

KERNEL DENSITY ESTIMATION USING LOW DISCREPANCY SAMPLING

FRED J. HICKERNELL

ABSTRACT. This project is where all of the files and commands go that are needed elsewhere

Let $Y = f(\mathbf{X})$, where $\mathbf{X} \sim \mathcal{U}[0, 1]^d$. Low discrepancy sequences are used for computing $\mu = \mathbb{E}(Y)$. Can they be used for estimation of ϱ , the probability density function of Y ?

Let ν be a probability mass or density function. A generalized kernel density estimator (KDE), $\widehat{\varrho}(\cdot, \nu, k)$, can be defined as

$$(1) \quad \widehat{\varrho}(y, \nu, k) := \int_{-\infty}^{\infty} k(z, y) \nu(z) dz, \quad y \in \mathbb{R},$$

where we call the $k : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ the *density kernel*. Here, $k(\cdot, y)$ is a weight function used to estimate the density at y , so $k(\cdot, y)$ should typically be larger near y than away from y . Moreover, we assume that

$$(2) \quad \int_{-\infty}^{\infty} k(z, y) dz = 1 \quad \forall y \in \mathbb{R}.$$

For (1) to be a good estimator, the density ν should approximate ϱ well.

Let $\varrho_{\mathbf{y}}$ be the empirical probability mass function of a vector of sampled Y values, $\mathbf{y} = (y_1, \dots, y_n)$. Then a *practical KDE* is

$$\widehat{\varrho}(y, \varrho_{\mathbf{y}}, k) = \int_{-\infty}^{\infty} k(z, y) \varrho_{\mathbf{y}}(z) dz = \frac{1}{n} \sum_{i=1}^n k(y_i, y) = \frac{1}{n} \sum_{i=1}^n k(f(\mathbf{x}_i), y).$$

Moreover, a *smoothed density* can be defined as

$$\widehat{\varrho}(y, \varrho, k) = \int_{-\infty}^{\infty} k(z, y) \varrho(z) dz = \int_{[0, 1]^d} k(f(\mathbf{x}), y) d\mathbf{x}.$$

The absolute error of the practical KDE can be bounded as the sum of two terms:

$$\begin{aligned} |\varrho(y) - \widehat{\varrho}(y, \varrho_{\mathbf{y}}, k)| &= |\varrho(y) - \widehat{\varrho}(y, \varrho, k) + \widehat{\varrho}(y, \varrho, k) - \widehat{\varrho}(y, \varrho_{\mathbf{y}}, k)| \\ &\leq |\varrho(y) - \widehat{\varrho}(y, \varrho, k)| + |\widehat{\varrho}(y, \varrho, k) - \widehat{\varrho}(y, \varrho_{\mathbf{y}}, k)|. \end{aligned}$$

The first term depends on the density kernel and measures how well the smoothed density estimates the true density. The second term measures how well the practical KDE estimates the smoothed density. We will analyze the two terms separately assuming that ϱ and $k(f(\cdot), y)$ lie in reproducing kernel Hilbert spaces.

Let $K^y : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be a reproducing kernel for a Hilbert space containing ϱ . Then an upper bound on the first term in the expression for the error can be

bounded via the Cauchy-Schwarz inequality and the Riesz representation theorem:

$$\begin{aligned}
 (3) \quad |\varrho(y) - \widehat{\varrho}(y, \varrho, k)| &= \left| \varrho(y) - \int_{-\infty}^{\infty} k(z, y) \varrho(z) dz \right| \\
 &= \left| \left\langle K^y(\cdot, y) - \int_{-\infty}^{\infty} k(z, y) K^y(z, \cdot) dz, \varrho \right\rangle_{K^y} \right| \\
 &\leq \left\| K^y(\cdot, y) - \int_{-\infty}^{\infty} k(z, y) K^y(z, \cdot) dz \right\|_{K^y} \|\varrho\|_{K^y},
 \end{aligned}$$

where the part of the error bound depending on density kernel can be rewritten in terms of integrals:

$$\begin{aligned}
 (4) \quad &\left\| K^y(\cdot, y) - \int_{-\infty}^{\infty} k(z, y) K^y(z, \cdot) dz \right\|_{K^y}^2 \\
 &= K^y(y, y) - 2 \int_{-\infty}^{\infty} k(z, y) K^y(z, y) dz \\
 &\quad + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} k(z, y) K^y(z, t) k(t, y) dz dt.
 \end{aligned}$$

Equation (3) is an upper bound on the first term in the expression for the error that separates the part depending on the choice of the density kernel from the part depending on the probability density.

