

Entropic graphs: theory

Alfred O. Hero

Dept. EECS, Dept Biomed. Eng., Dept. Statistics

University of Michigan - Ann Arbor

hero@eecs.umich.edu

<http://www.eecs.umich.edu/~hero>

Collaborators: O. Michel, B. Ma, H. Nemanchwa, J. Costa,

1. Motivating Examples in Imaging and Computer Vision
2. Entropic Feature Similarity/Dissimilarity Measures
3. Entropic Euclidean Graphs over Feature Space
4. Asymptotics of Entropic Graphs
5. Entropic Graphs for Pattern Matching

SOURCE: c:\hero\seminars\Disttinguished\Entropic\ipml_slides.tex

I. Motivating Examples in Imaging and Computer Vision

- Image retrieval and indexing
- Multimodality image fusion
- Inference on shape manifolds
- Image registration

Image Retrieval

QUERY

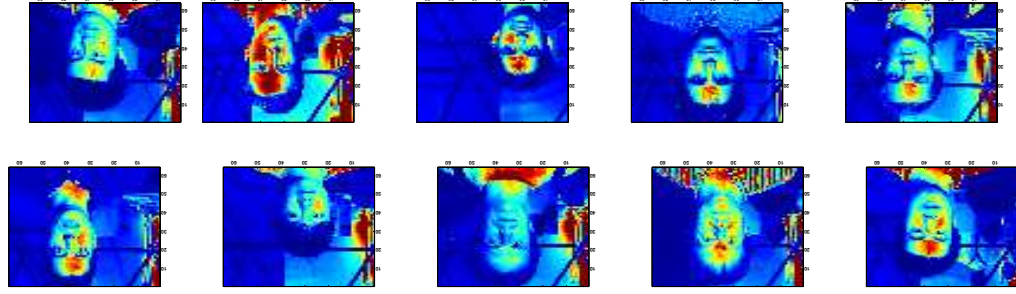
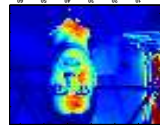


Figure 1: Yale face database <http://cvc.yale.edu/projects/yalefaces/yalefaces.html>

DATABASE

Feature Vectors in Feature Space

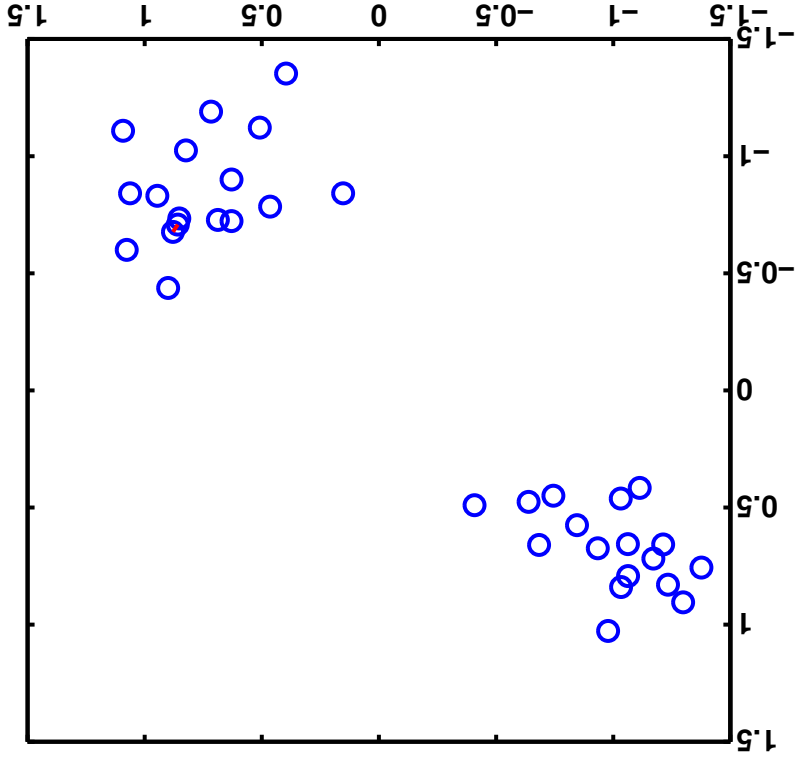
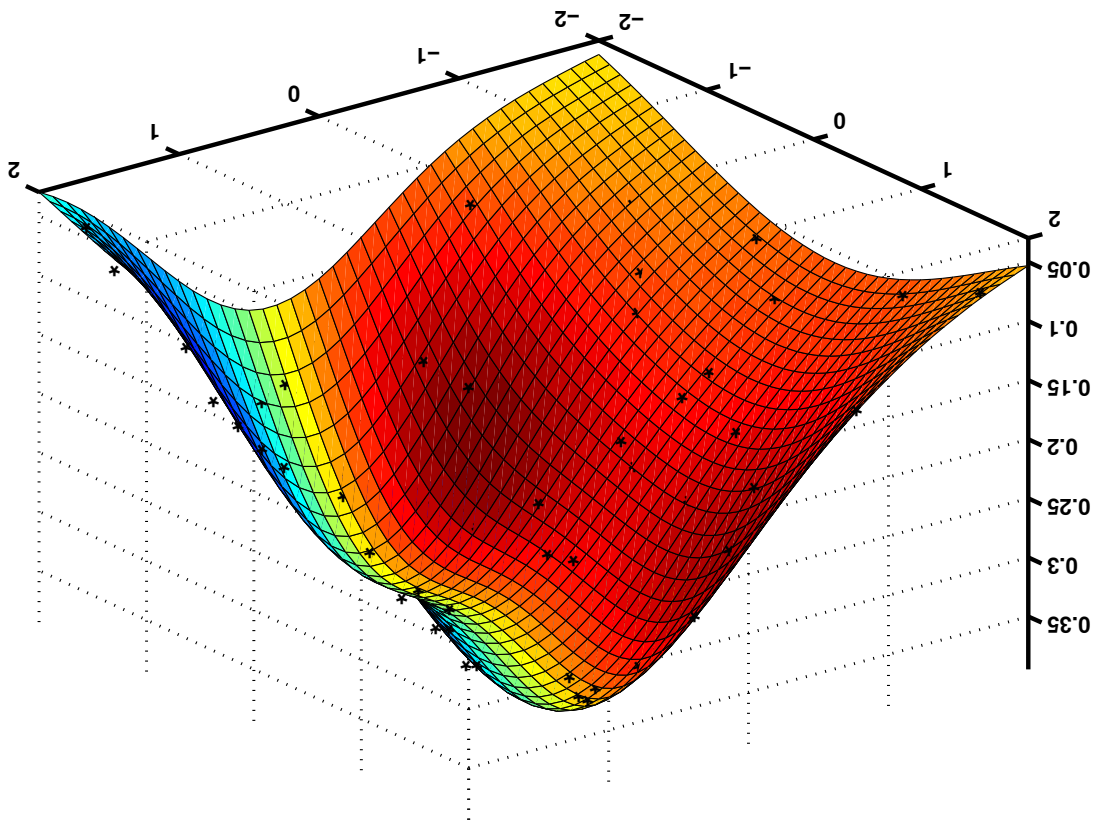


Figure 2: Vectors of projection coefficients extracted from two different images.

Inference on Shape Manifolds

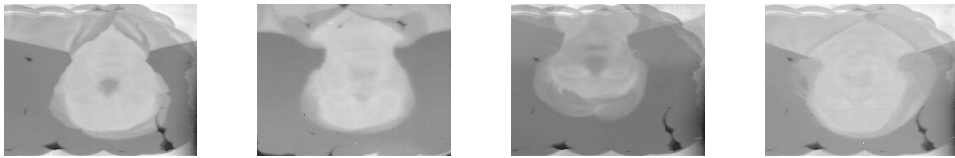


Multi-modality face Retrieval

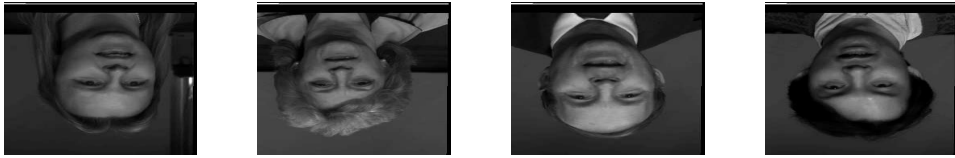
Database of Visible/IR images



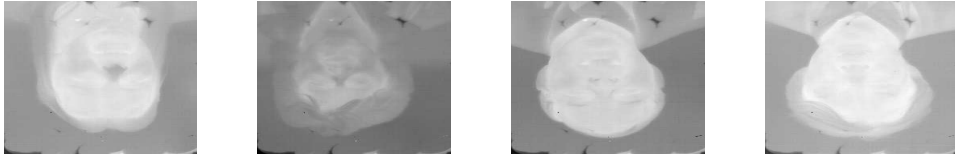
Visible



IR



Visible



IR

<http://www.equinoxsensors.com/products/HID.html>

Image Registration

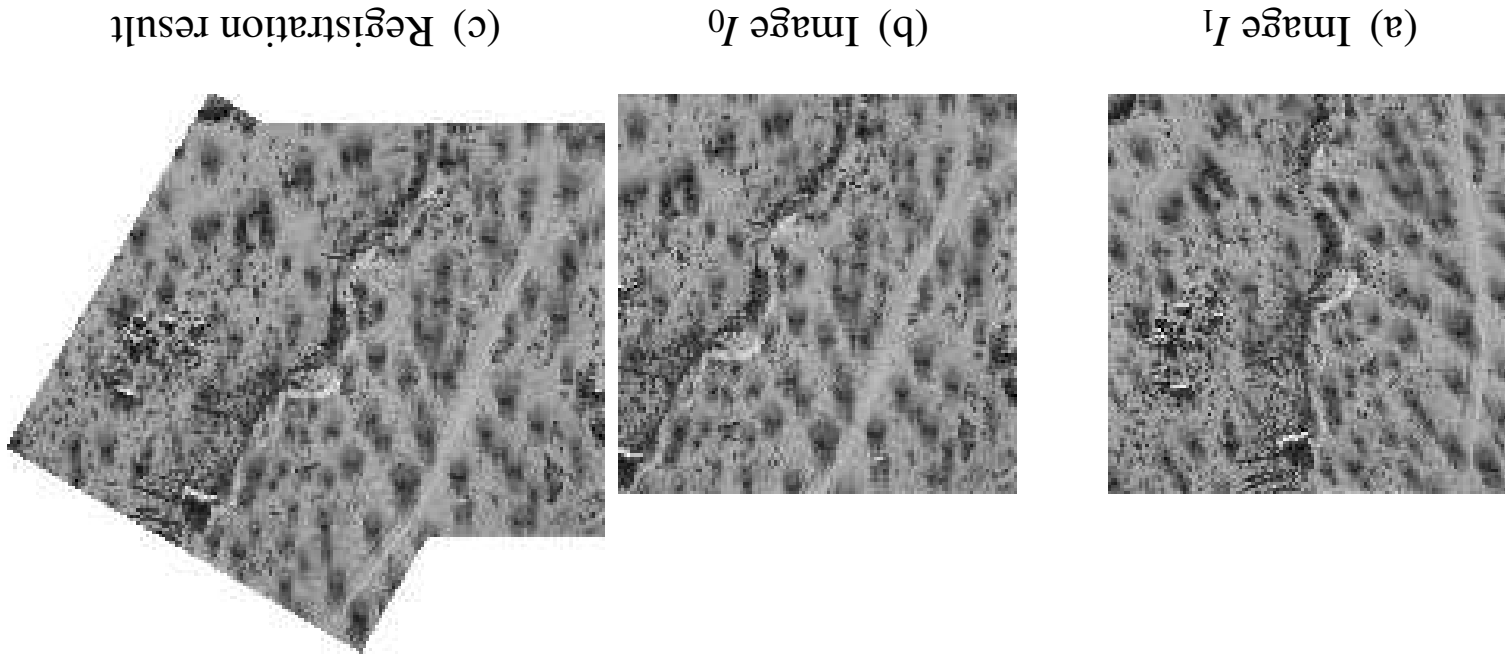
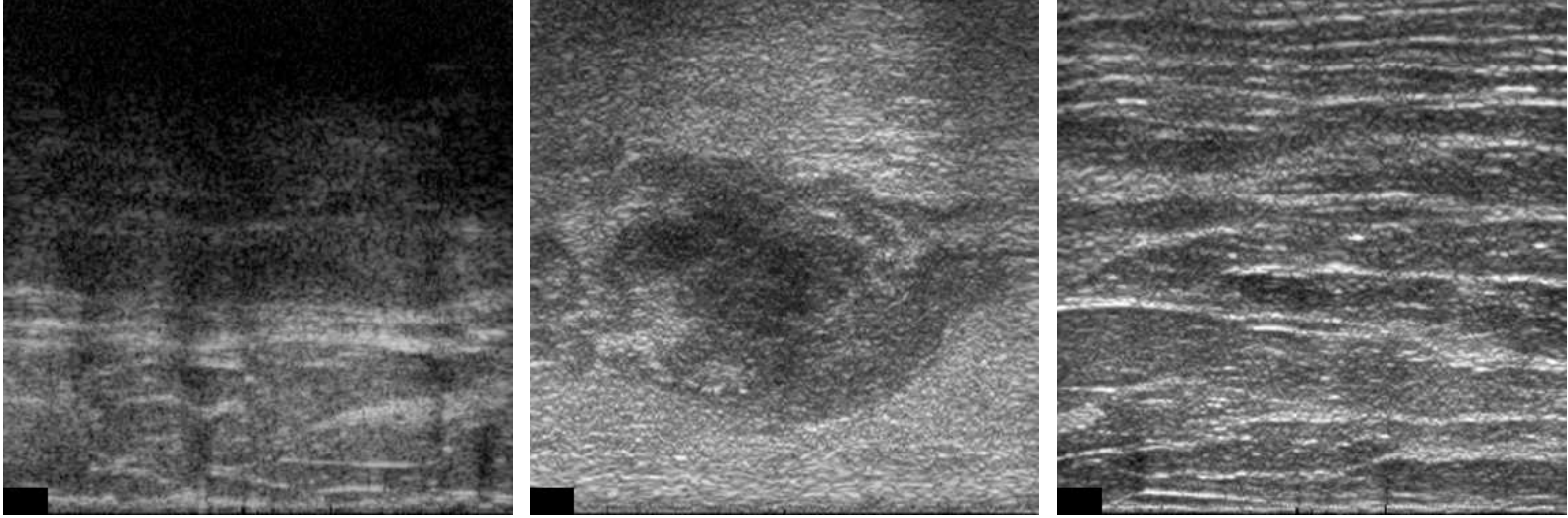


