MST length can then be plotted in Figure 9. The x-axis stands for the image index, which corresponds to the viewing angles from the aircraft. The minimum of MST length indicates the best matching of the EO image and the reference image, which corresponds to the registered pair in Figure 10.

## 6 Conclusion

In this report we have discussed and compared  $\alpha$ -divergence and  $\alpha$ -entropy estimation techniques using minimal graph estimation and density plug-in methods. We have also considered the  $\alpha$ -Jensen difference for performing indexing and image retrieval. We have investigated the estimation of  $\alpha$ -Jensen difference using density plug-in estimators and the MST minimal graph method. We demonstrated theoretical advantages of the latter method for indexing planar features or higher dimensional features with feature densities of bounded variation.

### Appendix A

**Proposition 5** Let  $f_{\frac{1}{2}} \stackrel{\text{def}}{=} \frac{1}{2}(f_0 + f_1)$ . The following local representation of the fractional Rényi entropy of a convex mixture  $\beta f_0 + (1 - \beta)f_1$  holds for all  $\alpha, \beta \in [0, 1]$ :

$$H_\alpha(\beta f_0 + (1 - \beta)f_1)$$

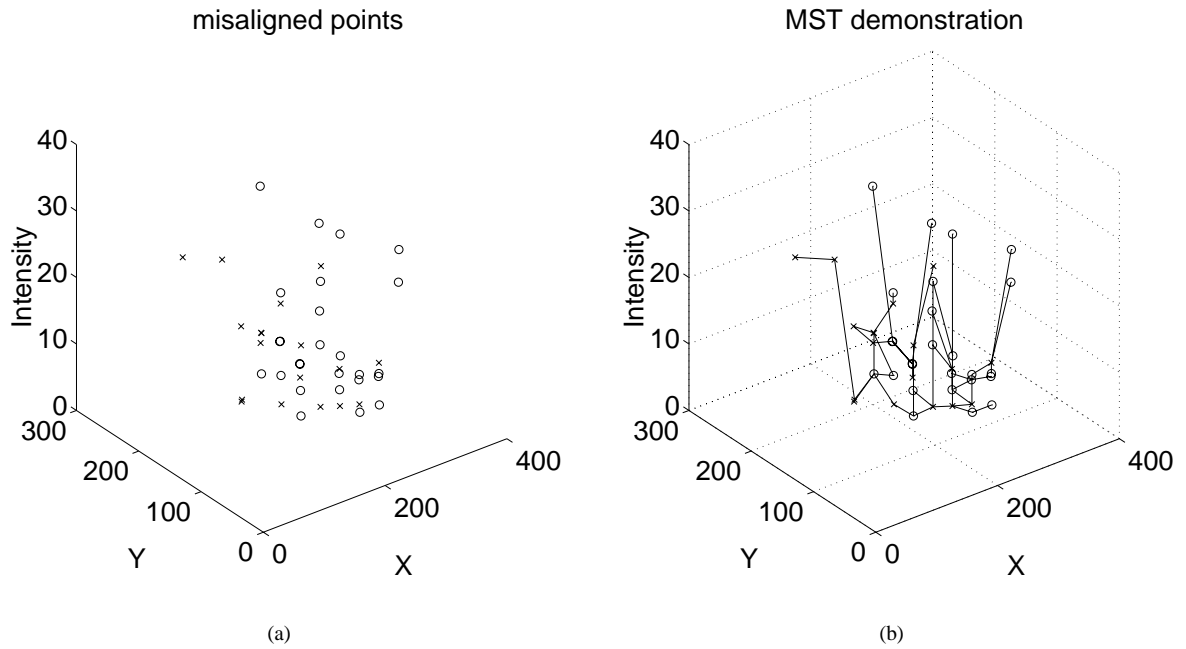


Figure 7: MST demonstration for misaligned images

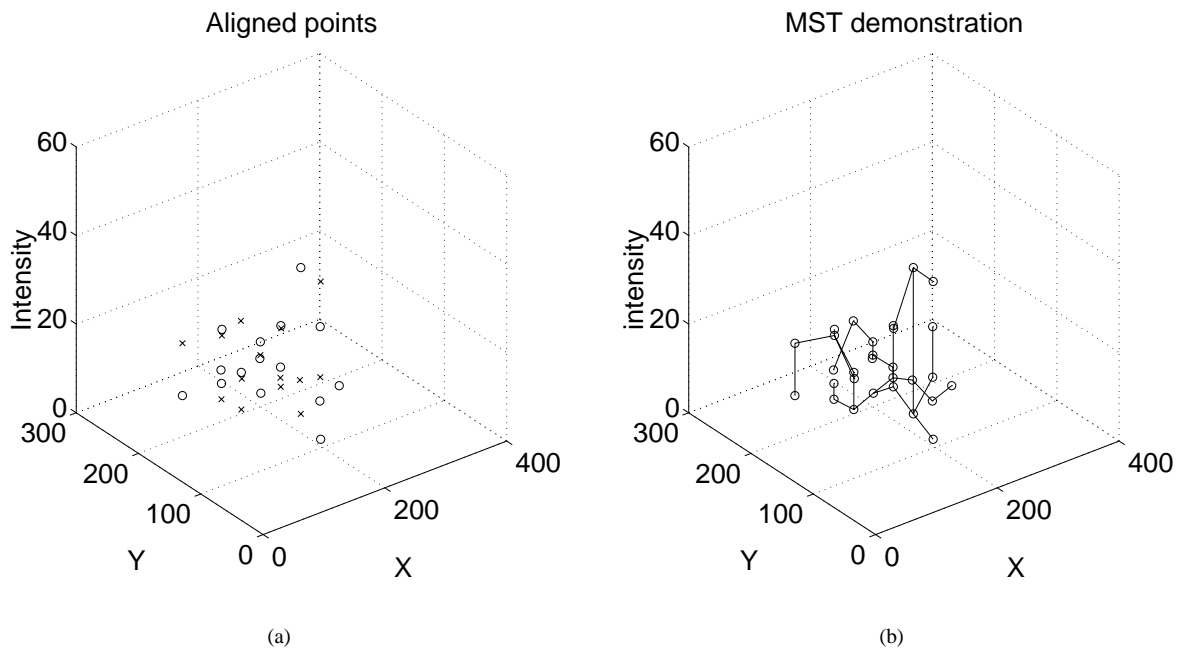


Figure 8: MST demonstration for aligned images

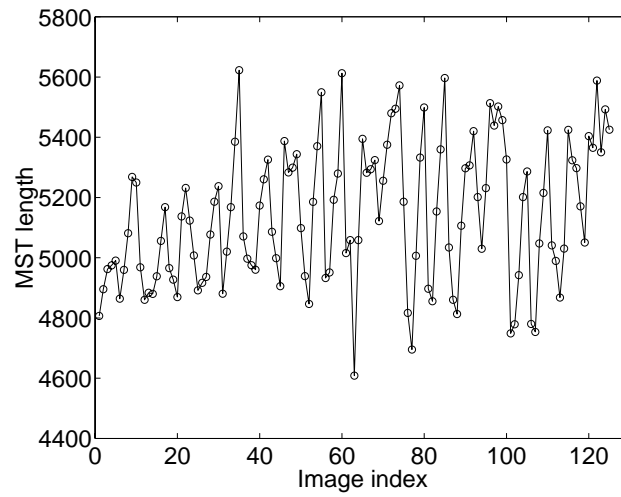


Figure 9: MST length for different test-reference image pairs

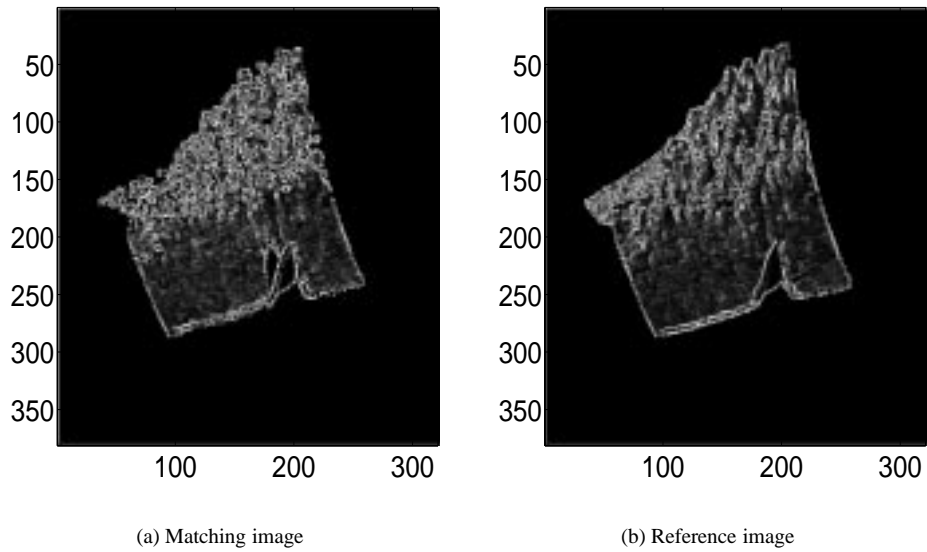


Figure 10: Co-registered EO-terrain maps

$$\begin{aligned}
= & H_\alpha(f_{\frac{1}{2}}) + \frac{\alpha}{1-\alpha} \left( \beta - \frac{1}{2} \right) \frac{\int f_{\frac{1}{2}}^\alpha(x) \left( \frac{f_0(x)-f_1(x)}{f_{\frac{1}{2}}(x)} \right) dx}{\int f_{\frac{1}{2}}^\alpha(x) dx} + \frac{\alpha}{2} \left( \frac{2\beta-1}{2} \right)^2 \frac{\int \left( \frac{f_0(x)-f_1(x)}{f_{\frac{1}{2}}(x)} \right)^2 f_{\frac{1}{2}}^\alpha(x) dx}{\int f_{\frac{1}{2}}^\alpha(x) dx} \\
& - \frac{\alpha^2}{2(1-\alpha)} \left( \frac{2\beta-1}{2} \right)^2 \left( \frac{\int \frac{f_0(x)-f_1(x)}{f_{\frac{1}{2}}(x)} f_{\frac{1}{2}}^\alpha(x) dx}{\int f_{\frac{1}{2}}^\alpha(x) dx} \right)^2 + o(\Delta^2)
