

Introduction

- ▶ For a fixed $\alpha \in [0, 1) \cup (1, \infty)$, we are interested in estimating the Rényi- α divergence

$$D_\alpha(\mathbf{p} \parallel \mathbf{q}) = \frac{1}{1-\alpha} \log \int_{\mathcal{X}} \mathbf{p}^\alpha(\mathbf{x}) \mathbf{q}^{1-\alpha}(\mathbf{x}) d\mathbf{x},$$

between two unknown, continuous, nonparametric probability densities \mathbf{p} and \mathbf{q} over $\mathcal{X} \subseteq \mathbb{R}^d$, using samples from each density.

- ▶ Applications of divergence estimation include
 - ▷ extending machine learning algorithms designed to operate on finite-dimensional feature vectors to the setting where inputs are sets or distributions.
 - ▷ estimating entropy and mutual information.
- ▶ Rényi- α Divergence has KL-Divergence as its $\alpha \rightarrow 1$ limit case, and is related to Tsallis- α , Jensen-Shannon, and Hellinger divergences.
- ▶ Few divergence estimators have known rates, and, to the best of our knowledge, none have known exponential concentration bounds.
- ▶ We propose and analyze a plug-in estimator based on kernel density estimation, for densities on the unit cube $\mathcal{X} = [0, 1]^d$. We prove
 - ▷ the estimator is exponentially concentrated about its mean.
 - ▷ for densities in a β -Hölder smoothness class with certain boundary conditions, the bias of the estimator is bounded by $\mathcal{O}\left(n^{-\frac{\beta}{\beta+d}}\right)$, where n is the number of samples from each density.

Assumptions

Let $\beta > 0$, and let $\ell := \lfloor \beta \rfloor$ be the greatest integer *strictly* less than β . We make the following four assumptions on the densities \mathbf{p} and \mathbf{q} , and the kernel \mathbf{K} :

- ▶ **(Smoothness)** All (mixed) ℓ -order partial derivatives of \mathbf{p} and \mathbf{q} exist and are $(\beta - \ell)$ -Hölder Continuous (i.e., $\exists L \in \mathbb{R}$ such that, $\forall \mathbf{x}, \mathbf{x} + \mathbf{v} \in \mathcal{X}$, $|\mathbf{i}| = \ell$,

$$|\mathbf{D}^{\mathbf{i}}\mathbf{p}(\mathbf{x} + \mathbf{v}) - \mathbf{D}^{\mathbf{i}}\mathbf{p}(\mathbf{x})|, |\mathbf{D}^{\mathbf{i}}\mathbf{q}(\mathbf{x} + \mathbf{v}) - \mathbf{D}^{\mathbf{i}}\mathbf{q}(\mathbf{x})| \leq L \|\mathbf{v}\|_2^{\beta-\ell}.$$
- ▶ **(Boundedness)** $\exists \kappa_1, \kappa_2 \in \mathbb{R}$ such that, $\forall \mathbf{x} \in \mathcal{X}$,

$$0 < \kappa_1 \leq \mathbf{p}(\mathbf{x}), \mathbf{q}(\mathbf{x}) \leq \kappa_2 < +\infty.$$
- ▶ **(Boundary)** All derivatives of \mathbf{p} and \mathbf{q} vanish at the boundary

$$\partial \mathcal{X} = \{\mathbf{x} \in \mathcal{X} : \mathbf{x}_i \in \{0, 1\} \text{ for some } i \in [d]\}$$
 (i.e., $\sup_{1 \leq i \leq \ell} |\mathbf{D}^{\mathbf{i}}(\mathbf{x})| \rightarrow 0$ as $\text{dist}(\mathbf{x}, \partial \mathcal{X}) \rightarrow 0$).
- ▶ **(Kernel)** The kernel $\mathbf{K} : \mathbb{R} \rightarrow \mathbb{R}$ has support in $[-1, 1]$,

$$\int_{-1}^1 \mathbf{K}(\mathbf{u}) d\mathbf{u} = 1 \quad \text{and} \quad \int_{-1}^1 \mathbf{u}^j \mathbf{K}(\mathbf{u}) d\mathbf{u} = 0, \quad \forall j \in \{1, \dots, \ell\}.$$