Figure 3: A multirate image registration example

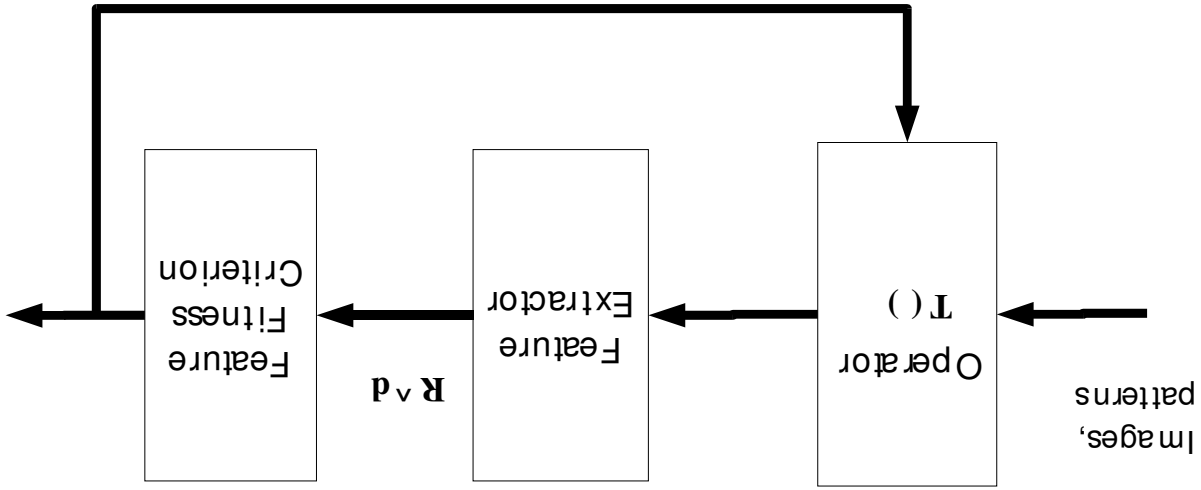
Figure 4: Three ultrasound breast scans. From top to bottom are: case 151, case 142 and case 162.



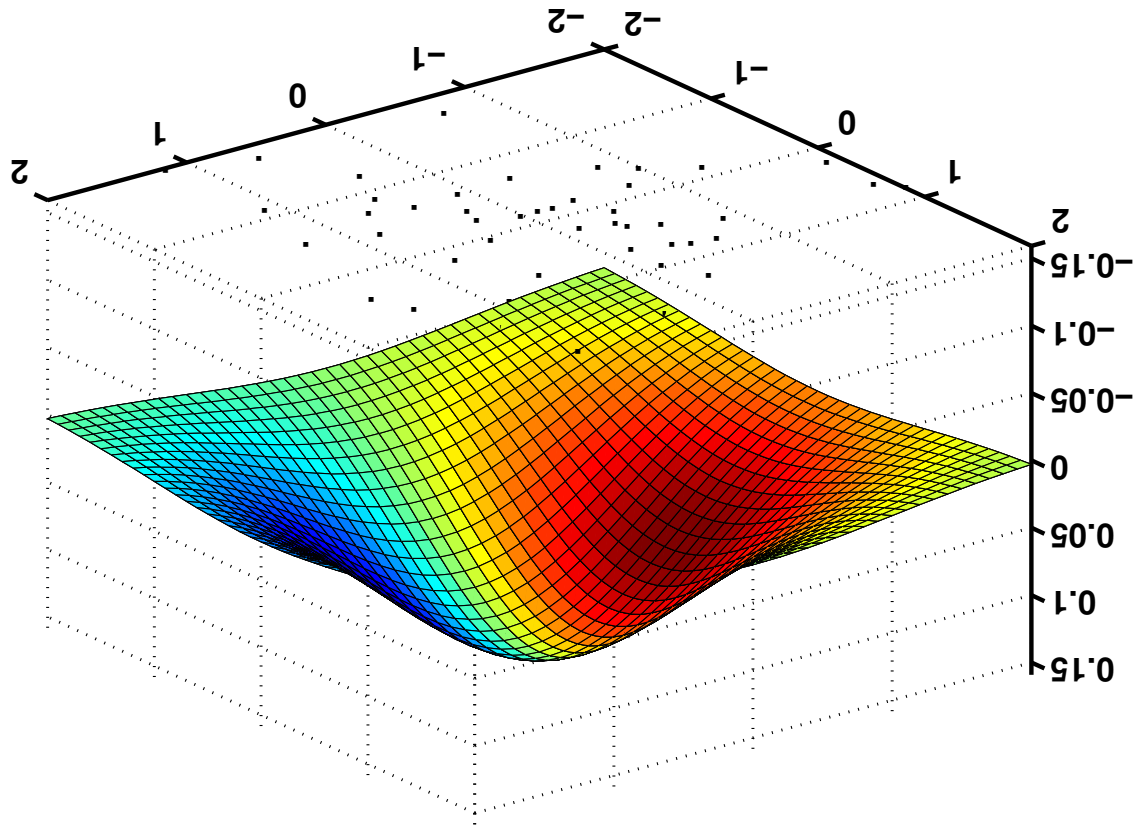
Objective: For given fitness criterion \mathcal{Q} , find operator T which minimizes/maximizes \mathcal{Q}

Our focus: entropic fitness criterion $\mathcal{Q}(f)$

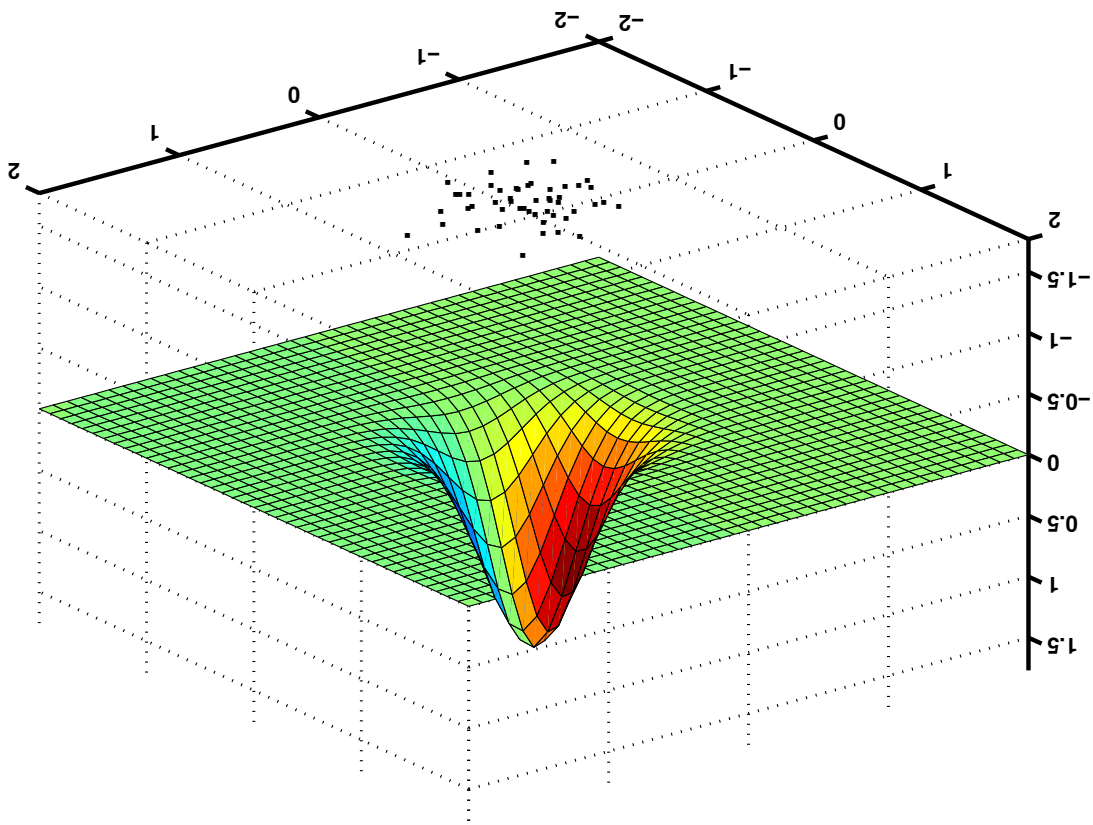
f : feature density over $x \in [0, 1]^d$



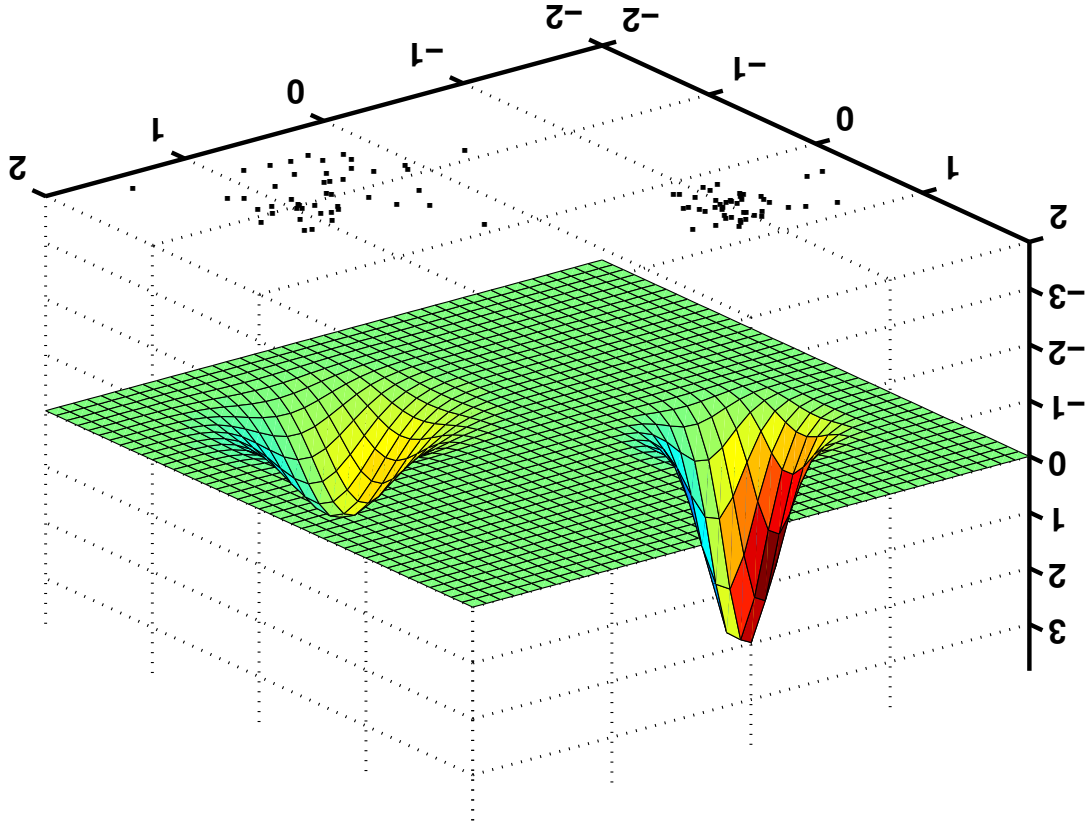
System Block Diagram



High Entropy Feature Density



Low Entropy Feature Density



High Entropy Feature Density

II. Entropic Similarity/Dissimilarity Measures

1. Shannon Entropy of feature density f

$$\tilde{Q}(f) = H(f) = - \int f(x) \ln f(x) dx$$

2. Jensen difference between feature densities f, g :

$$\tilde{Q}(f, g) = H(\epsilon f + (1 - \epsilon)g) - \epsilon H(f) - (1 - \epsilon)H(g)$$

3. KL Divergence between feature densities f, g

$$\tilde{Q}(f, g) = D(f \| g) = \int f(x) \ln \left(\frac{f(x)}{g(x)} \right) dx$$

4. Mutual information between feature sets $f_{X,Y}$

$$\tilde{Q}(f_{X,Y}) = \text{MI}(X, Y) = \iint f_{X,Y}(x, y) \ln \left(\frac{f_{X,Y}(x, y)}{f_X(x) f_Y(y)} \right) dx dy$$

Issue: How to estimate entropic \mathcal{Q} from measured data?

