

Chapter 2.

References:

↳ Chap. 24 rdV

↳ Sec. 6.3 All of Stats.

- We observe $X_1, \dots, X_n \stackrel{iid}{\sim} Q$

- Let $F(x) = Q\{X \leq x\}$ denote the CDF
 $f(x)$ denote the density.

- Our goal: Estimate f at a point x_0 .

- Given that $f(x_0) = \frac{d}{dx} F(x) \big|_{x=x_0}$ we might consider estimating $f(x_0)$ w/ a plug-estimator

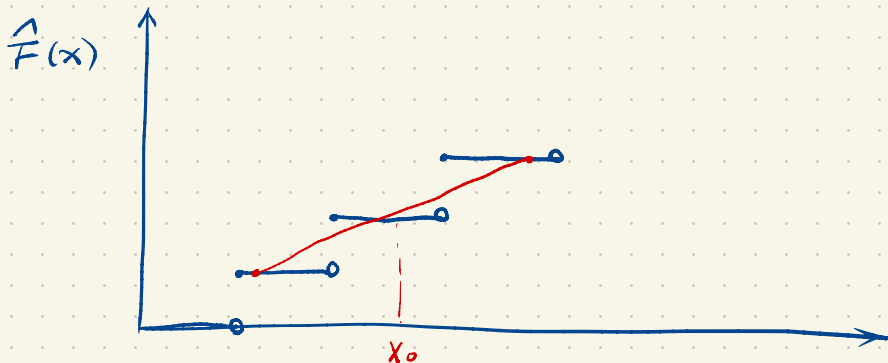
$$\hat{f}(x_0) = \frac{d}{dx} \hat{F}(x) \big|_{x=x_0} \quad \text{empirical CDF}$$

- One way to estimate F is via the empirical CDF

$$(1) \quad \hat{F}(x_0) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq x_0\}$$

- We know that this is a "good" estimator of CDF.
 \hookrightarrow e.g.: $\sqrt{n} [\hat{F}(x_0) - F(x_0)] \rightsquigarrow N(0, \sigma^2)$

- But: Estimating $f(x_0)$ via (1) turns out to be a very bad idea



A lot of things will be 0.

A better one will lead us to KDE.

- Another option uses that

$$f(x_0) = \lim_{h \rightarrow 0} \frac{F(x_0+h) - F(x_0-h)}{2h}$$

and so, when h is ^{small}

$$f(x_0) \approx \frac{F(x_0+h) - F(x_0-h)}{2h}$$

- This suggests an estimator

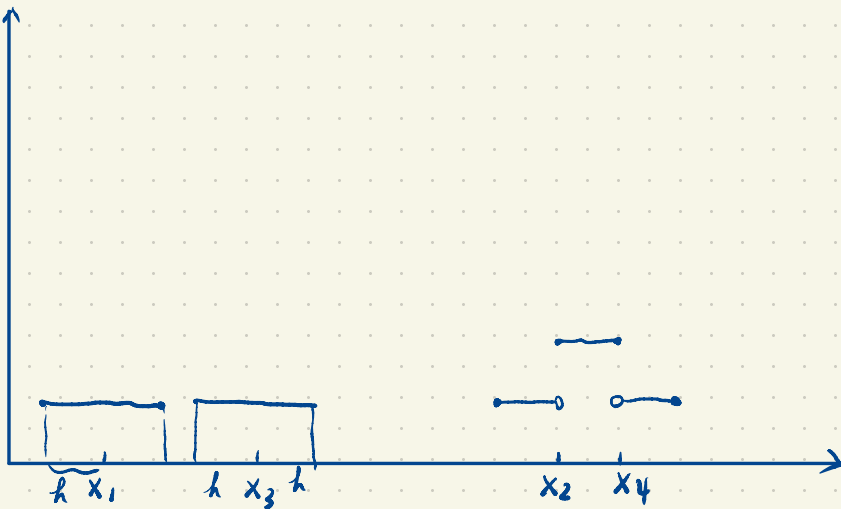
$$\hat{f}_h(x_0) := \frac{\hat{F}(x_0+h) - \hat{F}(x_0-h)}{2h}$$

$$= \frac{1}{2nh} \sum_{i=1}^n \mathbb{1}\{x_0-h < x_i \leq x_0+h\}$$

(a.s.)

$$= \frac{1}{nh} \sum_{i=1}^n \frac{1}{2} \mathbb{1}\left\{\frac{|x_i - x_0|}{h} \leq 1\right\} \star$$

$\hat{f}_h(x)$



Note: This estimate of f is not smooth.

Question: Can we define a smoother estimate of f ?

Ans: Yes!

- Let $K: \mathbb{R} \rightarrow \mathbb{R}$ be a kernel, that is, a function satisfying $\int K(u) du = 1$.

Def: An s -th order kernel is a kernel K that satisfies

$$\int u^r K(u) du = 0 \quad \text{for } r=1, 2, \dots, s-1$$

$$|\int u^s K(u) du| < \infty$$

↳ If K is symmetric about zero, [$K(u) = K(-u)$] then K is always at least a 2nd order kernel ($s=2$)

↳ Using higher order kernels ($s > 2$) can lead to estimators with lower bias.