Mirrored Kernel Density Estimator

Given a bandwidth \mathbf{h} , our Rényi- α divergence estimate is computed in 3 steps:

1. Mirror data over subsets of edges of \mathcal{X} .
2. Compute clipped kernel density estimates $\hat{\mathbf{p}}$ and $\hat{\mathbf{q}}$ from the mirrored data, using product kernel \mathbf{K}_d and bandwidth \mathbf{h} , and clipping the kernel density estimates pointwise below at κ_1 and above at κ_2 .
3. Estimate $D_\alpha(\mathbf{p} \parallel \mathbf{q})$ by the plug-in estimator $D_\alpha(\hat{\mathbf{p}} \parallel \hat{\mathbf{q}})$.

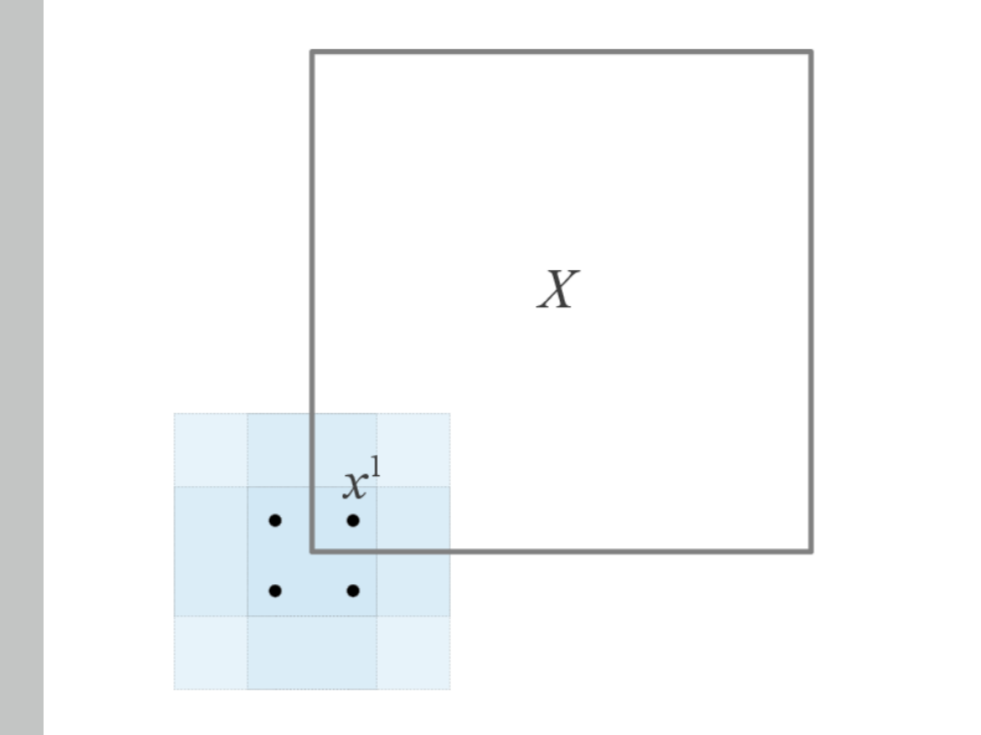


Figure : Four kernels centered on a single data point and its three reflected copies, in the case $d = 2$.

Results: Exponential Concentration Bound

- ▶ We show that, $\forall \varepsilon > 0$,

$$\mathbb{P}(|D_\alpha(\hat{\mathbf{p}} \parallel \hat{\mathbf{q}}) - \mathbb{E}D_\alpha(\hat{\mathbf{p}} \parallel \hat{\mathbf{q}})| > \varepsilon) \leq 2 \exp(-C_V \varepsilon^2 n),$$

where

$$C_V = \frac{|\alpha - 1|}{2C_f C_L \|\mathbf{K}\|_1^d}$$

is constant in n and \mathbf{h} .