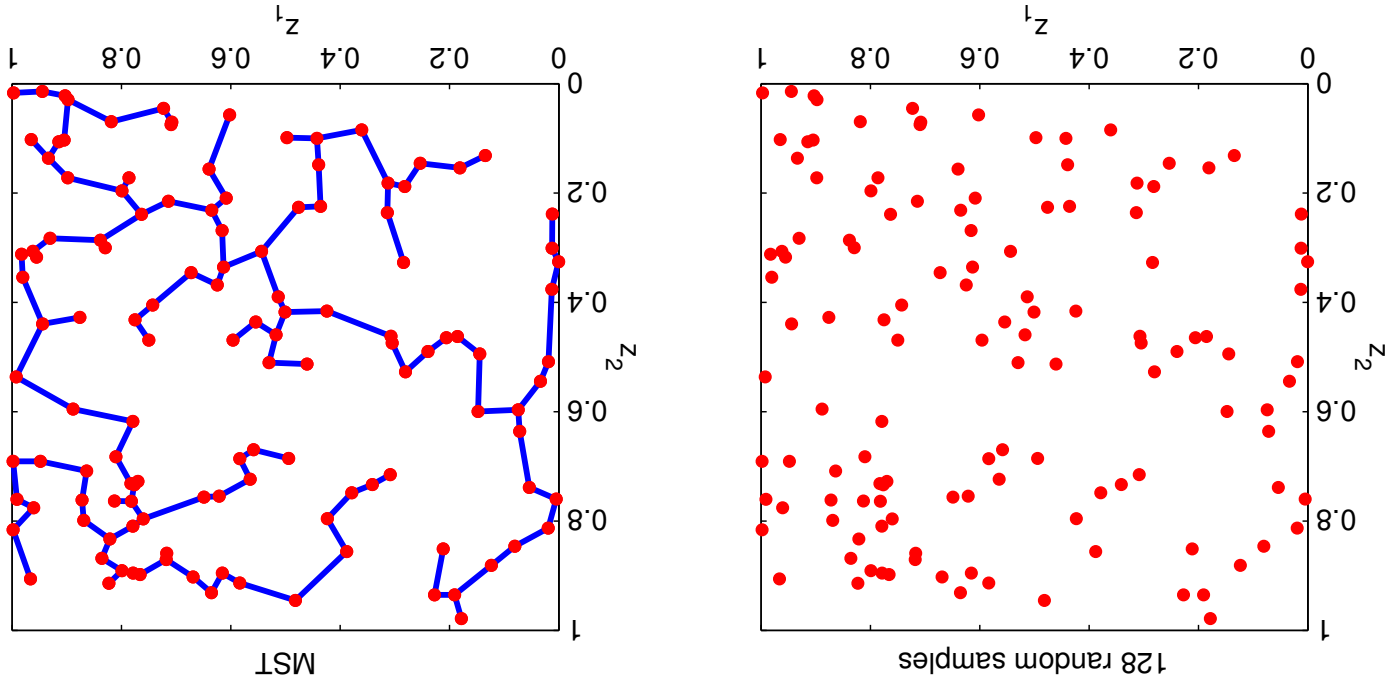
Some possibilities:

1. Assume parametric models for $f, g, f_{X,Y}$
(Vasconcelos&Lipman:2000,Stoica&etal:1998)
2. Substitute non-parametric density estimates of $f, g, f_{X,Y}$
 - (a) Quantize feature space and use histogram estimates
(Beirlant&etal:1997)
 - (b) Use adaptive partitioning density estimates (Vasicek:1976, Miller:2002, Gray&etal:2000)
3. Use “entropic graphs” which emulate/estimate \mathcal{Q}
(Hero&Michel:1997,Nemwuchwala,Hero&Carson:2002)

III. Entropic Euclidean Graphs

1. The minimal spanning tree (MST)
2. The k-nearest neighbor (k-NN) graph
3. Asymptotic trends

A Set of Feature Samples and a Euclidean Spanning Graph

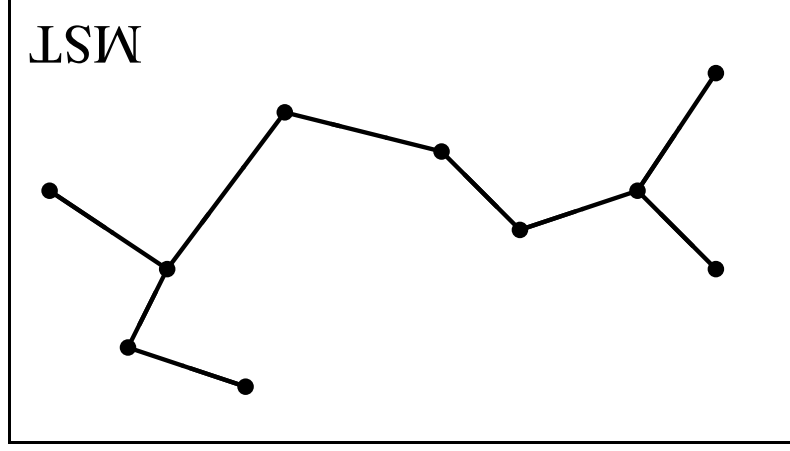


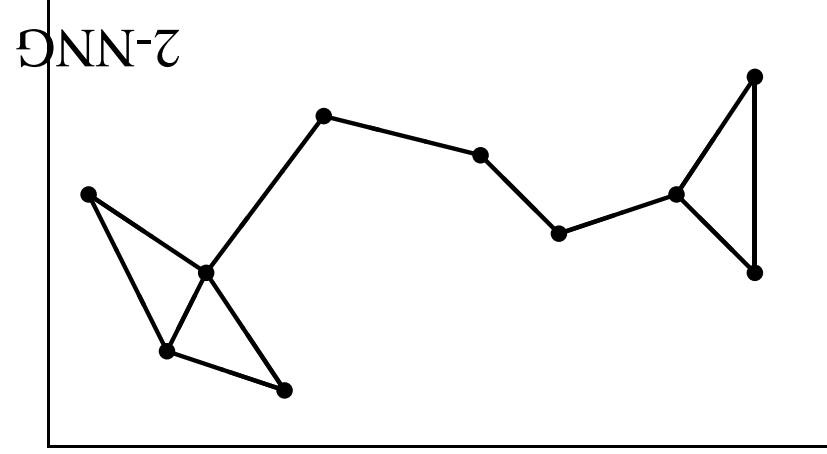
Minimal Euclidean Graphs: MST

Let $T_n = \mathcal{T}(X_n)$ denote the possible sets of edges in the class of acyclic graphs spanning X_n (spanning trees).

The Euclidean Power Weighted MST achieves

$$L_{\text{MST}}^\gamma(X_n) = \min_{T_n \in \mathcal{T}_n} \sum_{e \in T_n} \|e\|^\gamma.$$





$$L_{k-NG}^{\gamma}(X_n) = \sum_n \min_{\mathcal{N}_{k,i}(X_n)} \sum_{e \in \mathcal{N}_{k,i}(X_n)} |e|^{\gamma}$$

The Euclidean Power Weighted k - NG is

other points in X_n .

Let $\mathcal{N}_{k,i}(X_n)$ denote the possible sets of k edges connecting point x_i to all

Minimal Euclidean graphs: k - NG

MST for Two Different Samples

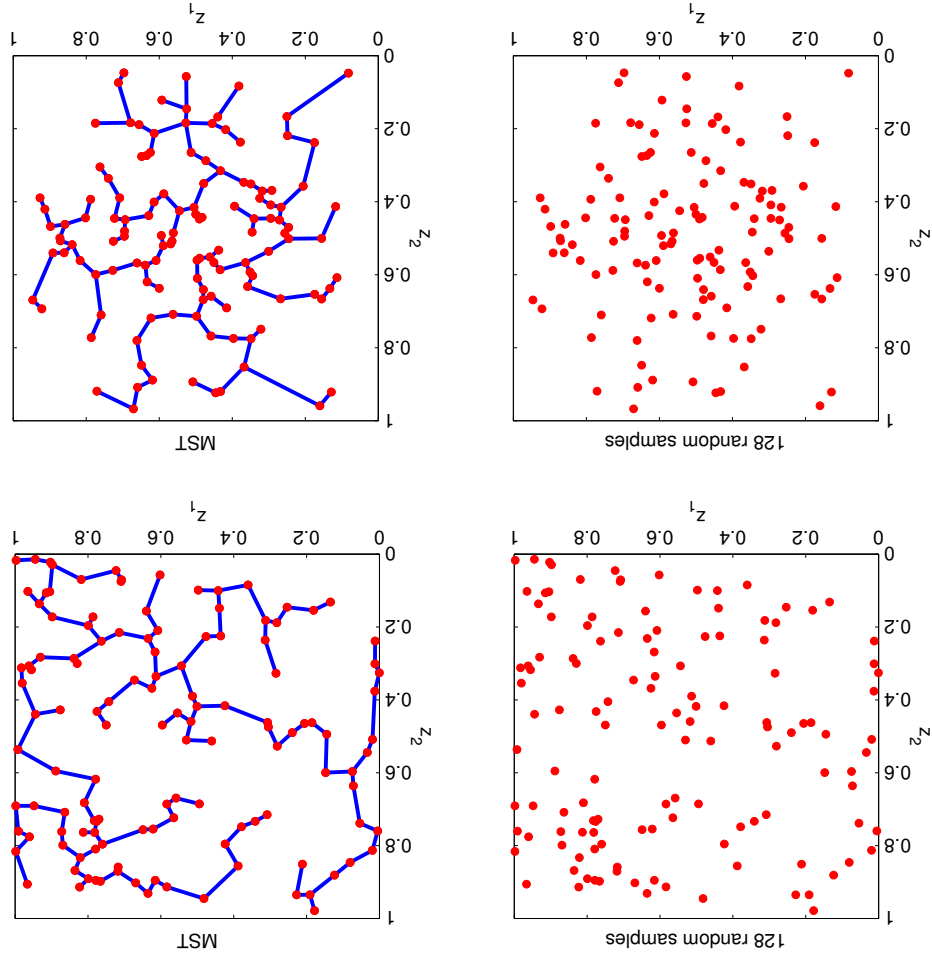


Figure 5:

Large n behavior of MST

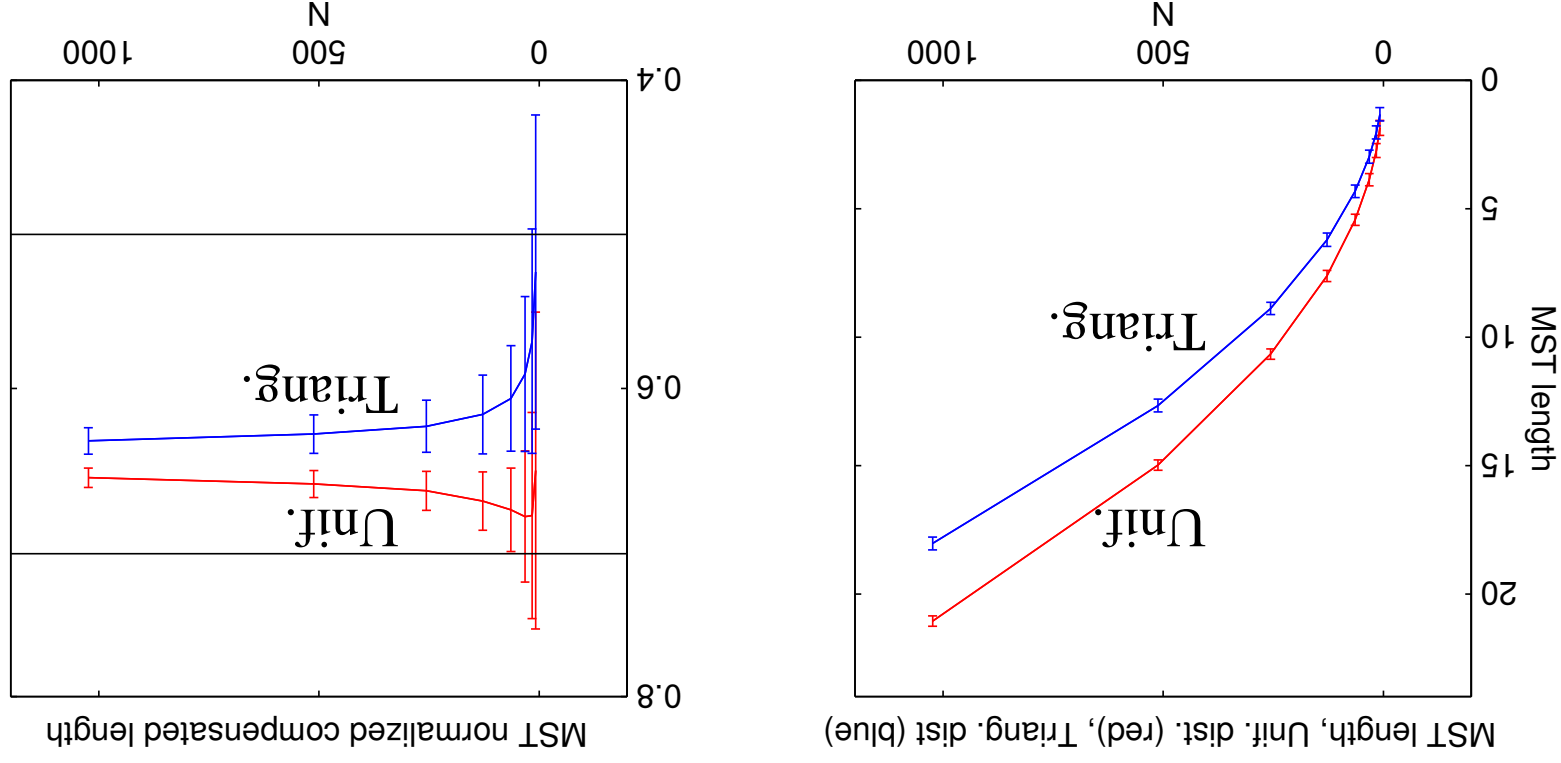


Figure: MST and log MST weights as function of the number of samples.

IV. Entropic Euclidean Graph Theory

1. The Beardwood, Halton Hammersley Theorem
2. Extension to divergence estimation
3. Extension to greedy algorithms
4. Extension to K-MST

Asymptotics: the BHH Theorem

Define the MST length functional

$$L_\gamma(X_n) = \min_{T_n} \sum_{e \in T_n} \|e\|_\gamma.$$

Theorem 1 [Beardwood, Halton&Hammerley:1959] Let

$X_n = \{X_1, \dots, X_n\}$ be an i.i.d. realization from a Lebesgue density f with support $S \subset [0, 1]^d$.

$$\lim_{n \rightarrow \infty} L_\gamma(X_n)/n^{d/(d-\gamma)} = \beta_{L_\gamma, d} \int_S f(x)^{d/(d-\gamma)} dx, \quad (a.s.)$$

Or, letting $\alpha = d - \gamma/d$

$$\frac{1}{1-\alpha} \ln(L_\gamma(X_n)/n^\alpha) \rightarrow H^\alpha(f) + c \quad (a.s.)$$

Rényi Entropy and Divergence

- Rényi Entropy of order α [Rényi:61,70]

$$H_\alpha(f) = \frac{1}{1-\alpha} \ln \int_S f^\alpha(x) dx$$

- Rényi α -divergence of fractional order $\alpha \in [0, 1[$

$$D_\alpha(f_1 \parallel f_0) = \int_S \ln \frac{1}{1-\alpha} \frac{f_1^\alpha}{f_0^\alpha} dx = \int_S \ln \frac{1}{1-\alpha} f_1^\alpha f_0^{1-\alpha} dx$$

– α -Divergence vs. Kullback-Liebler divergence

$$\lim_{\alpha \rightarrow 1} D_\alpha(f_1 \parallel f_0) = \int f_1 \ln \frac{f_1}{f_0} dx.$$

α -Divergence and Decision Theoretic Error Exponents

Let Z_i be i.i.d.:

$$H_0 : Z_i \sim f$$

$$H_1 : Z_i \sim g$$

Bayes probability of error

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

Sanov bound (Blahut:1987, Dembo&Zeitouni:98)

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_F(n) = - \sup_{\alpha \in [0,1]} \{ (1 - \alpha) D^\alpha(g \| f) \}$$

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_M(n) = - \sup_{\alpha \in [0,1]} \{ (1 - \alpha) D^\alpha(f \| g) \}.$$

Entropic Graphs for Clustering and Outlier Rejection: k-MST

Assume f is a mixture density of the form

$$f = (1 - \epsilon)f_1 + \epsilon f_0,$$

where

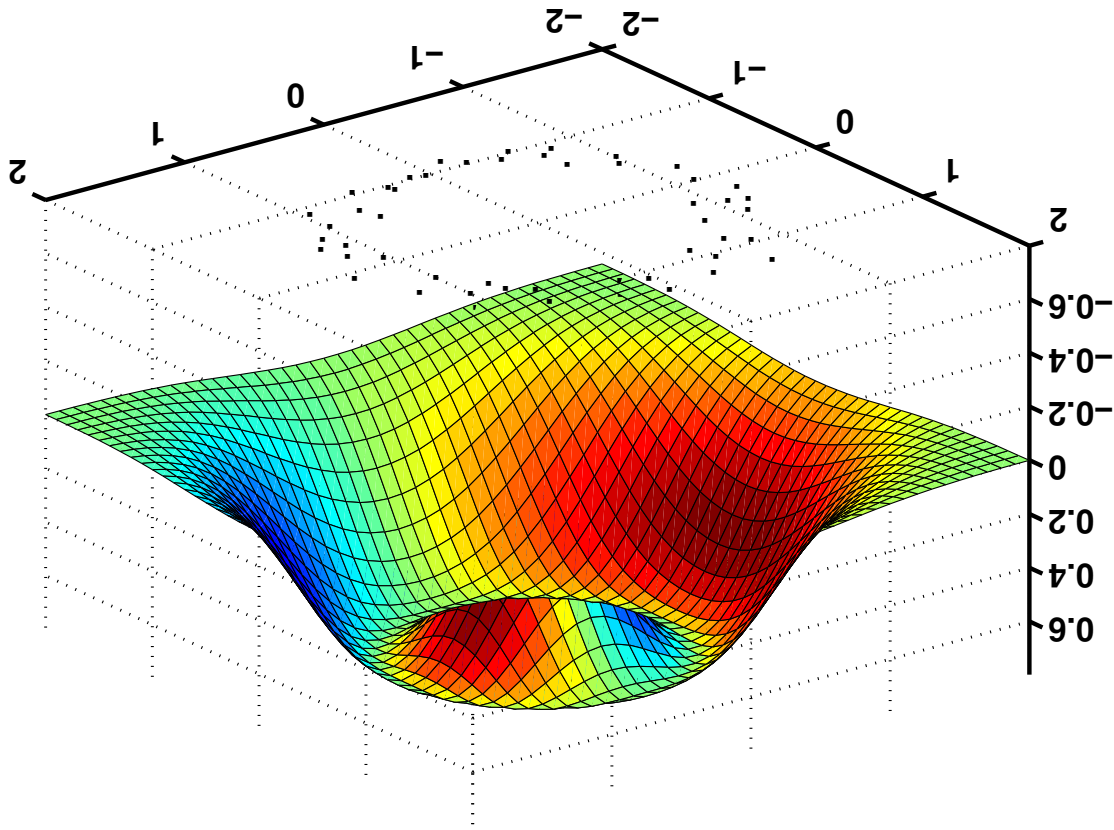
- f_0 is a known "outlier" density
- f_1 is an unknown target density
- $\epsilon \in [0, 1]$ is unknown mixture parameter

Objective: given realization X_n from f cluster the realizations from f_1 .

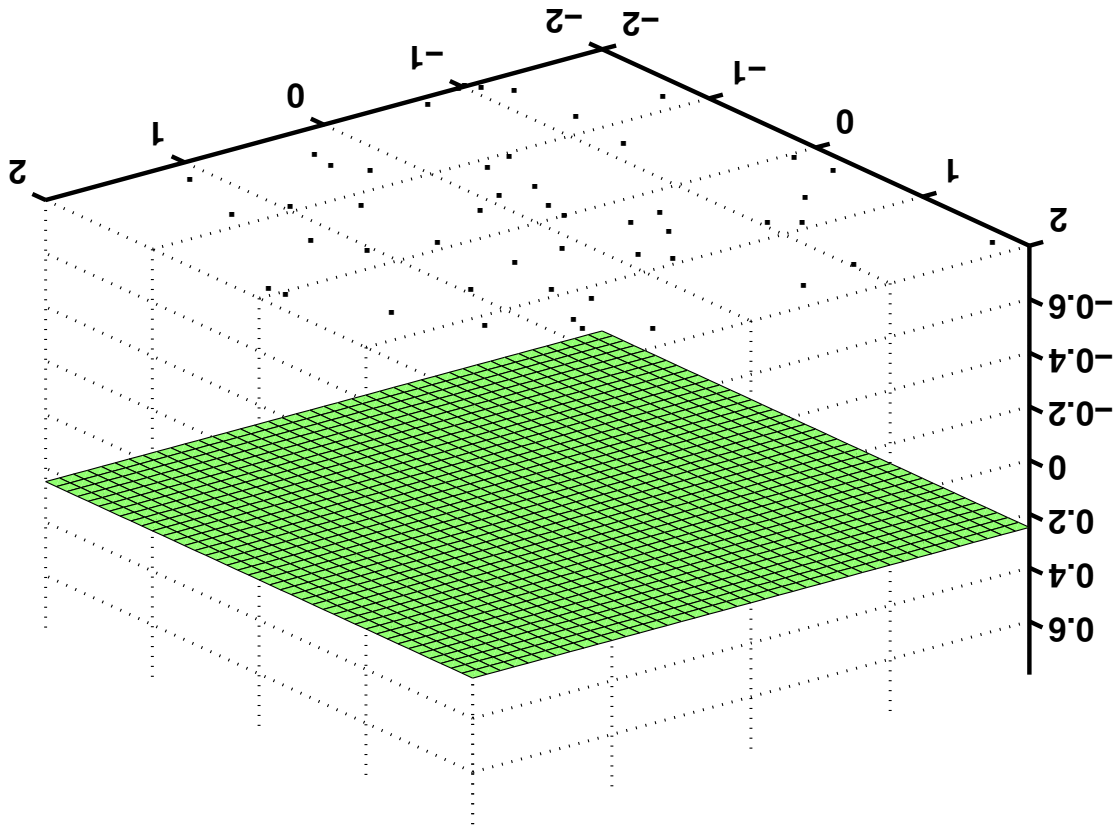
Two-step k-MST procedure:

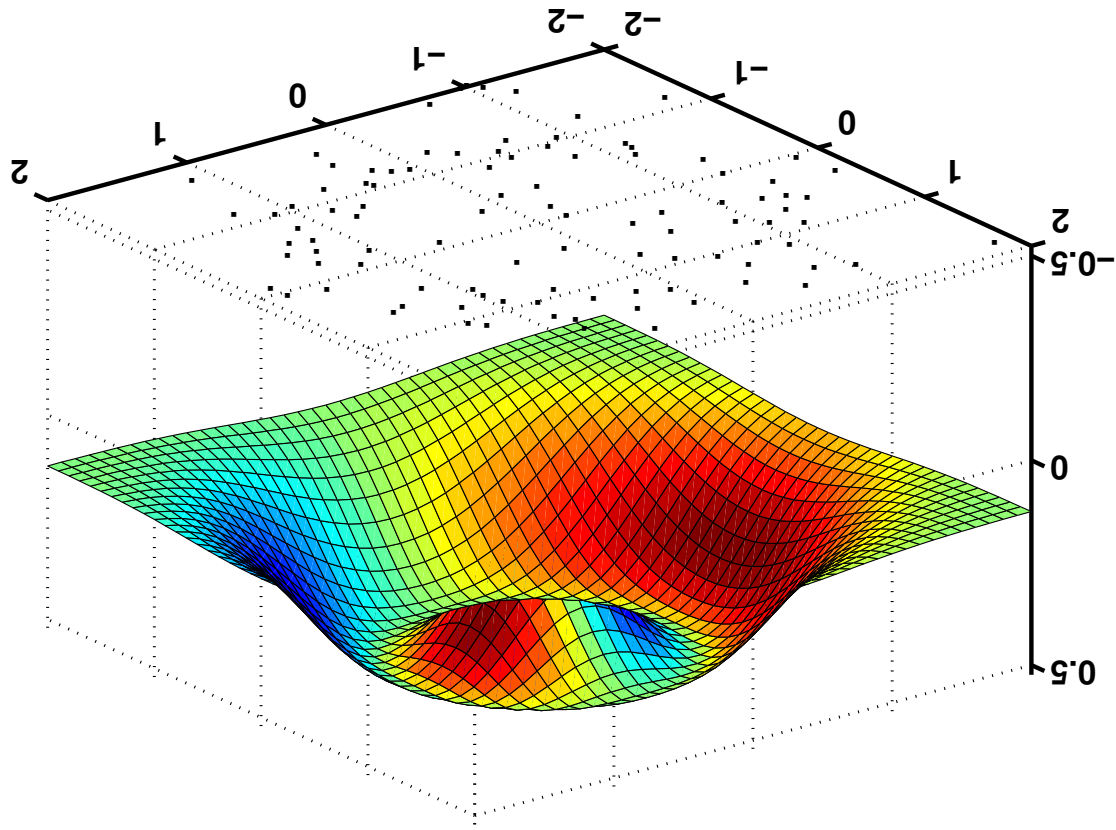
1. Convert f_0 to maxent (uniform) density via measure transformation
2. "Prune" the MST on transformed X_n to eliminate vertices arising from maxent density

Example: Annulus Target Density f_1



Uniform Outlier Density f_o





Mixture Density

k -point Minimal Spanning Tree (k -MST)

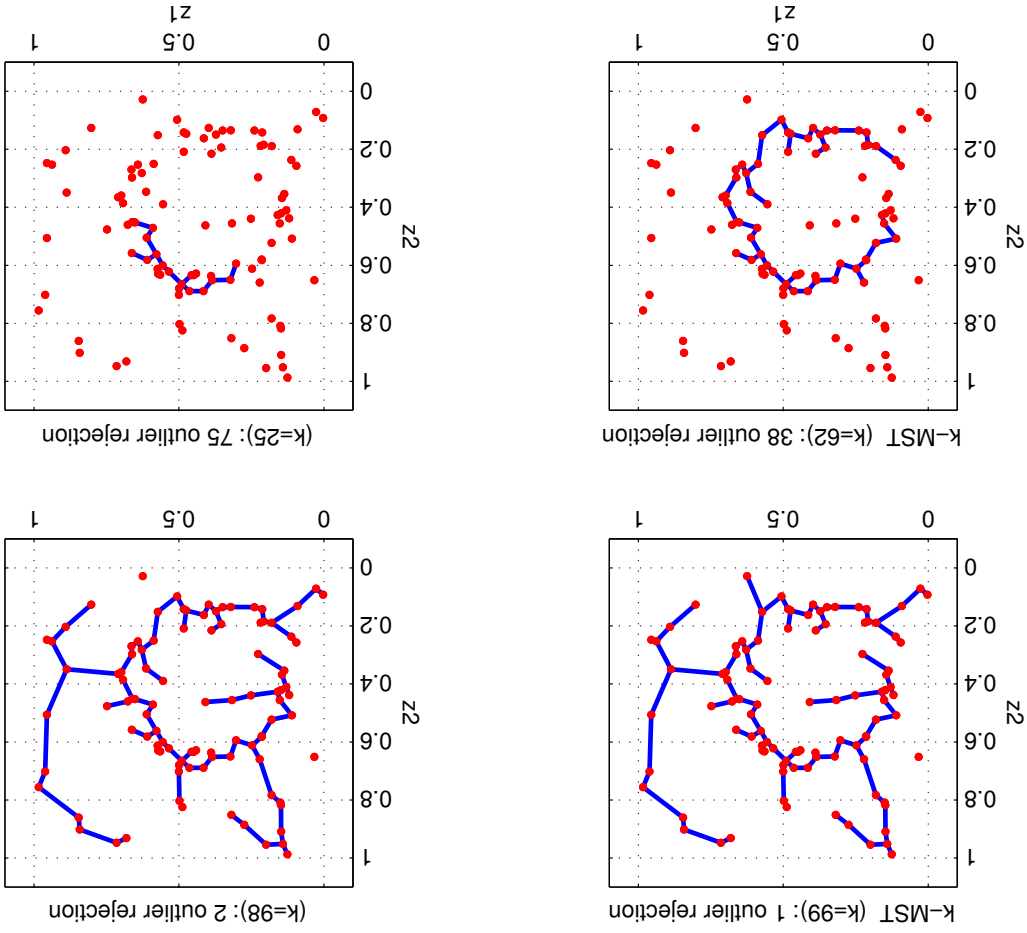


Figure 6: Clustering an annulus density from uniform noise via k -MST.

k-MST Stopping Rule (Hero&Michel:1997)

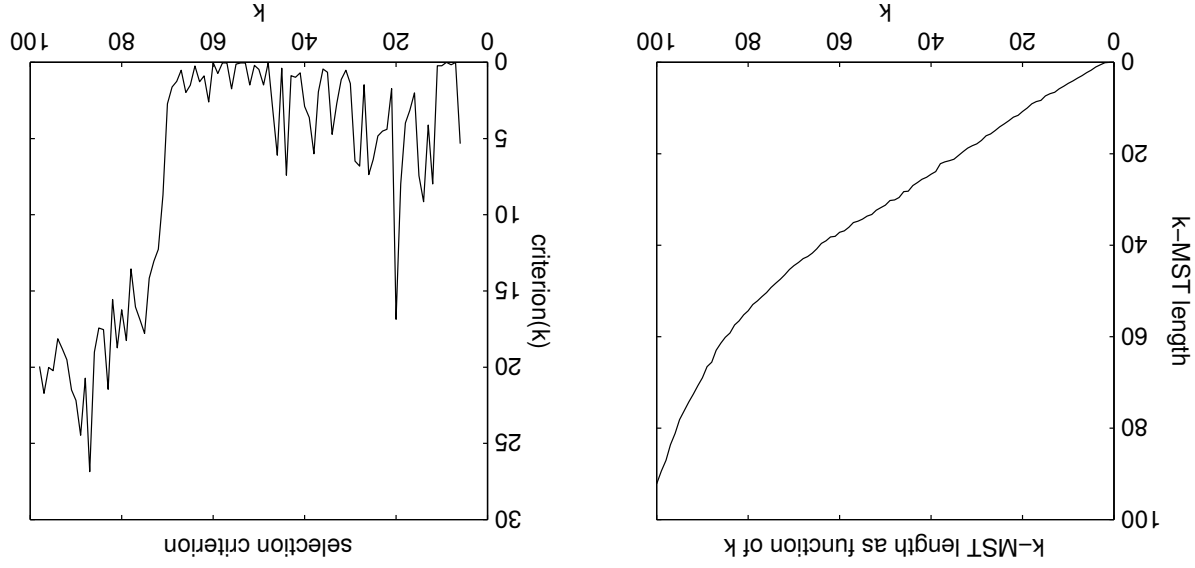


Figure 7: Left: k -MST curve for 2D annulus density with addition of uniform "outliers" has a knee in the vicinity of $n - k = 35$.

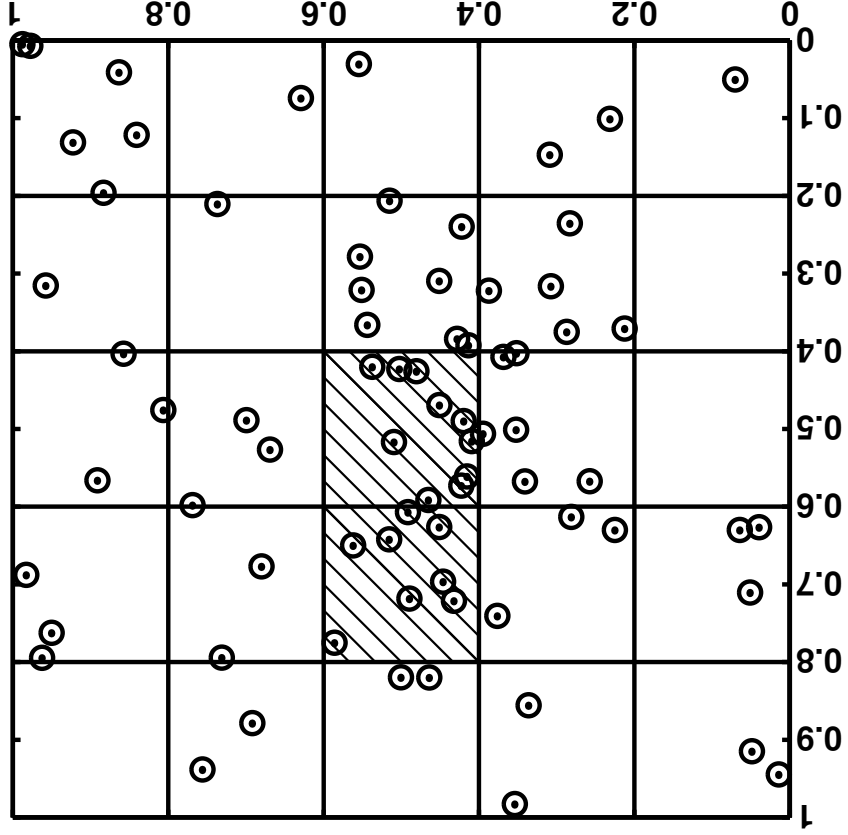


Figure 8: A smallest subset B_m^k is the union of the two cross hatched cells shown for the case of $m = 5$ and $k = 17$.