\end{aligned} \tag{27}$$

where  $\Delta \in [0, 1]$  is

$$\Delta \stackrel{\text{def}}{=} \sup_x \frac{|f_0(x) - f_1(x)|}{f_{\frac{1}{2}}(x)}. \tag{28}$$

*proof*

Let  $f_{1-\beta}(x) = \beta f_0(x) + (1-\beta)f_1(x)$ . It can be written as

$$\begin{aligned}
f_{1-\beta}(x) &= \frac{1}{2}[f_0(x) + (1-\beta)(f_1(x) - f_0(x))] + \frac{1}{2}[f_0(x) + \beta(f_0(x) - f_1(x))] \\
&= f_{\frac{1}{2}}(x) + \frac{1}{2}(2\beta-1)(f_0(x) - f_1(x)) \\
&= f_{\frac{1}{2}}(x) \left( 1 + \frac{(2\beta-1)\Delta_x}{2f_{\frac{1}{2}}(x)} \right)
\end{aligned}$$

where  $\Delta_x = (f_0(x) - f_1(x))/f_{\frac{1}{2}}(x)$ . A Taylor series expansion of  $f_{1-\beta}^\alpha(x)$  yields

$$\begin{aligned}
f_{1-\beta}^\alpha(x) &= f_{\frac{1}{2}}^\alpha(x) \left( 1 + \frac{2\beta-1}{2} \Delta_x \right)^\alpha \\
&= f_{\frac{1}{2}}^\alpha(x) \left( 1 + \frac{\alpha(2\beta-1)}{2} \Delta_x + \frac{\alpha(\alpha-1)}{2} \left( \frac{2\beta-1}{2} \Delta_x \right)^2 + o(\Delta_x^2) \right)
\end{aligned} \tag{29}$$