- ▶ Main tool in proof is McDiarmid's Inequality, by which it suffices to bound the change in the estimate when resampling a single data point by C_V/n .
- ▶ This is achieved by combining a smoothness property of D_α with the observation that the integral of the mirrored kernel density estimate changes by at most $\frac{2}{n} \int_{[-1, 1]^d} |\mathbf{K}^d(\mathbf{u})| d\mathbf{u}$.

Bias Lemma

- ▶ **Bias Lemma:** Writing the pointwise bias of the clipped and mirrored kernel density as $\mathbf{b}_p(\mathbf{x}) = \mathbb{E}\hat{\mathbf{p}}(\mathbf{x}) - \mathbf{p}(\mathbf{x})$, we show

$$\int_{\mathcal{X}} \mathbf{b}_p^2(\mathbf{x}) d\mathbf{x} \leq C_b \mathbf{h}^{2\beta}.$$

- ▶ For somewhat small \mathbf{h} and large β (in particular, $\mathbf{h} \leq \sqrt{\frac{1}{3d^{1/\ell}}}$ and $\beta \geq 6d + 2$ suffices), one can show $C_b \leq 3L$.
- ▶ Away from the boundary of \mathcal{X} (i.e., in $[\mathbf{h}, 1 - \mathbf{h}]^d$), there is no boundary bias, and so we simply cite well-known results in kernel density estimation, using the assumed symmetry properties of the kernel.
- ▶ For \mathbf{x} near (within \mathbf{h} of) the boundary of \mathcal{X} , we combine the Smoothness and Boundary Conditions via a Taylor bound to derive a pointwise bound $\mathbf{b}_p(\mathbf{x}) \leq C_b \mathbf{h}$.

Results: Convergence Rate

- ▶ We show there exists $C_B \in \mathbb{R}$ (constant in n and \mathbf{h}) such that

$$|\mathbb{E}D_\alpha(\hat{\mathbf{p}} \parallel \hat{\mathbf{q}}) - D_\alpha(\mathbf{p} \parallel \mathbf{q})| \leq C_B \left(\mathbf{h}^\beta + \frac{1}{n\mathbf{h}^d} \right).$$
- ▶ Proven by making a second-order Taylor estimate and then using Hölder's Inequality to reduce terms to the Bias Lemma and the integrated mean squared error of a standard kernel density estimator.

Discussion

- ▶ The exponential concentration bound gives a bound on the variance of the estimator:

$$\mathbb{V}[F(\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k)] \leq C_V^2 n^{-1}.$$
- ▶ This does not depend on \mathbf{h} , so pick \mathbf{h} to minimize the bias bound.
 - ▷ Asymptotically optimal \mathbf{h} is $\mathbf{h} \asymp n^{-\frac{1}{\beta+d}}$, so bias bound is $\mathcal{O}\left(n^{-\frac{\beta}{\beta+d}}\right)$.
- ▶ Hence MSE is $\mathcal{O}(n^{-\frac{\beta}{\beta+d}} + n^{-1})$, which is the parametric rate $\mathcal{O}(n^{-1})$ if $\beta \geq d$ and $\mathcal{O}(n^{-\frac{\beta}{\beta+d}})$ otherwise.
- ▶ Kernel assumptions for the bias bound necessitate $\|\mathbf{K}\|_1 > 1$ when $\beta \geq 2$ and C_V includes $\|\mathbf{K}\|_1^d$, which is exponential in d .
 - ▷ Lower bounds in d are unknown; whether dependence is necessarily exponential is an important open problem.

Experimental Results on Synthetic Data

- ▶ $\vec{\mu}_1 = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.3 \end{bmatrix}, \vec{\mu}_2 = \begin{bmatrix} 0.7 \\ 0.7 \\ 0.7 \end{bmatrix}$
- ▶ $\Sigma = \begin{bmatrix} 0.2 & 0 & 0 \\ 0 & 0.2 & 0 \\ 0 & 0 & 0.3 \end{bmatrix}$
- ▶ $\mathbf{p}_1 = \mathcal{N}(\vec{\mu}_1, \Sigma), \mathbf{p}_2 = \mathcal{N}(\vec{\mu}_2, \Sigma)$
- ▶ In each trial, n points were drawn from \mathbf{p}_1 and \mathbf{p}_2 restricted to $[0, 1]^3$. $D_\alpha(\hat{\mathbf{p}} \parallel \hat{\mathbf{q}})$ was computed directly.