Fix $p \in [0, 1]$. If $k/n \rightarrow p$ then the length of the greedy partitioning k -MST satisfies [Heró&Miché:IT99]

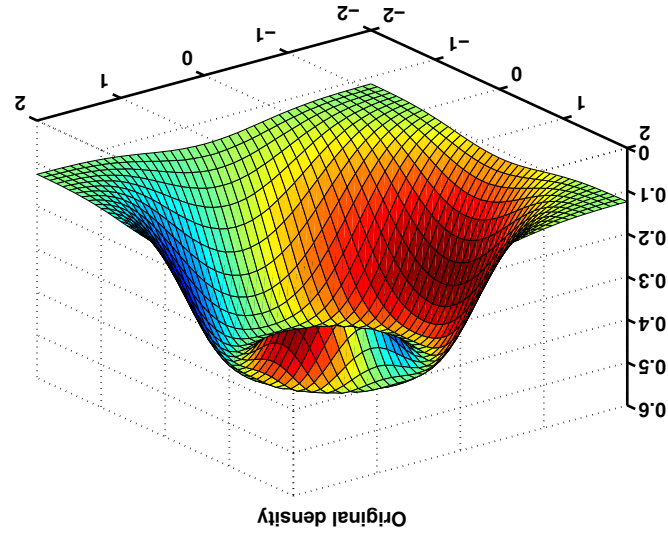
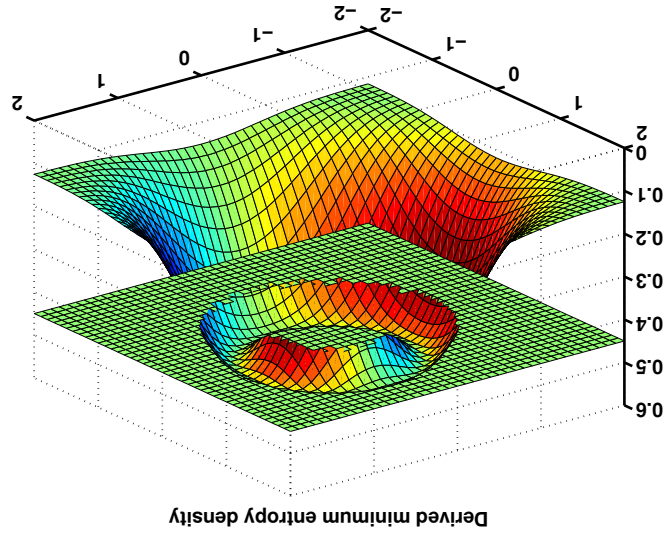
$$L_{\gamma}(X_{n,k}^*) / \lfloor pn \rfloor^{\alpha} \rightarrow \beta_{L_{\gamma,d}} \int_S f^{\alpha}(x|x \in A_0) dx \quad (a.s.)$$

where A_0 is level set of f which satisfies $\int_{A_0} f = p$. Alternatively, with

$$H_{\alpha}(f|x \in A_0) = \frac{1}{1-\alpha} \ln \int_S f^{\alpha}(x|x \in A_0) dx$$

$$\frac{1}{1-\alpha} \ln L_{\gamma}(X_{n,k}^*) / \lfloor pn \rfloor^{\alpha} \rightarrow \beta_{L_{\gamma,d}} H_{\alpha}(f|x \in A_0) + c \quad (a.s.)$$

Figure 9: Waterpouring construction of minimum entropy density.



k-MST Influence Function

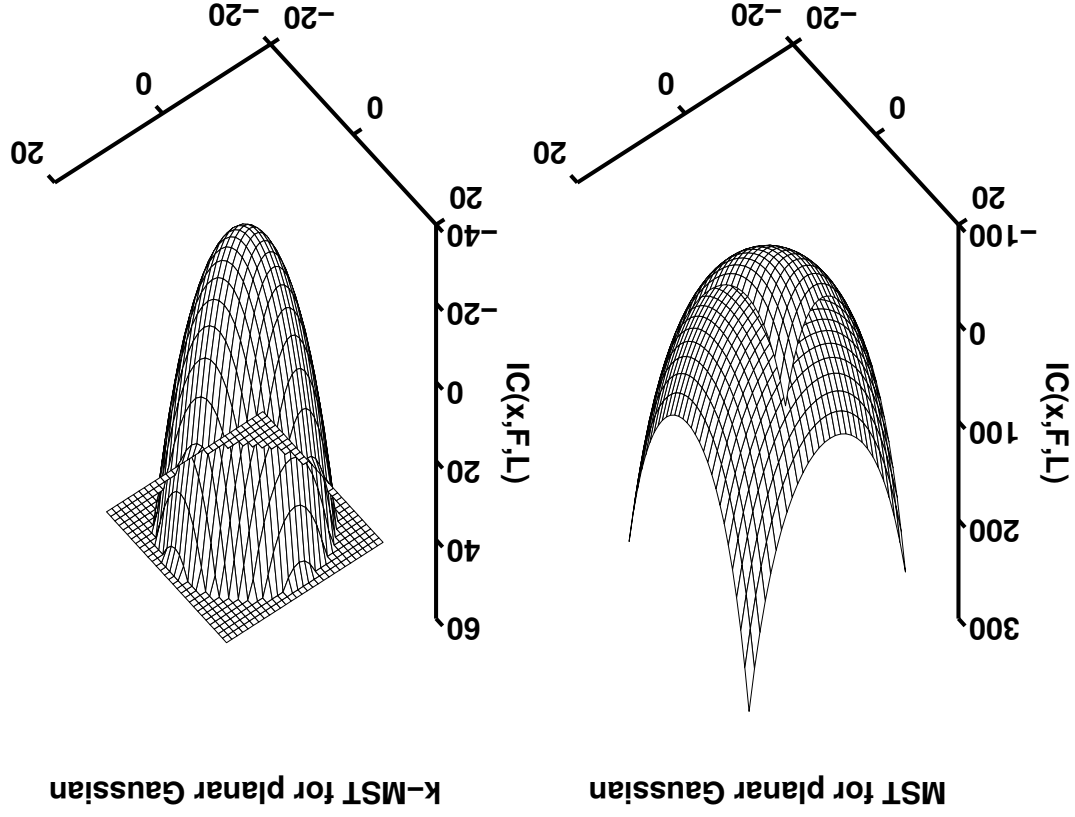
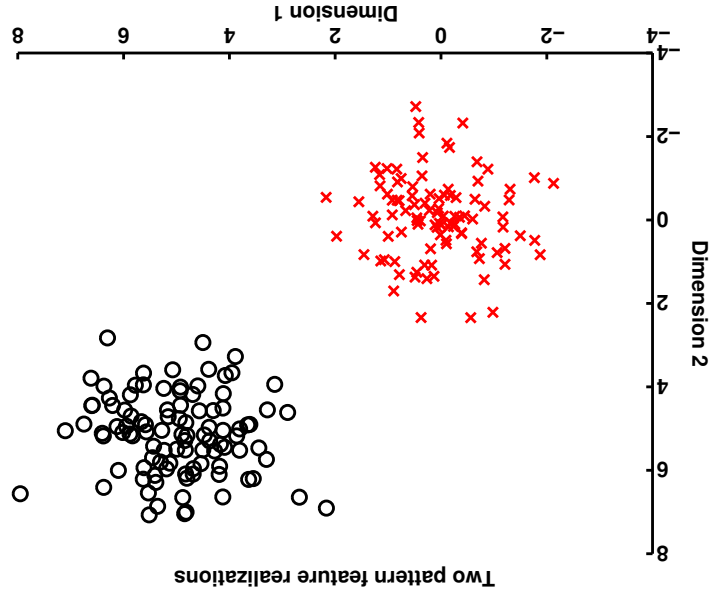


Figure 10: *MST and k-MST influence curves for Gaussian density on the plane.*

V. Entropic Graphs for Pattern Matching



Two groups of i.i.d. feature realizations on $[0, 1]^d$:

- $\mathcal{X}^m = \{X_1, \dots, X_m\}, X_i \sim f$
- $\mathcal{Y}^n = \{Y_1, \dots, Y_n\}, Y_i \sim g$
- $d - 1 = b, (u + m)/m = d$

Objective: estimate separation of f and g using \mathcal{X}_m and \mathcal{Y}_n

Some entropic graph estimation possibilities

Option 1. construct MST/k-NG on pooled data $\mathcal{X}_m \cup \mathcal{Y}_n$

(Hero, Ma, Michel&Gorman:2001):

$$\ln L_\gamma(\mathcal{X}_m \cup \mathcal{Y}_n) / N^\alpha \leftarrow (1 - \alpha) H_\alpha(p f + q g) + c, \quad (a.s.)$$

If subsequently subtract $\ln L(\mathcal{X}_m) / N^\alpha$ and $\ln L(\mathcal{Y}_n) / N^\alpha$ obtain estimator of α -Jensen difference (Basseville:1989, He&etal:2001)

$$\Delta(f, g) = H_\alpha(p f + q g) - p H_\alpha(f) - q H_\alpha(g)$$

Option 2: prune all single-class connections from pooled MST and compute normalized length

$$L_\gamma(\mathcal{X}_m \Delta \mathcal{Y}_n) = \frac{1}{N^\alpha} \sum_{e_{xy} \in \mathcal{E}} |e_{xy}|_\gamma$$

- for $\gamma = 0$ obtain "Multivariate runs statistic" Friedman&Rafsky:1979 (FR).

- for $0 < \gamma < d$ obtain generalized FR statistic (Costa&Hero:2003)

- FR($\gamma = 0$) statistic converges a.s. to affinity (Henze&Penrose:1998)

$$A_{FR}(f, g) = 2pq \int \frac{f(x)g(x)}{pf(x) + qg(x)} dx$$

This affinity is related to divergence measure:

$$D_{FR}(f \| g) = 1 - A_{FR}(f, g) = \int \frac{p^2 f^2(x) + q^2 g^2(x)}{pf(x) + qg(x)} dx$$

Option 3: implement entropic graph approximation of adaptive partition estimators of different divergence functionals (example below).

Illustration: Jensen Difference estimator

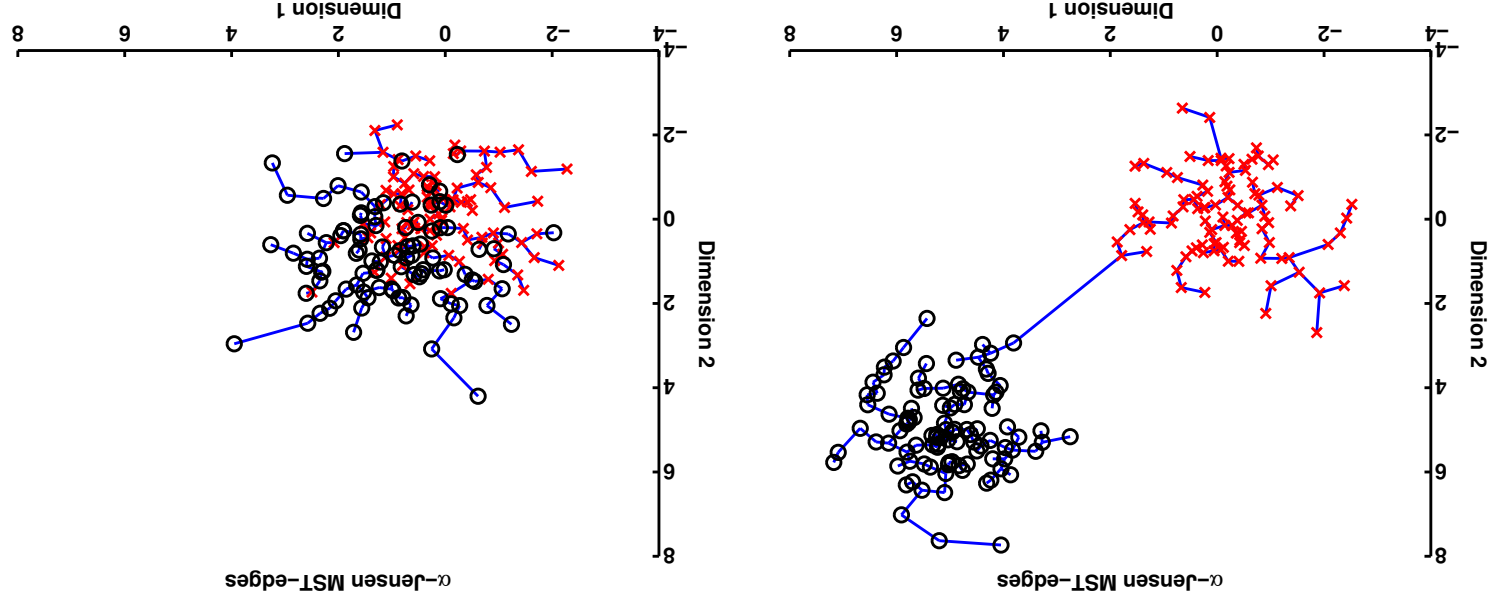
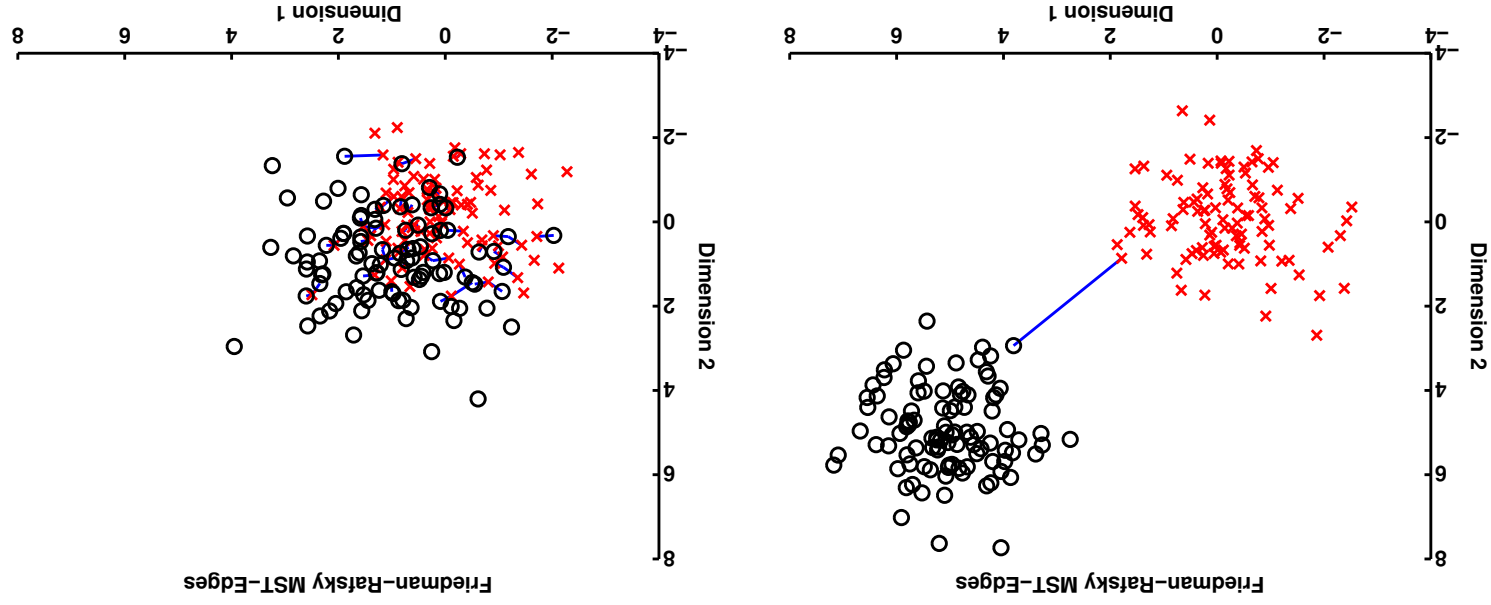


Illustration: Friedman-Ratsky Statistic



Define

$$I(f) = \int_S f(x) dx$$

For N i.i.d. realizations $\{x_i\}_{i=1}^N$ from f define:

1. Π : an M -cell partition of $[0, 1]^d$.