Taking the logarithm of both sides of (29) and dividing by  $1-\alpha$

$$\begin{aligned}
& \frac{1}{1-\alpha} \ln \int f_{1-\beta}^\alpha(x) dx \\
= & \frac{1}{1-\alpha} \ln \int \left[ f_{\frac{1}{2}}^\alpha(x) + \alpha f_{\frac{1}{2}}^\alpha(x) \left( \frac{2\beta-1}{2} \Delta_x \right) \right. \\
& \quad \left. + \frac{\alpha(\alpha-1)}{2} f_{\frac{1}{2}}^\alpha(x) \left( \frac{2\beta-1}{2} \Delta_x \right)^2 + f_{\frac{1}{2}}^\alpha(x) o(\Delta_x^2) \right] dx \\
= & \frac{1}{1-\alpha} \ln \left\{ \int f_{\frac{1}{2}}^\alpha(x) dx \left[ 1 + \frac{\alpha \int f_{\frac{1}{2}}^\alpha(x) \left( \frac{2\beta-1}{2} \Delta_x \right) dx}{\int f_{\frac{1}{2}}^\alpha(x) dx} \right. \right. \\
& \quad \left. \left. + \frac{\frac{\alpha(\alpha-1)}{2} \int f_{\frac{1}{2}}^\alpha(x) \left( \frac{2\beta-1}{2} \Delta_x \right)^2 dx}{\int f_{\frac{1}{2}}^\alpha(x) dx} + o(\Delta^2) \right] \right\}
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{1-\alpha} \ln \int f_{\frac{1}{2}}^{\alpha}(x) dx + \frac{1}{1-\alpha} \ln \left[ 1 + \frac{\alpha \int f_{\frac{1}{2}}^{\alpha}(x) \left( \frac{2\beta-1}{2} \Delta_x \right) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} \right. \\
&\quad \left. + \frac{\frac{\alpha(\alpha-1)}{2} \int f_{\frac{1}{2}}^{\alpha}(x) \left( \frac{2\beta-1}{2} \Delta_x \right) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} + o(\Delta^2) \right].
\end{aligned}$$

Since  $\ln(1+x) = x - \frac{x^2}{2} + o(x^2)$ , we have

$$\begin{aligned}
H_{\alpha}(\beta f_0 + (1-\beta)f_1) &= \frac{1}{1-\alpha} \ln \int f_{1-\beta}^{\alpha}(x) dx \\
&= H_{\alpha}(f_{\frac{1}{2}}) + \frac{\alpha}{1-\alpha} \frac{2\beta-1}{2} \frac{\int \Delta_x f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} + \frac{\alpha}{2} \left( \frac{2\beta-1}{2} \right)^2 \frac{\int \Delta_x^2 f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} \\
&\quad - \frac{\alpha^2}{2(1-\alpha)} \left( \frac{2\beta-1}{2} \right)^2 \frac{\left( \int \Delta_x f_{\frac{1}{2}}^{\alpha}(x) dx \right)^2}{\left( \int f_{\frac{1}{2}}^{\alpha}(x) dx \right)^2} + o(\Delta^2). \tag{30}
\end{aligned}$$

This completes the proof of Proposition 5.  $\square$

**Proposition 6** *The following asymptotic representation of the fractional Jensen difference of two densities  $f_0$  and  $f_1$  holds for all  $\alpha, \beta \in [0, 1]$ :*

$$\begin{aligned}
&\Delta H_{\alpha}(\beta, f_0, f_1) \\
&= \frac{\alpha\beta(1-\beta)}{2} \left( \frac{\int \left( \frac{f_0(x)-f_1(x)}{f_{\frac{1}{2}}(x)} \right)^2 f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} + \frac{\alpha}{1-\alpha} \left( \frac{\int \left( \frac{f_0(x)-f_1(x)}{f_{\frac{1}{2}}(x)} \right) f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} \right)^2 \right) \\
&\quad + o(\Delta^2) \tag{31}
\end{aligned}$$

where  $f_{\frac{1}{2}}$  and  $\Delta$  are as defined in Proposition 5.