Mean squared error and standard deviation of our estimator were computed from 100 trials and plotted below.

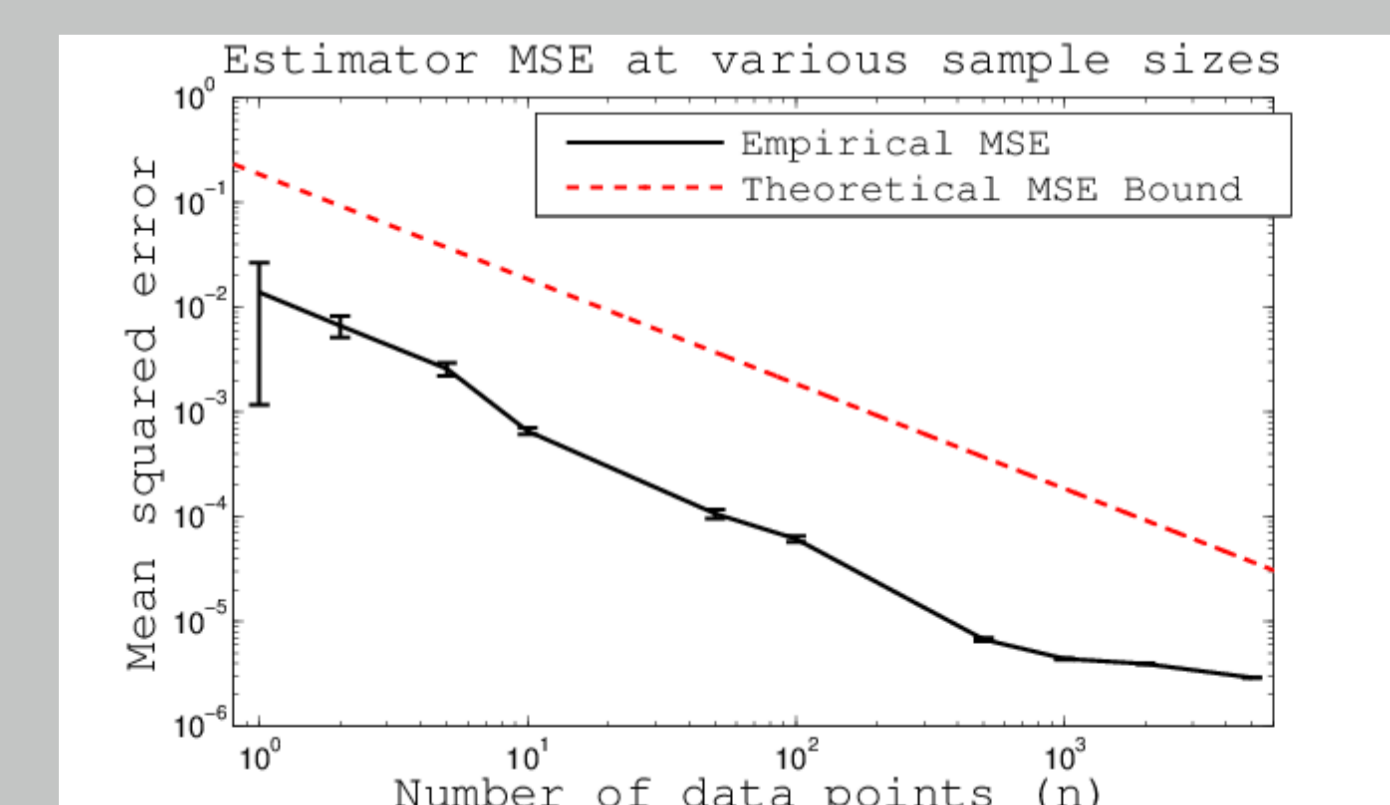


Figure : Log-log plot of mean squared (computed over 100 trials) for various sample sizes n , alongside our theoretical bound. Error bars indicate standard deviation of estimator over 100 trials.