2. $\Pi(x)$: the cell in Π containing point $x \in [0, 1]^d$

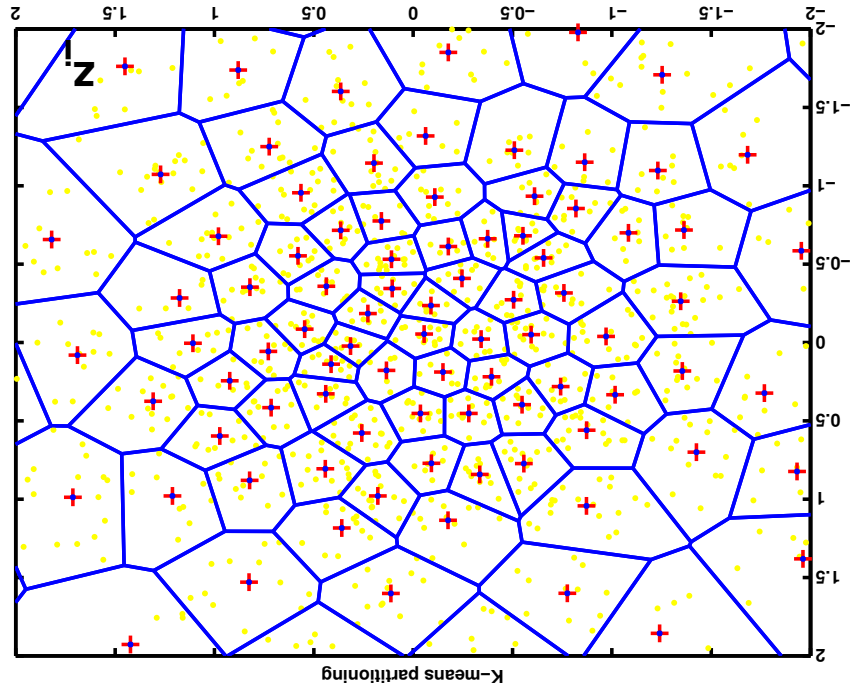
3. \hat{f}_Π : a partition estimator of f

$$\hat{f}_\Pi(x) = \frac{\lambda(\Pi(x))}{n(\Pi(x))}, \quad x \in [0, 1]^d$$

$$\hat{\Pi}_\alpha = \sum_{i=1}^N \frac{1}{N} f_{\Pi}^{\alpha-1}(z_i) = \sum_{i=1}^N \frac{1}{N} \frac{\chi(\Pi(z_i))}{n(\Pi(z_i))}$$

α -entropy estimator

For $\{z_i\}_{i=1}^N$ an i.i.d. realization **independent** of $\{x_i\}_{i=1}^N$ consider the



Under weak conditions on Π (Lugosi & Nobel: 1996), this estimator converges a.s. to

$$E[f^{\alpha-1}(z_i)] = \int_S f^\alpha(x) dx = I(f)$$

as $N \rightarrow \infty$. Equivalently,

$$\beta_{L_\gamma, d} I_\Pi \rightarrow \int_S f^\alpha(x) dx$$

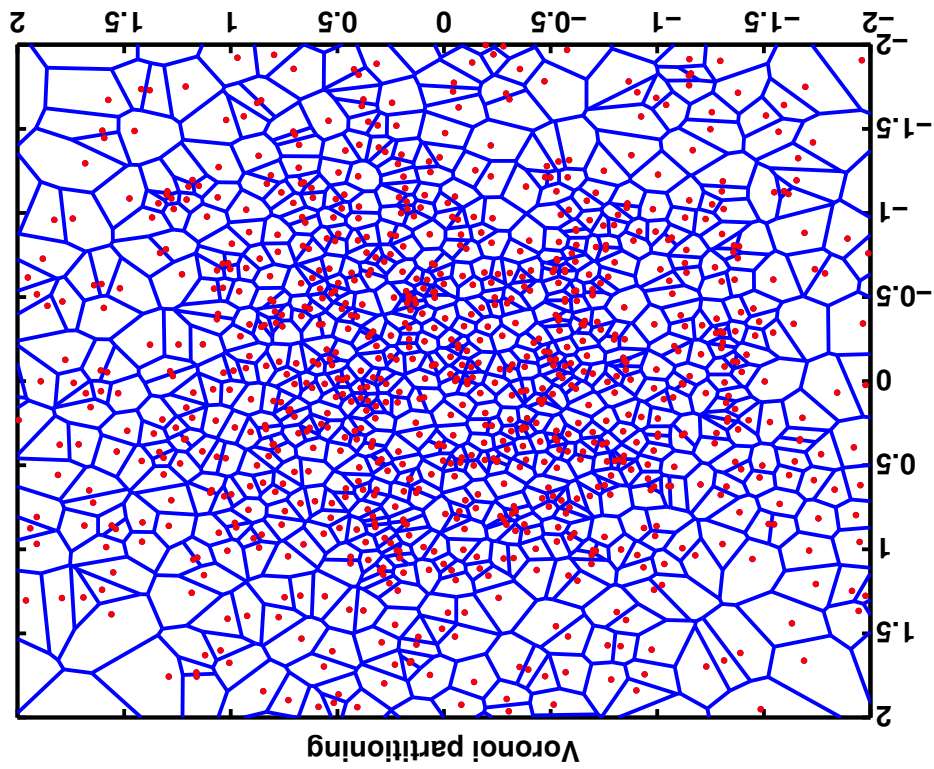
which corresponds to the a.s. limit of $L_\gamma(X_N)/N^\alpha$.

To exploit this correspondence, (formally) specialize to:

1. $\alpha = (d - \gamma)/d$
2. Π is Voronoi partition ($\mu(\Pi(x)) \equiv 1$)
3. $z_i = x_i, i = 1, \dots, N$

$$\beta_{L_{\gamma,p}}^{\lambda} \left(\frac{N}{p} \right) = \beta_{L_{\gamma,p}}^{\lambda} \left(\frac{N}{p} \right) = \beta_{L_{\gamma,p}}^{\lambda} \left(\frac{N}{p} \right)$$

In this case we have:



Q. What relation between $\lambda_{1/d}(\Pi(z_i))$ and e_i would make $\beta_{L\gamma,d} \Pi$ equal to $L_\gamma(X_N)/N\alpha$?

A. When

$$\beta_{L\gamma,d} \sum_{i=1}^l \frac{N}{1} = \left(\lambda_{1/d}(\Pi(z_i)) \right) \sum_{i=1}^l \frac{N}{\gamma}$$

which occurs if we identify

$$(1) \quad \lambda_{1/d} \Pi(z_i) = \frac{\beta_{L\gamma,d}}{n_{1/d}} e_i$$

Heuristic: can use formal relation (1) to obtain entropic graph implementations of divergence estimators.

$$A(f, g) = \int (p f(x) + q g(x))^\alpha (f p(x) g^q(x))^{1-\alpha} dx$$

1. Pooled sample $Z_{m+n} = X_m \cup \mathcal{G}_m$ has density $h = p f + q g$

2. Adaptive-partition plug-in estimator of $A(f, g)$ is

$$\hat{A}_{ap} = \frac{1}{N} \sum_{z_i=1}^{z_i=N} \left(\frac{\hat{f}_p(z_i) \hat{g}_q(z_i)}{h(z_i)} \right)^{1-\alpha} \rightarrow A(f, g) \text{ (a.s.)}$$

3. Specialize partition to Voronoi and substitute (1):

$$\hat{A}_{eg} = \frac{1}{N} \sum_{i=1}^N \min_{p \in R_\gamma(X_m \cup \mathcal{G}_n)} \left\{ \left(\frac{e^i(\mathcal{G}_n)}{e^i(X_m)} \right)^{p_\gamma}, \left(\frac{e^i(X_m)}{e^i(\mathcal{G}_n)} \right)^{q_\gamma} \right\}$$

Planar Pattern Matching Simulation

- X_m realization from $\mathcal{N}_2(\bar{0}, \mathbf{I})$
- \mathcal{Y}_n realization from $\mathcal{N}_2(\bar{D}, \mathbf{I})$
- Four pattern separation measures Δ investigated

$$\Delta = \frac{L_\gamma(X_m \cup \mathcal{Y}_n) / N^\alpha, L_0(X_m \Delta \mathcal{Y}_n) / N}{L_\gamma(X_m \Delta \mathcal{Y}_n) / N^\alpha, R_\gamma(X_m \cup \mathcal{Y}_n) / N}$$

- $\gamma = 1, \alpha = 1/2$

- Local resolution measure

$$p(\Delta) = \frac{|E[\Delta | D = 0] - E[\Delta | D = 1]|}{\sqrt{\sigma_z^2(D = 0) + \sigma_z^2(D = 1)}}$$