*proof*

Specializing to  $\beta = 0$  and  $\beta = 1$  in (27) of Proposition 5 we obtain

$$\begin{aligned}
H_{\alpha}(f_0) &= H_{\alpha}(f_{\frac{1}{2}}) + \frac{\alpha}{1-\alpha} \frac{\int \frac{1}{2} \Delta_x f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} + \frac{\alpha}{2} \frac{\int \left( \frac{1}{2} \Delta_x \right)^2 f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} \\
&\quad - \frac{\alpha^2}{2(1-\alpha)} \left( \frac{\int \frac{1}{2} \Delta_x f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} \right)^2 + o(\Delta^2) \tag{32}
\end{aligned}$$

$$\begin{aligned}
H_{\alpha}(f_1) &= H_{\alpha}(f_{\frac{1}{2}}) - \frac{\alpha}{1-\alpha} \frac{\int \frac{1}{2} \Delta_x f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} + \frac{\alpha}{2} \frac{\int \left( \frac{1}{2} \Delta_x \right)^2 f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} \\
&\quad - \frac{\alpha^2}{2(1-\alpha)} \left( \frac{\int \frac{1}{2} \Delta_x f_{\frac{1}{2}}^{\alpha}(x) dx}{\int f_{\frac{1}{2}}^{\alpha}(x) dx} \right)^2 + o(\Delta^2) \tag{33}
\end{aligned}$$

where  $\Delta_x = (f_0(x) - f_1(x))/f_{\frac{1}{2}}(x)$ . Substituting (27), (32) and (33) into (20), we obtain the expression (31) for the Jensen difference. This completes the proof of Proposition 6  $\square$

**Proposition 7** *The  $\alpha$ -divergence of fractional order  $\alpha \in (0, 1)$  between two densities  $f_0$  and  $f_1$  has the local representation*

$$D_\alpha(f_0 \| f_1) = \frac{\alpha}{4} \int f_{\frac{1}{2}}(x) \left( \frac{f_0(x) - f_1(x)}{f_{\frac{1}{2}}(x)} \right)^2 dx + o(\Delta^2) \quad (34)$$

where  $f_{\frac{1}{2}}$  and  $\Delta$  are as defined in Proposition 5.

*proof*

Rewrite the density  $f_0$  as

$$f_0(x) = \frac{1}{2}(f_0(x) + f_1(x)) + \frac{1}{2}(f_0(x) - f_1(x)) = f_{\frac{1}{2}}(x)(1 + \frac{1}{2}\Delta_x), \quad (35)$$

where  $\Delta_x = (f_0(x) - f_1(x))/f_{\frac{1}{2}}(x)$ . Similarly,

$$f_1(x) = f_{\frac{1}{2}}(x)(1 - \frac{1}{2}\Delta_x). \quad (36)$$

Thus, by Taylor series expansion, we have

$$\begin{aligned} f_1^\alpha(x) &= f_{\frac{1}{2}}^\alpha(x) - \alpha f_{\frac{1}{2}}^{\alpha-1}(x) \left( \frac{\Delta_x}{2} \right) + \frac{\alpha(\alpha-1)}{2} f_{\frac{1}{2}}^{\alpha-2}(x) \left( \frac{\Delta_x}{2} \right)^2 + o(\Delta_x^2) \\ f_0^{1-\alpha}(x) &= f_{\frac{1}{2}}^{1-\alpha}(x) + (1-\alpha) f_{\frac{1}{2}}^{1-\alpha-1}(x) \left( \frac{\Delta_x}{2} \right) + \frac{\alpha(1-\alpha)}{2} f_{\frac{1}{2}}^{1-\alpha-2}(x) \left( \frac{\Delta_x}{2} \right)^2 + o(\Delta_x^2). \end{aligned}$$

Therefore

$$f_1^\alpha(x) f_0^{1-\alpha}(x) = f_{\frac{1}{2}}(x) \left[ 1 - (2\alpha-1) \frac{\Delta_x}{2} - \alpha(1-\alpha) \left( \frac{\Delta_x}{2} \right)^2 + o(\Delta_x^3) \right]$$