General form of the KDE

For $h > 0$,

$$\begin{aligned} \hat{f}_h(x_0) &:= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right) \\ &= \frac{1}{n} \sum_{i=1}^n K_h\left(\frac{x_i - x_0}{h}\right) \end{aligned}$$

where $K_h(u) = \frac{1}{h} K\left(\frac{u}{h}\right)$

Example of kernels:

1) Uniform : $K(u) = \frac{1}{2} \mathbb{1}_{\{ |u| \leq 1 \}}$ ★
bounded support

2) Epanechnikov : $K(u) = \frac{3}{4} (1 - u^2) \mathbb{1}_{\{ |u| \leq 1 \}}$
bounded support

3) Gaussian : $K(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right)$

Note: All of these kernels are 2nd-order.

Nice fact: If K is non-negative, then for any $h > 0$,

\hat{f}_h is a PDF.

$$\begin{aligned} & \int \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right) dx \\ &= \frac{1}{n} \sum \int \frac{1}{h} K\left(\frac{x_i - x_0}{h}\right) dx \quad \left(u = \frac{x_i - x}{h}\right) \\ &= \frac{1}{n} \sum_{i=1}^n \int K(u) du = 1. \end{aligned}$$

Studying $\hat{f}_h(x_0)$:

We now study the performance of $\hat{f}_h(x_0)$ as an estimator of $f(x_0)$.

We'll quantify the performance in terms of the MSE:

$$\mathbb{E}[\{\hat{f}_h(x_0) - f(x_0)\}^2] = \underbrace{\{\mathbb{E}[\hat{f}_h(x_0)] - f(x_0)\}^2}_{\text{bias}^2} + \underbrace{\text{var}[\hat{f}_h(x_0)]}_{\text{variance}}$$

Here, we'll suppose that f belongs to (β, L) Hölder class w/ $\beta=2$, and $L > 0$.

↳ Note: All of the calculations we do still go through if the restriction of f to a nbhd. of x_0 is (β, L) -Hölder.

Recall: Saying that f is $(2, L)$ Hölder means that

$$|f'(x_1) - f'(x_2)| \leq L|x_1 - x_2| \quad \forall x_1, x_2$$

We focus on the case where the kernel K is:

- ↳ bounded
- ↳ non-negative
- ↳ 2nd order
- ↳ bounded support

We're going to see that choosing a small bandwidth $h > 0$ yields low bias and high variance and vice versa.

Bias of KDE

$$\text{Recall: Bias} = \mathbb{E}[\hat{f}_h(x_0)] - f(x_0)$$

$$\mathbb{E}[\hat{f}_h(x_0)] = \frac{1}{nh} \sum_{i=1}^n \mathbb{E}\left[K\left(\frac{x_i - x_0}{h}\right)\right]$$

$$\text{(iid)} = \frac{1}{h} \mathbb{E}\left[K\left(\frac{x_i - x_0}{h}\right)\right]$$

$$= \frac{1}{h} \int K\left(\frac{x_i - x_0}{h}\right) f(x_i) dx_i \quad \left(u = \frac{x_i - x_0}{h}\right)$$

$$= \int K(u) f(x_0 + uh) du$$

Recalling $\int K(u) du = 1$, and so

$$\begin{aligned} \text{Bias} &= \mathbb{E}[\hat{f}_h(x_0)] - f(x_0) \\ &= \int K(u) \underbrace{[f(x_0 + uh) - f(x_0)]}_{\star} du \end{aligned} \quad \text{Smooth!}$$

By the mean-value Thm., we know there exists

\tilde{x}_{uh} such that

$$\star = uh f'(\tilde{x}_{uh})$$

Hence

$$\text{Bias} = \int K(u) uh f'(\tilde{x}_{uh}) du$$

$$= \int K(u) uh f'(x_0) du + \int K(u) uh [f'(\tilde{x}_{uh}) - f'(x_0)] du$$

$$\text{2nd order} \leftarrow = h f'(x_0) \int K(u) u du$$

$$= 0$$

Hence

$$|\text{Bias}| = \left| \int K(u) u h [f'(\tilde{x}_{uh}) - f'(x_0)] du \right|$$

$$\text{Jensen} \leq h \int K(u) |u| |f'(\tilde{x}_{uh}) - f'(x_0)| du$$

$$\text{Hölder (Lipschitz)} \leq L h \int K(u) |u| |\tilde{x}_{uh} - x_0| du$$

$$|\tilde{x}_{uh} - x_0| \leq L h^2 \underbrace{\int K(u) u^2 du}_{\sigma_K^2}$$

$$= L \sigma_K^2 h^2$$

$$\text{Hence, } \boxed{\text{Bias}^2 \leq L^2 \sigma_K^4 h^4}$$

Variance of the KDE

$$\text{Var}(\hat{f}_h(x_0)) = \text{Var}\left[\frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x_0}{h}\right)\right]$$

