

Lecture # 22

①

To summarize from last time, if we consider estimating the density f_X of a r.v. on $[0,1]$ with a histogram with uniform bins of width h , then for an iid sample of size n from X , we have the following error:

$$\text{MSE} = \text{Bias}^2 + \text{Variance}$$

$$= \int_0^1 \left[\mathbb{E}(\hat{f}_n(x)) - f_X(x) \right]^2 dx + \int_0^1 \mathbb{E} \left(\mathbb{E}(\hat{f}_n(x)) - \hat{f}_n(x) \right)^2 dx$$

$$\approx \frac{h^2}{12} \cdot \int_0^1 [f'_X(x)]^2 dx + \frac{1}{h} \cdot \frac{1}{n}$$

So, small h (tiny bins) favor low bias at the expense of high variance.
Vice versa for large bins.

Q = What is the optimal h ? We can phrase this as a calculus problem, because $g(h) := \frac{h^2}{12} \cdot \int_0^1 [f'_X(x)]^2 dx + \frac{1}{h} \cdot \frac{1}{n}$ is smooth in h (any

from $h=0$). Let's solve $g'(h) = 0$:

$$\Gamma \quad g'(h) = \frac{h}{6} \cdot \int_0^1 [f'_X(x)]^2 dx - \frac{1}{h^2} \cdot \frac{1}{n} = 0$$

$$\Leftrightarrow h^3 \cdot \int_0^1 [f'_X(x)]^2 dx - \frac{6}{n} = 0$$

$$\Leftrightarrow h = \left(\frac{6}{n \cdot \int_0^1 [f'_X(x)]^2 dx} \right)^{1/3}$$

$$= \frac{1}{n^{1/3}} \cdot \left(\frac{6}{\int_0^1 [f'_X(x)]^2 dx} \right)^{1/3}$$

So, as a function of the sample size, h should scale like $h \sim n^{-1/3}$ to have optimal error. Plugging this back into the MSE gives an optimal MSE of

$$\text{MSE} \approx n^{-2/3} \cdot \left(\frac{3}{4} \right)^{2/3} \cdot \left[\int_0^1 [f'_X(x)]^2 dx \right]^{1/3}$$

So, as $n \rightarrow \infty$, the error $\rightarrow 0$ at a rate of $n^{-2/3}$.

Q: Is this the best any estimator for f_X can do? No! In fact, a simple extension of histogram estimation improves on that convergence rate: Kernel Density Estimation (KDE).

Defn: A function $K: \mathbb{R} \rightarrow \mathbb{R}$ is a kernel if:

(1.) K is infinitely differentiable.

(2.) $\int_{\mathbb{R}} K(x) dx = 1$

(3.) $\int_{\mathbb{R}} x K(x) dx = 0$

(4.) $\int_{\mathbb{R}} x^2 K(x) dx = \sigma^2 > 0$.

ex: $K(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$ or more generally

$K_\sigma(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x^2}{2\sigma^2}\right)$

is the Gaussian kernel.

ex: $K(x) = \begin{cases} \frac{3}{4} \left(1 - \frac{x^2}{5}\right) & , |x| < \sqrt{5} \\ 0, & \text{else} \end{cases}$

ex: One can take $\tilde{K}(x) = \begin{cases} 1, & |x| \leq \frac{1}{2} \\ 0, & \text{else} \end{cases}$

and mollify/smooth it to be C^∞ while preserving its area under the curve and symmetry. in certain places

The idea of KDE is to put a kernel at each sample, instead of a rectangle as in a histogram. So, KDE can be understood as even more data adaptive than histogram estimation.

Data: Let K be a kernel. Given a sample $\{x_i\}_{i=1}^n \subset \mathbb{R}$ the KDE associated to the sample, K , and bandwidth $h > 0$ is $\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n h^{-1} K\left(\frac{x-x_i}{h}\right)$.