Let $K^{\mathbf{x}} : [0, 1]^d \times [0, 1]^d \rightarrow \mathbb{R}$ be a reproducing kernel for a Hilbert space containing $k(y, f(\cdot))$ for all $y \in \mathbb{R}$ and all f of interest. Then

$$\begin{aligned}
 |\widehat{\varrho}(y, \varrho, k) - \widehat{\varrho}(y, \varrho_{\mathbf{y}}, k)| &= \left| \int_{-\infty}^{\infty} k(z, y) \varrho(z) dz - \frac{1}{n} \sum_{i=1}^n k(y_i, y) \right| \\
 &= \left| \int_{[0,1]^d} k(f(\mathbf{x}), y) d\mathbf{x} - \frac{1}{n} \sum_{i=1}^n k(f(\mathbf{x}_i), y) \right| \\
 &= \left| \left\langle \int_{[0,1]^d} K^{\mathbf{x}}(\cdot, \mathbf{x}) d\mathbf{x} - \frac{1}{n} \sum_{i=1}^n K^{\mathbf{x}}(\cdot, \mathbf{x}_i), k(f(\cdot), y) \right\rangle_{K^y} \right| \\
 &\leq \left\| \int_{[0,1]^d} K^{\mathbf{x}}(\cdot, \mathbf{x}) d\mathbf{x} - \frac{1}{n} \sum_{i=1}^n K^{\mathbf{x}}(\cdot, \mathbf{x}_i) \right\|_{K^{\mathbf{x}}} \|k(f(\cdot), y)\|_{K^{\mathbf{x}}},
 \end{aligned}$$

where

$$\begin{aligned}
 &\left\| \int_{[0,1]^d} K^{\mathbf{x}}(\cdot, \mathbf{x}) d\mathbf{x} - \frac{1}{n} \sum_{i=1}^n K^{\mathbf{x}}(\cdot, \mathbf{x}_i) \right\|_{K^{\mathbf{x}}}^2 \\
 &= \int_{[0,1]^d \times [0,1]^d} K^{\mathbf{x}}(\mathbf{x}, t) d\mathbf{x} dt - \frac{2}{n} \sum_{i=1}^n \int_{[0,1]^d} K^{\mathbf{x}}(\mathbf{x}, \mathbf{x}_i) d\mathbf{x} \\
 &\quad + \frac{1}{n^2} \sum_{i,j=1}^n K^{\mathbf{x}}(\mathbf{x}_i, \mathbf{x}_j).
 \end{aligned}$$

This is an upper bound on the second term in the expression for the error that separates the part depending on the choice of the kernel evaluated at $(f(\cdot), y)$ from the part depending on the sample nodes.

1. ISOTROPIC DENSITY KERNELS

Suppose that $k(z, y) = h^{-1}\tilde{k}((z - y)/h)$ where h is a bandwidth, and $K^y(y, z) = \tilde{K}^y(y - z)$. Then the quantity measuring the quality of the density kernel in (4) may be written as

$$\begin{aligned} & \left\| K^y(\cdot, y) - \int_{-\infty}^{\infty} k(z, y) K^y(z, \cdot) dz \right\|_{K^y}^2 \\ &= \tilde{K}^y(0) - \frac{2}{h} \int_{-\infty}^{\infty} \tilde{k}((z - y)/h) \tilde{K}^y(z - y) dz \\ & \quad + \frac{1}{h^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \tilde{k}((z - y)/h) \tilde{K}^y(z - t) \tilde{k}((t - y)/h) dz dt \\ &= \tilde{K}^y(0) - 2 \int_{-\infty}^{\infty} \tilde{k}(w) \tilde{K}^y(hw) dw \\ & \quad + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \tilde{k}(w) \tilde{K}^y(h(w - v)) \tilde{k}(v) dw dv. \end{aligned}$$

via the variable transformations $w = (z - y)/h$ and $v = (t - y)/h$.

Let $\tilde{K}^y(y) = \exp(-y^2/2)$ and $\tilde{k}(y) = \exp(-y^2/2)/\sqrt{2\pi}$.

REFERENCES