Pattern Matching Simulation: Convergence

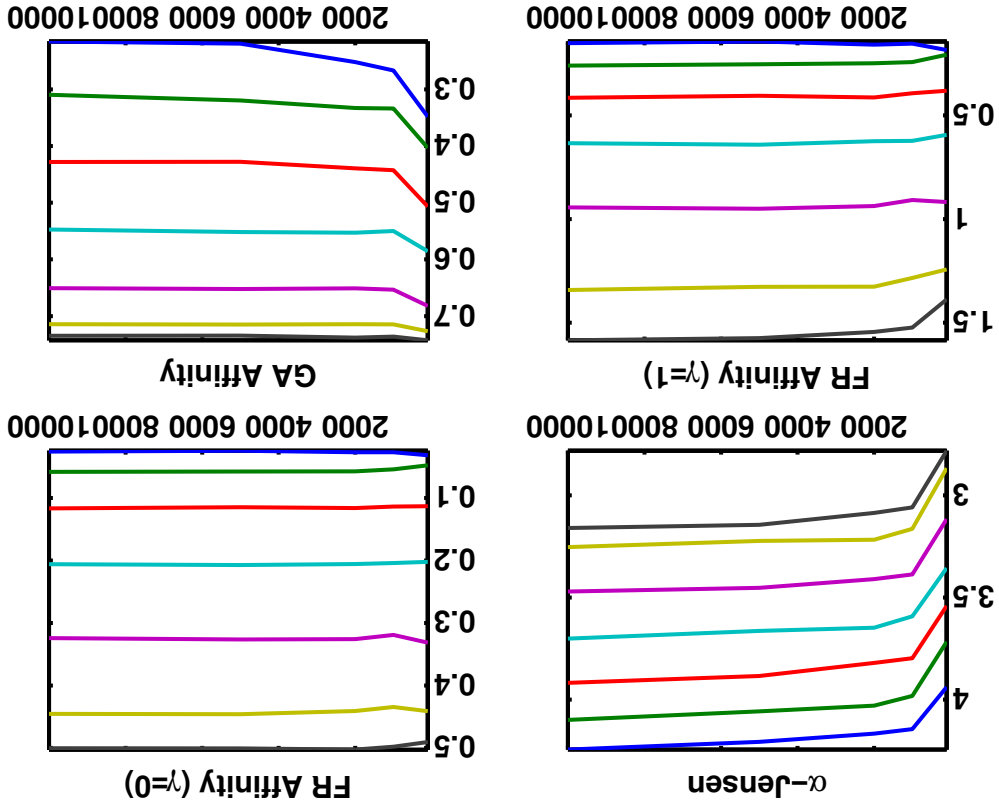
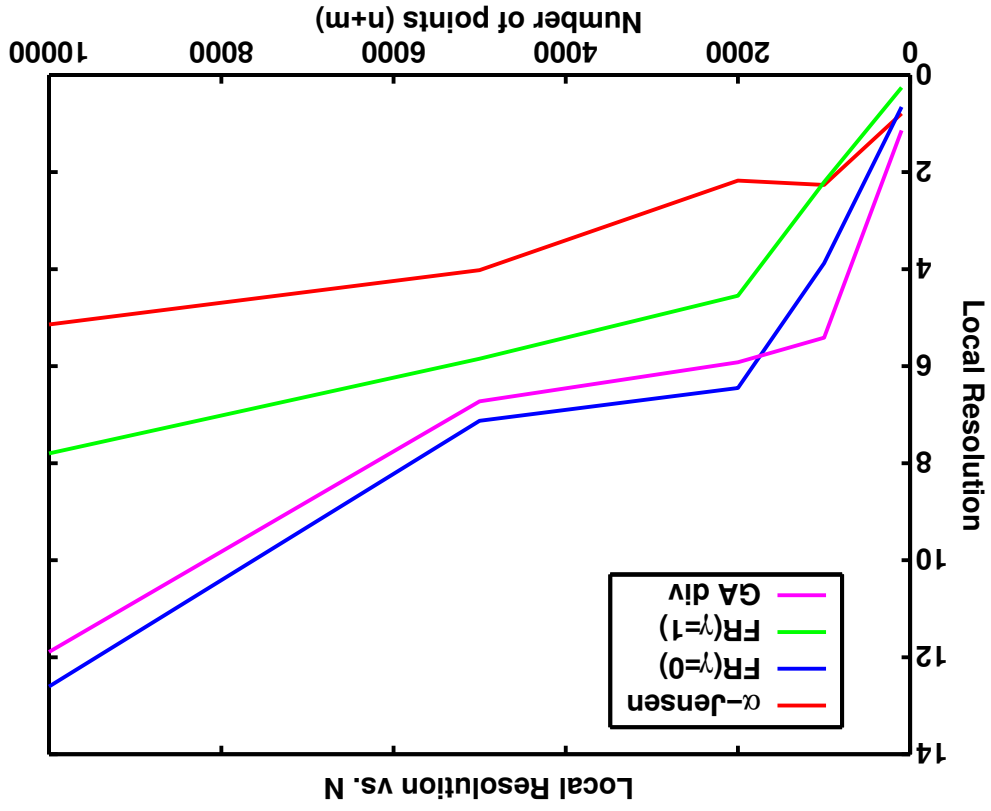


Figure 11: Convergence rates $1/\sqrt{N}$ (left) and $1/N$ (right)

Pattern Matching Simulation: Local Resolution



Conclusions

1. Entropic graphs can be used to estimate α -entropy and α -divergence
2. These methods can be applied to high dimensional feature-spaces
3. Clustering can be performed using entropic K -point graphs
4. Asymptotic theory can be used to motivate new entropic graph measures

References

1. "Robust entropy estimation strategies based on edge weighted random graphs," A.O. Hero and O. Michel, Proc. of Int. Soc. for Optical Engineering (SPIE) Symposium on Optical Science, San Diego, July 1998.
2. A. O. Hero, B. Ma, O. Michel and J. Gorman, "Applications of entropic spanning graphs," IEEE Signal Proc. Magazine (Special Issue on Mathematics in Imaging), Vol 19, No. 5, pp 85-95, Sept. 2002.
3. "Image registration using alpha-entropy measures and entropic graphs," H. Neemuchwala, A. O. Hero and P. L. Carson, European Journal of Signal processing, to appear Sept. 2003.
4. "Parametric and non-parametric approaches for multisensor data fusion," Bing Ma, PhD thesis in Dept. EECS, Univ. of Michigan, Jan. 2001.
5. "Asymptotic theory of greedy approximations to minimal K-point random graphs," A. O. Hero and O. Michel, IEEE Trans. on Information Theory, Vol. IT-45, pp. 1921-1939, Sept. 1999.

6. "Asymptotic rates of convergence of random minimal graphs," A. O. Hero, J. Costa and B. Ma, IEEE Trans. on Inform. Theory, Submitted Aug. 2001.
7. O. Michel and A. O. Hero "Entropic graph applications", Proc. XI European Signal Processing Conference, Toulouse France, Sept 2002.
8. "Estimation of Rényi Information Divergence via Pruned Minimal Spanning Trees," A. O. Hero and O. Michel, Proc. of 1999 IEEE Workshop on Higher Order Statistics, Caesaria Israel, June 1999.

Extension of BHH to Divergence Estimation?

Question: How to generalize entropic graph estimates of

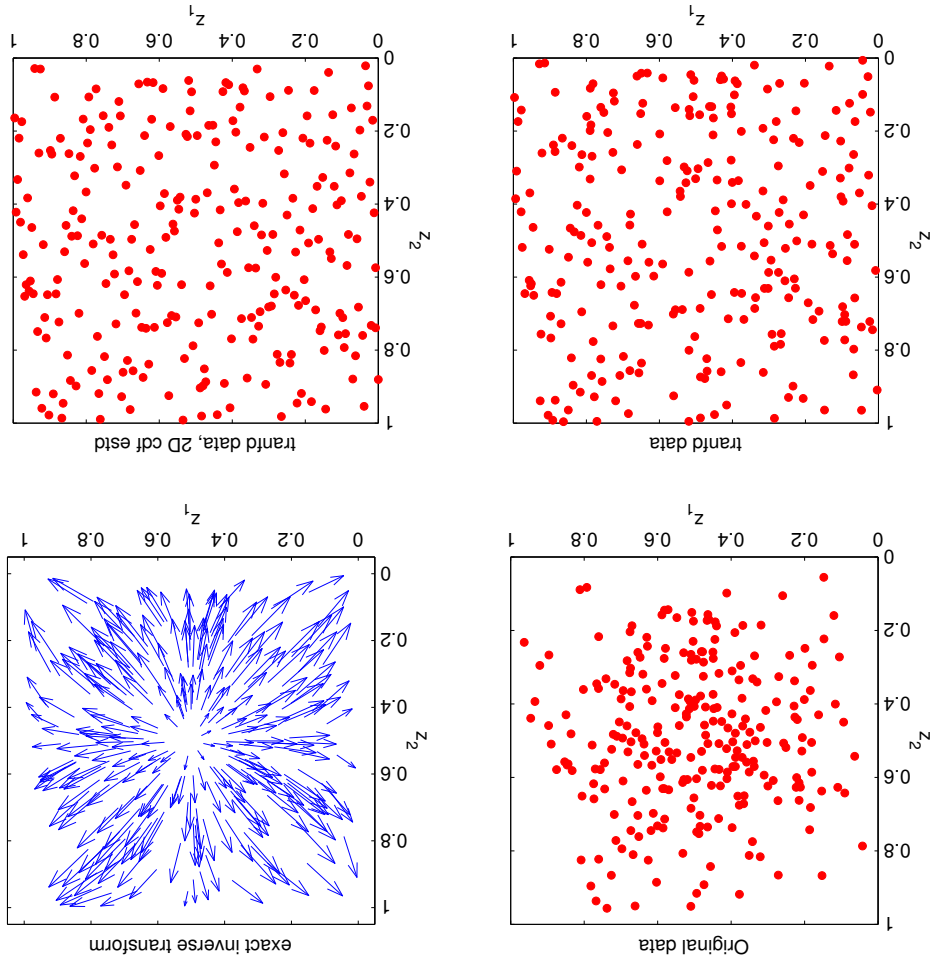
$$\frac{1}{1-\alpha} \ln \int f^\alpha(x) dx \quad \text{to} \quad \frac{1}{1-\alpha} \ln \int f^\alpha(x) g^{1-\alpha}(x) dx ?$$

One possibility:

- $g(x)$: a **known** reference density on $[0, 1]^d$
- Assume $f \ll g$, i.e. for all x such that $g(x) = 0$ we have $f(x) = 0$.
- Make measure transformation $M(x)$ such that $dx \mapsto g(x)dx$ on $[0, 1]^d$.
Then for $\mathcal{G}^n = M(\mathcal{X}^n)$

$$L_{\mathcal{G}^n} / n^\alpha \leftarrow \beta_{L_{\mathcal{G}^n, d}} \int \left(\frac{g(x)}{f(x)} \right)^\alpha g(x) dx, \quad (\text{a.s.})$$

Figure 12: Top Left: i.i.d. sample from triangular distribution, Top Right: exact transformation, Bottom: after application of exact and empirical transformations.



What is the entropic graph's convergence rate?

Theorem 2 (Hero, Costa&Ma:2001) Let $d \geq 2$ and $1 \leq \gamma \leq d - 1$.

Assume X_1, \dots, X_n are i.i.d. random vectors over $[0, 1]^d$ with density $f \in \Sigma_d(\beta, l)$, $\beta, l > 0$, having support $S \subset [0, 1]^d$. Assume also that $f^{\frac{1}{2} - \frac{d}{\gamma}}$ is integrable. Then,

$$O\left(n^{-r_1(d, \beta)}\right) \leq$$

$$E \sup_{f \in \Sigma_d(\beta, l)} \left[L^\gamma(X_1, \dots, X_n)^{d/\lambda - d} - \beta^{L^\gamma, d} \int_S f^{(d-\lambda)/d} dx \right]^{1/d}$$

$$\leq O\left(n^{-r_2(d, \beta)}\right),$$

where

$$r_1(d, \beta) = \min\left\{ \frac{4\beta}{4\beta + d}, 1/2 \right\} \quad r_2(d, \beta) = \frac{\alpha\beta}{\alpha\beta + 1} \frac{d}{1}$$

and $\alpha = \frac{d}{d-\gamma}$.

Extension to Partition Approximations

$$L_m^\lambda(X_n) = \sum_{i=1}^m L_i^\lambda(X_n \cup \mathcal{O}^i) + b(m),$$

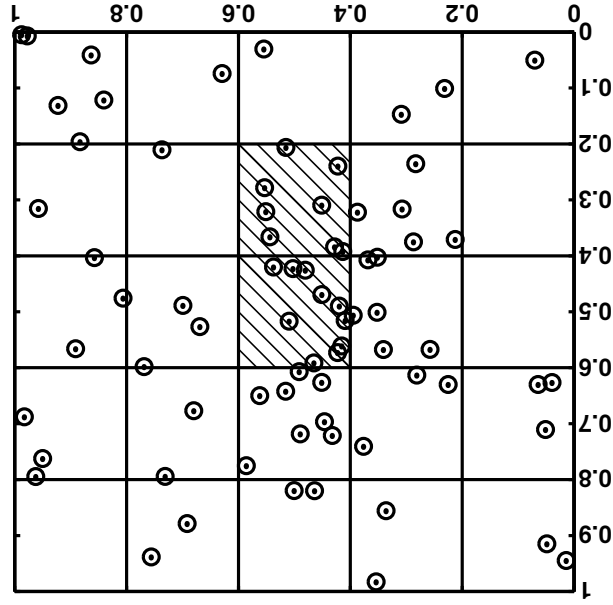


Figure 13: *Partition approximation.*

This bound is attained by choosing the progressive-resolution sequence

$$m = m(n) = \lfloor (1 + \beta \alpha)^{\frac{\lambda}{1-p}} \rfloor.$$

$$r_{\beta, d} = \frac{\alpha \beta \frac{\lambda}{1-p}}{\alpha \beta + 1} \cdot \frac{1}{p}.$$

where

$$\left(n^{-r_{\beta, d}} \right) \leq O,$$

$$\sup_{f \in \Sigma_{\beta, d}^p} E \left[\left| L_m^\gamma(X_1, \dots, X_n) - \int_S f(x) dx \right| \right] \leq O \left(n^{-r_{\beta, d}} \right)$$

$$\leq O \left(n^{-r_{\beta, d}} \right)$$

proposition, if $b(m) = O(m^{-\gamma})$

approximation to $L_m^\gamma(X_n)$. Under the same hypotheses as in the previous

Theorem 3 (Hero, Costa & Ma:2001) Let $L_m^\gamma(X_n)$ be a partition