and

$$\begin{aligned} D_\alpha(f_0 \| f_1) &= \frac{1}{\alpha-1} \ln \int f_1^\alpha(x) f_0^{1-\alpha}(x) dx \\ &= \frac{1}{\alpha-1} \ln \int \left( f_{\frac{1}{2}}(x) - (2\alpha-1) f_{\frac{1}{2}}(x) \frac{\Delta_x}{2} - \alpha(1-\alpha) f_{\frac{1}{2}}(x) \left( \frac{\Delta_x}{2} \right)^2 + f_{\frac{1}{2}}(x) o(\Delta_x^2) \right) dx \\ &= \frac{1}{\alpha-1} \ln \left( 1 - \alpha(1-\alpha) \int f_{\frac{1}{2}}(x) \left( \frac{\Delta_x}{2} \right)^2 dx + o(\Delta^2) \right) \\ &= \alpha \int f_{\frac{1}{2}}(x) \left( \frac{\Delta_x}{2} \right)^2 dx + o(\Delta^2) \\ &= \frac{\alpha}{4} \int f_{\frac{1}{2}}(x) \left( \frac{f_0(x) - f_1(x)}{f_{\frac{1}{2}}(x)} \right)^2 dx + o(\Delta^2). \end{aligned}$$

This completes the proof of Prop. 7.  $\square$

## Appendix B

**Proposition 8** Let  $f_1(x) = f(x; \mu_1, \Lambda_1)$  and  $f_0(x) = f(x; \mu_0, \Lambda_0)$  be multivariate  $d$ -dimensional Gaussian densities with vector means  $\mu_1, \mu_0$  and positive definite covariance matrices  $\Lambda_1, \Lambda_0$ . The Rényi divergence of order  $\alpha$  between  $f_1$  and  $f_0$  is

$$D_\alpha(f_1||f_0) = \frac{1/2}{\alpha - 1} \ln \frac{|\Lambda_0|^\alpha |\Lambda_1|^{1-\alpha}}{|\alpha\Lambda_0 + (1-\alpha)\Lambda_1|} + \frac{\alpha}{2} \Delta\mu^T (\alpha\Lambda_0 + (1-\alpha)\Lambda_1)^{-1} \Delta\mu \quad (37)$$

where  $\Delta\mu = \mu_1 - \mu_0$ .

*Proof*

Start from the definition

$$D_\alpha(f_1||f_0) = \frac{1}{\alpha - 1} \ln \int f^\alpha(x; \mu_1, \Lambda_1) f^{1-\alpha}(x; \mu_0, \Lambda_0) dx$$

and make a change of variable  $y = \Lambda_0^{-\frac{1}{2}}(x - \mu_0)$  in the integral to obtain

$$\int f^\alpha(x; \mu_1, \Lambda_1) f^{1-\alpha}(x; \mu_0, \Lambda_0) dx = |\Lambda_0|^{\frac{1}{2}} \int f^\alpha(y; \Lambda_0^{-\frac{1}{2}} \Delta\mu, \Lambda_0^{-\frac{1}{2}} \Lambda_1 \Lambda_0^{-\frac{1}{2}}) f^{1-\alpha}(y; 0, I_d) dy, \quad (38)$$

where  $I_d$  is the  $d \times d$  identity matrix.

By completion of the square and elementary matrix manipulations it is straightforward to show that for any  $d$ -element vector  $m$  and positive definite  $d \times d$  covariance matrix  $A$

$$\begin{aligned} & \int f^\alpha(y; m, A) f^{1-\alpha}(y; 0, I_d) dy \\ &= \left( \frac{|A|^{1-\alpha}}{|(1-\alpha)A + \alpha I_d|} \right)^{\frac{1}{2}} \exp\left(-\frac{1}{2}\alpha(1-\alpha) m^T [(1-\alpha)A + \alpha I_d]^{-1} m\right) \end{aligned}$$

Finally, identifying  $m = \Lambda_0^{-\frac{1}{2}} \Delta\mu$  and  $\Lambda = \Lambda_0^{-\frac{1}{2}} \Lambda_1 \Lambda_0^{-\frac{1}{2}}$ , substitution of the above into (38) and performing some matrix algebra we obtain (37).

This completes the proof of Prop. 8. □

## References

- [1] I. Ahmad and P.-E. Lin, “A nonparametric estimation of the entropy for absolutely continuous distributions,” *IEEE Trans. on Inform. Theory*, vol. IT-22, pp. 664–668, 1976.
- [2] M. Basseville, “Distance measures for signal processing and pattern recognition,” *Signal Processing*, vol. 18, pp. 349–369, 1989.
- [3] J. Beardwood, J. H. Halton, and J. M. Hammersley, “The shortest path through many points,” *Proc. Cambridge Philosophical Society*, vol. 55, pp. 299–327, 1959.
- [4] J. Beirlant, E. J. Dudewicz, L. Györfi, and E. van der Meulen, “Nonparametric entropy estimation: an overview,” *Intern. J. Math. Stat. Sci.*, vol. 6, no. 1, pp. 17–39, june 1997.



