

Introduction

- ▶ Many important statistical quantities can be written as

$$\mathbf{F}(\mathbf{p}_1, \dots, \mathbf{p}_k) = \int_{\mathcal{X}_1 \times \dots \times \mathcal{X}_k} \mathbf{f}(\mathbf{p}_1(\mathbf{x}_1), \dots, \mathbf{p}_k(\mathbf{x}_k)) \mathbf{d}(\mathbf{x}_1, \dots, \mathbf{x}_k),$$

where each $\mathbf{p}_i : \mathcal{X}_i \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^+$ is a probability density, $\mathbf{f} : \mathbb{R}^k \rightarrow \mathbb{R}$ is smooth.

- ▶ Examples of such quantities, which we call *Density Functionals*, are:
 - ▷ Shannon/KL, Rényi- α , and Tsallis- α entropies, mutual informations, and divergences
 - ▷ \mathbf{f} -divergences (e.g., Hellinger, Jensen-Shannon, etc.)
 - ▷ \mathbf{L}_p norms and distances
 - ▷ Conditional versions of the above quantities
- ▶ For many of these quantities, few consistent estimators are known, and almost none of these have finite-sample convergence of concentration guarantees.
- ▶ We propose and study a nonparametric estimator for such quantities, based on plugging in a boundary-corrected kernel density estimate.
- ▶ We prove that, when each $\mathcal{X}_i = [0, 1]^d$ is a unit cube:
 - ▷ our estimator is exponentially concentrated about its mean.
 - ▷ for densities in a β -Hölder smoothness class with certain boundary conditions, the bias of the estimator decays as $\mathcal{O}\left(n^{-\frac{\beta}{\beta+d}}\right)$, 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 \mathbf{f} , the densities $\mathbf{p}_1, \dots, \mathbf{p}_k$, the kernel \mathbf{K} :

- ▶ **(\mathbf{f} -Smoothness)** \mathbf{f} is twice continuously differentiable.
- ▶ **(Density Smoothness)** All (mixed) ℓ -order partial derivatives of $\mathbf{p}_1, \dots, \mathbf{p}_k$ exist and are $(\beta - \ell)$ -Hölder Continuous (i.e., there exists $\mathbf{L} \in \mathbb{R}$ such that, $\forall \mathbf{x}, \mathbf{x} + \mathbf{v} \in \mathcal{X}$, $|\vec{i}| = \ell$, each

$$|\mathbf{D}^{\vec{i}} \mathbf{p}_i(\mathbf{x} + \mathbf{v}) - \mathbf{D}^{\vec{i}} \mathbf{p}_i(\mathbf{x})| \leq \mathbf{L} \|\mathbf{v}\|_2^{\beta - \ell}.$$
- ▶ **(Density Boundaries)** All derivatives of $\mathbf{p}_1, \dots, \mathbf{p}_k$ of order up to ℓ vanish at the boundary

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

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

Mirrored Kernel Density Estimator

Given a bandwidth \mathbf{h} , our density functional estimate is computed in 3 steps:

1. Augment data from \mathbf{p}_i with reflections over each subset of edges of \mathcal{X}_i .
2. Compute kernel density estimates $\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k$ from the augmented data, using a product kernel and bandwidth \mathbf{h} .
3. Estimate $\mathbf{F}(\mathbf{p}_1, \dots, \mathbf{p}_k)$ by the plug-in estimator $\mathbf{F}(\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k)$.