$$\text{(independence)} = \frac{1}{nh^2} \sum_{i=1}^n \text{var}\left[K\left(\frac{X_i - x_0}{h}\right)\right]$$

$$\text{(identical)} = \frac{1}{nh^2} \text{var}\left[K\left(\frac{X_1 - x_0}{h}\right)\right]$$

$$\leq \frac{1}{nh^2} \mathbb{E}\left[K\left(\frac{X_1 - x_0}{h}\right)^2\right] \leftarrow \text{2nd moment}$$

$$= \frac{1}{nh^2} \int K\left(\frac{X_1 - x_0}{h}\right)^2 f(x_1) dx_1, \quad u = \frac{x_1 - x_0}{h}$$

$$= \frac{1}{nh} \underbrace{\int K(u)^2 f(x_0 + uh) du}_{\text{★★}} \quad (2)$$

We'll study $\star\star$ in what follows.

To do this, we'll make use of two facts:

- 1) f is Hölder continuous \Rightarrow continuous.
- 2) K has bounded support.

Let $K_1 = \inf\{u : k(u) > 0\}$, $K_2 = \sup\{u : k(u) > 0\}$.
We have that

$$\begin{aligned}\star\star &= \int K(u)^2 f(x_0 + uh) du \\ &= \int_{K_1}^{K_2} K(u)^2 f(x_0 + uh) du \\ &\leq \left[\sup_{u \in [K_1, K_2]} f(x_0 + uh) \right] \int_{K_1}^{K_2} K(u)^2 du\end{aligned}$$

If $h \leq 1$, then this shows that

$$\leq \underbrace{\left[\sup_{t \in [K_1, K_2]} f(x_0 + t) \right] \int_{K_1}^{K_2} K(u)^2 du}_{=: \tilde{C}}$$

Hence, we have shown that (by plugging in)

$$\text{Var}(\hat{f}_h(x_0)) \leq \frac{\tilde{C}}{nh}$$

Plugging in our bound on the bias^2 and variance , we find that MSE

$$\text{MSE} \leq L^2 \sigma_K^4 h^4 + \frac{\tilde{c}}{nh}$$

By setting $L^2 \sigma_K^4 h^4 = \frac{\tilde{c}}{nh}$

$$\Rightarrow h = c n^{-1/5}$$

$$\therefore \text{MSE} \leq O(n^{-4/5})$$

Generalization

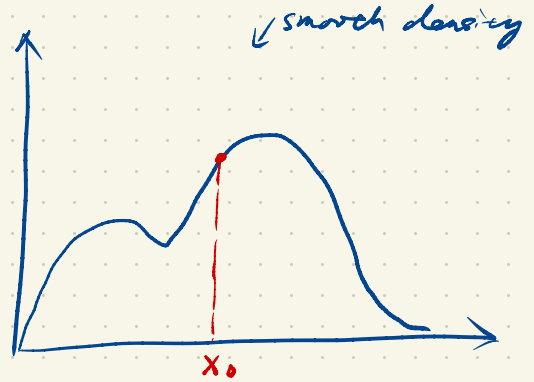
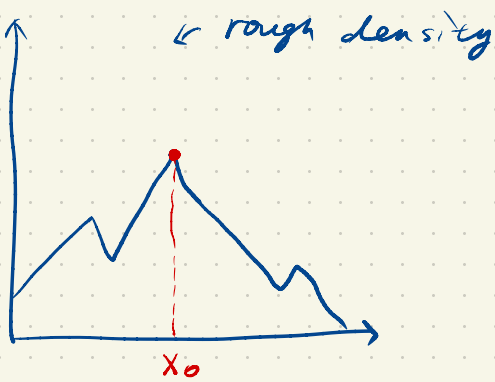
1) 1-dimensional setting w/ different amounts of smoothness

If f belongs to a (β, L) Hölder class, then similar arguments show that

$$\text{MSE} = O\left(n^{-\frac{2\beta}{2\beta+1}}\right)$$

if a kernel has sufficiently high order.

We saw this is class w/ $\beta = 2$.



sufficiently high order kernel leads to very smooth density which makes x_0 easier to be estimated.

2) d -dimensional probs.

Suppose X is d -dimensional and we want to estimate $f(x_0)$ at a fixed $x_0 \in \mathbb{R}^d$.

A KDE in this setting takes the form

$$\hat{f}_h(x_0) = \frac{1}{nh^d} \sum_{i=1}^n \prod_{j=1}^d K\left(\frac{x_{ij} - x_{0j}}{h}\right)$$

If f is β -times differentiable and all partial derivatives up to order β is bounded, then

$$MSE = O\left(n^{-\frac{2\beta}{2\beta+d}}\right)$$

Note: The dimension d appears in the denominator of exponent.

d	MSE UB
1	$n^{-0.8}$
2	$n^{-0.67}$
4	$n^{-0.5}$
10	$n^{-0.29}$

One way of thinking about the exponents on n :

If $\text{MSE} \propto n^{-\alpha}$, then, when n is large, to have MSE, you need to collect about $2^{1/\alpha}$ more data.