- [5] R. J. Beran, "Minimum Hellinger distance estimates for parametric models," *Annals of Statistics*, vol. 5, pp. 445–463, 1977.
- [6] T. Cover and J. Thomas, *Elements of Information Theory*, Wiley, New York, 1987.
- [7] I. Csiszár, "Information-type measures of divergence of probability distributions and indirect observations," *Studia Sci. Math. Hung.*, vol. 2, pp. 299–318, 1967.
- [8] A. Dembo and O. Zeitouni, *Large deviations techniques and applications*, Springer-Verlag, NY, 1998.
- [9] R. A. Devore and V. A. Popov, "Interpolation of Besov spaces," *Trans. Amer. Math. Soc.*, vol. 305, no. 1, pp. 397–414, 1988.
- [10] M. N. Do and M. Vetterli, "Texture similarity measurement using kullback-liebler distance on wavelet subbands," in *IEEE Int. Conf. on Image Processing*, pp. 367–370, Vancouver, BC, 2000.
- [11] D. L. Donoho, "One-sided inference about functionals of a density," *Annals of Statistics*, vol. 16, pp. 1390–1420, 1988.
- [12] D. L. Donoho and I. M. Johnstone, "Asymptotic minimaxity of wavelet estimators with sampled data," Technical report, Dept. Statistics, Stanford University, Dec. 1997. <http://www-stat.stanford.edu/~donoho/Reports/>.
- [13] D. L. Donoho, I. M. Johnstone, G. Kerkyacharian, and D. Picard, "Adapting to unknown smoothing via wavelet shrinkage," *J. Am. Statist. Assoc.*, vol. 90, no. 432, pp. 1200–1224, 1995.
- [14] D. L. Donoho, I. M. Johnstone, G. Kerkyacharian, and D. Picard, "Wavelet shrinkage: asymptopia?," *J. Royal Statistical Society, Ser. B*, vol. 57, pp. 301–309, 1995.
- [15] A. Gersho, "Asymptotically optimal block quantization," *IEEE Trans. on Inform. Theory*, vol. IT-28, pp. 373–380, 1979.
- [16] L. Györfi and E. van der Meulen, "Density-free convergence properties of various estimators of entropy," *Comput. Statist. Data Anal.*, vol. 5, pp. 425–436, 1987.
- [17] P. Hall, "On powerful distributional tests based on sample spacings," *Journ. Multivar. Anal.*, vol. 19, pp. 201–225, 1986.
- [18] P. Hall and S. C. Morton, "On the estimation of entropy," *Ann. Inst. Statist. Math.*, vol. 45, pp. 69–88, 1993.
- [19] Y. He, A. Ben-Hamza, and H. Krim, "An information divergence measure for ISAR image registration," *Signal Processing*, Submitted, 2001.
- [20] A. Hero, B. Ma, and O. Michel, "Convergence rates of random minimal graphs," *IEEE Trans. on Inform. Theory*, Submitted Aug. 2001.
- [21] A. Hero and O. Michel, "Robust entropy estimation strategies based on edge weighted random graphs," in *Proc. of Meeting of Intl. Soc. for Optical Engin. (SPIE)*, San Diego, CA, July 1998.
- [22] A. Hero and O. Michel, "Estimation of Rényi information divergence via pruned minimal spanning trees," in *IEEE Workshop on Higher Order Statistics*, Caesaria, Israel, June 1999.
- [23] A. Hero and O. Michel, "Asymptotic theory of greedy approximations to minimal k-point random graphs," *IEEE Trans. on Inform. Theory*, vol. IT-45, no. 6, pp. 1921–1939, Sept. 1999.
- [24] R. Hoffman and A. K. Jain, "A test of randomness based on the minimal spanning tree," *Pattern Recognition Letters*, vol. 1, pp. 175–180, 1983.