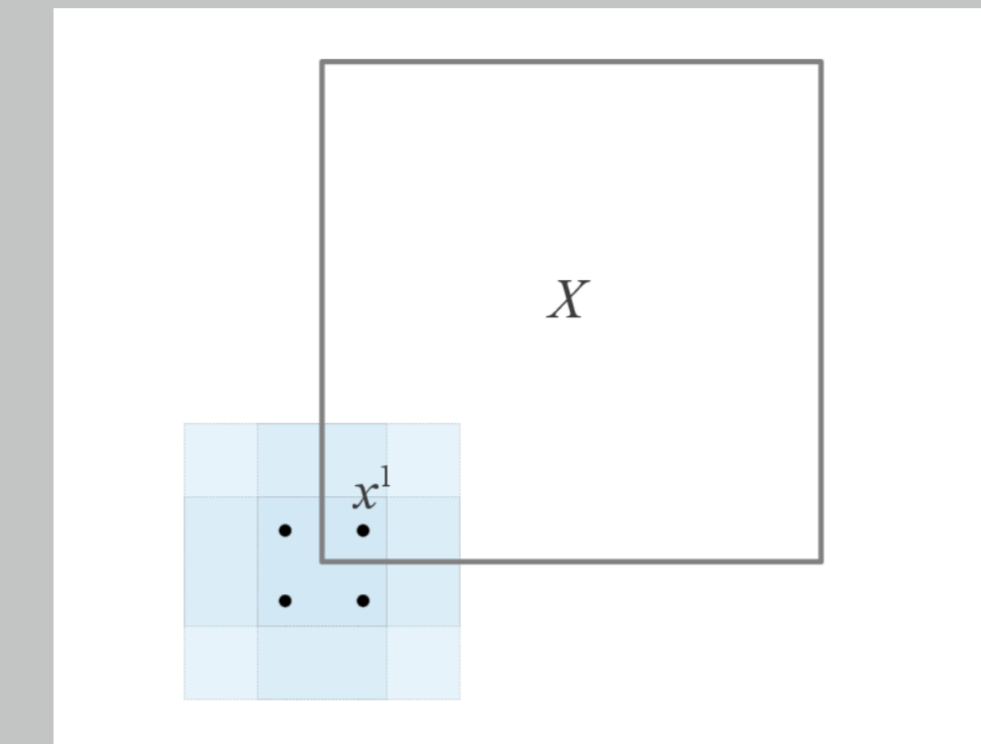


Figure : Four kernels centered on a single data point and its three reflected copies, in the case $d = 2$.

Results: Exponential Concentration Bound

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

$$\mathbb{P}\left(|\mathbf{F}(\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k) - \mathbb{E}\mathbf{F}(\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k)| > \varepsilon\right) \leq 2 \exp\left(-\frac{\varepsilon^2 n}{\mathbf{C}_V^2}\right),$$

where $\mathbf{C}_V = 2\mathbf{C}_f \|\mathbf{K}\|_1^d$ 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 \mathbf{C}_V/n .
- ▶ This is achieved by combining the smoothness of \mathbf{f} with the observation that the integral of the mirrored kernel density estimate changes by at most $\frac{2}{n} \|\mathbf{K}\|_1^d$.

Results: Convergence Rate

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

$$|\mathbb{E}\mathbf{F}(\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k) - \mathbf{F}(\mathbf{p}_1, \dots, \mathbf{p}_k)| \leq \mathbf{C}_B \left(\mathbf{h}^\beta + \frac{1}{n\mathbf{h}^d}\right).$$

- ▶ Previous work [2] bounded the integral of each mirrored kernel density estimator's pointwise squared bias:

$$\int_{\mathcal{X}_i} (\mathbb{E}\hat{\mathbf{p}}_i(\mathbf{x}) - \mathbf{p}(\mathbf{x}))^2 \mathbf{d}\mathbf{x} \leq \mathbf{C}_b \mathbf{h}^{2\beta} \quad (1)$$

- ▶ To derive our convergence rate, we make a second-order Taylor estimate of \mathbf{f} and then use Hölder's Inequality to reduce the resulting terms to (1) and the integrated mean squared error of a standard kernel density estimator.

Condition Density Functionals

- ▶ It is often useful to condition density functionals on one or more additional variables; e.g., to estimate

$$\mathbf{F}(\mathbf{P}) = \int_{\mathcal{Z}} \mathbf{P}_Z(\mathbf{z}) \mathbf{f}\left(\int_{\mathcal{X}} \mathbf{g}\left(\frac{\mathbf{P}_{\mathbf{X},\mathbf{Z}}(\mathbf{x}, \mathbf{z})}{\mathbf{P}_Z(\mathbf{z})}\right) \mathbf{d}\mathbf{x}\right) \mathbf{d}\mathbf{z}.$$

- ▶ For example, conditional entropy estimation has applications to clustering [3] and conditional mutual information is useful for learning graphical models [1].
- ▶ As long as the density of the conditioned variable (e.g., \mathbf{P}_Z) has a positive lower bound, our results extend to the mirrored kernel density plug-in estimator for such *conditional density functionals*.

Discussion

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

$$\mathbb{V}[\mathbf{F}(\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k)] \leq \mathbf{C}_V^2 n^{-1}.$$

- ▶ This does not depend on \mathbf{h} , so pick \mathbf{h} to minimize the bias bound.
 - ▷ Asymptotically, the 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 $\mathcal{O}\left(n^{-\frac{2\beta}{\beta+d}} + n^{-1}\right)$, which is the parametric rate $\mathcal{O}(n^{-1})$ if $\beta \geq d$ and $\mathcal{O}\left(n^{-\frac{\beta}{\beta+d}}\right)$ otherwise.
- ▶ Kernel assumptions for the bias bound necessitate $\|\mathbf{K}\|_1 > 1$ when $\beta \geq 2$ and \mathbf{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.
- ▶ For divergences and information theoretic density functionals, the integral in the plug-in estimator can be well estimated by a sample mean. In this case, our estimator can be computed in $\mathcal{O}(2^d n \log n)$, making it effective for large, low-dimensional datasets.

References

- ▷ D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, Cambridge, MA, 2009.
- ▷ S. Singh and B. Póczos. Generalized exponential concentration inequality for Rényi divergence estimation. In *International Conference on Machine Learning (ICML)*, 2014.
- ▷ G. ver Steeg, A. Galstyan, F. Sha, and S. DeDeo. Demystifying information-theoretic clustering. In *International Conference on Machine Learning (ICML)*, 2014.