- [25] I. A. Ibragimov and R. Z. Has'minskii, *Statistical estimation: Asymptotic theory*, Springer-Verlag, New York, 1981.
- [26] A. Jain, P. Moulin, M. I. Miller, and K. Ramchandran, "Information-theoretic bounds on target recognition performance based on degraded image data," *preprint*, Dec 1999.
- [27] H. Joe, "On the estimation of entropy and other functionals of a multivariate density," *Ann. Inst. Statist. Math.*, vol. 41, pp. 683–697, 1989.
- [28] A. P. Korostelev and A. B. Tsybakov, *Minimax theory of image reconstruction*, Springer-Verlag, New York, 1993.
- [29] S. Kullback and R. Liebler, "On information and sufficiency," *Ann. Math. Statist.*, vol. 22, pp. 79–86, 1951.
- [30] B. Ma, *Parametric and non-parametric approaches for multisensor data fusion*, PhD thesis, University of Michigan, Ann Arbor, MI 48109-2122, 2001.
- [31] B. Ma, A. O. Hero, J. Gorman, and O. Michel, "Image registration with minimal spanning tree algorithm," in *IEEE Int. Conf. on Image Processing*, Vancouver, CA, October 2000.
- [32] A. Mokkadem, "Estimation of the entropy and information of absolutely continuous random variables," *IEEE Trans. on Inform. Theory*, vol. IT-35, no. 1, pp. 193–196, 1989.
- [33] H. Neemuchwala, A. Hero, and P. Carson, "Feature coincidence trees for registration of ultrasound breast images," in *IEEE Int. Conf. on Image Processing*, Thessaloniki, Greece, October 2001.
- [34] D. N. Neuhoff, "On the asymptotic distribution of the errors in vector quantization," *IEEE Trans. on Inform. Theory*, vol. IT-42, pp. 461–468, March 1996.
- [35] J. P. Pluim, J. B. A. Maintz, and M. A. Viergever, " $f$ -information measures in medical image registration," in *Proc. of Meeting of Intl. Soc. for Optical Engin. (SPIE)*, M. Sonka and K. M. Hanson, editors, volume 4322, pp. 579–587, 2001.
- [36] A. Rényi, "On measures of entropy and information," in *Proc. 4th Berkeley Symp. Math. Stat. and Prob.*, volume 1, pp. 547–561, 1961.
- [37] F. Riesz and B. Sz.-Nagy, *Functional analysis*, Ungar, New York, 1955.
- [38] H. Samet, *Applications of spatial data structures : computer graphics, image processing, and GIS*, Addison-Wesley, reading, MA, 1990.
- [39] J. M. Steele, *Probability theory and combinatorial optimization*, volume 69 of *CBMF-NSF regional conferences in applied mathematics*, Society for Industrial and Applied Mathematics (SIAM), 1997.
- [40] R. Stoica, J. Zerubia, and J. M. Francos, "Image retrieval and indexing: A hierarchical approach in computing the distance between textured images," in *IEEE Int. Conf. on Image Processing*, Chicago, 1998.
- [41] R. Stoica, J. Zerubia, and J. M. Francos, "The two-dimensional wold decomposition for segmentation and indexing in image libraries," in *Proc. IEEE Int. Conf. Acoust., Speech, and Sig. Proc.*, Seattle, 1998.
- [42] N. Vasconcelos and A. Lippman, "A Bayesian framework for content-based indexing and retrieval," in *IEEE Data Compression Conference*, Snowbird, Utah, 1998. <http://nuno.www.media.mit.edu/people/nuno/>.
- [43] N. Vasconcelos and A. Lippman, "Bayesian representations and learning mechanisms for content based image retrieval," in *SPIE Storage and Retrieval for Media Databases 2000*, San Jose, CA, 2000. <http://nuno.www.media.mit.edu/people/nuno/>.

- [44] O. Vasicek, "A test for normality based on sample entropy," *J. Royal Statistical Society, Ser. B*, vol. 38, pp. 54–59, 1976.
- [45] P. Viola and W. Wells, "Alignment by maximization of mutual information," in *Proc. of 5th Int. Conf. on Computer Vision, MIT*, volume 1, pp. 16–23, 1995.
- [46] J. E. Yukich, *Probability theory of classical Euclidean optimization*, volume 1675 of *Lecture Notes in Mathematics*, Springer-Verlag, Berlin, 1998.