

Non-parametric density estimation

Reading material: https://faculty.washington.edu/yenchic/18W_425/Lec6_hist_KDE.pdf

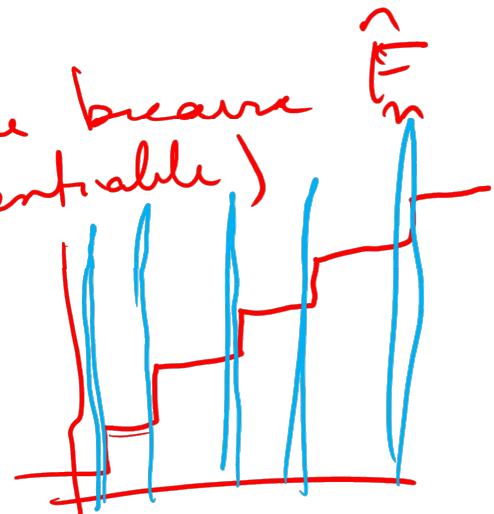
Motivation

- $D = \{X_1, X_2, \dots, X_n\}$ sampled i.i.d from an unknown $f(X)$
- Estimate $\hat{f}(x)$ without committing on a specific parametric form of $f(X)$
- Why not use empirical CDF?

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$$

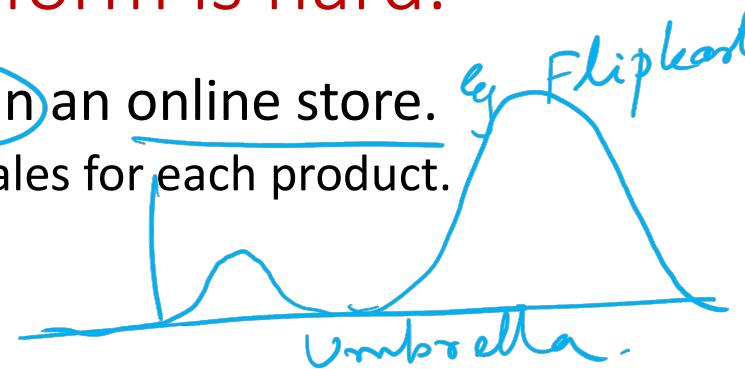
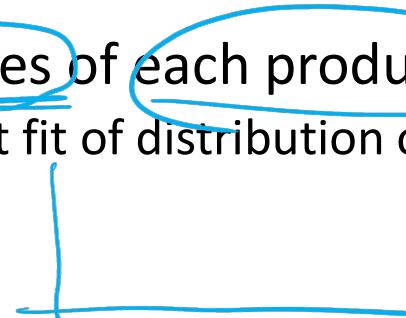
$$\hat{f}_n(x) = \frac{\partial}{\partial x} \hat{F}_n(x) \quad (\text{Not very usable because it is not differentiable})$$

- Too inefficient to maintain. Need to store entire data
- Too jerky. Non-zero density at observed points, zero elsewhere.
- Density estimation: assume some form of smoothness of $f(X)$

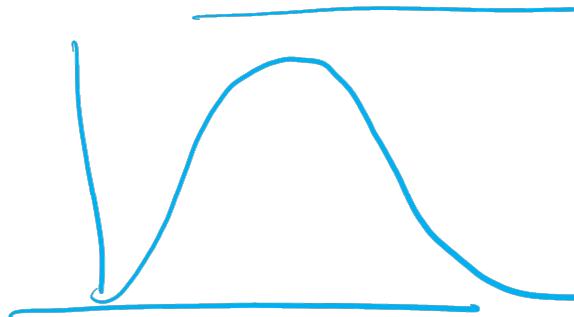
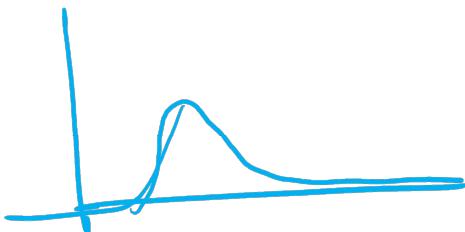


Real-life examples where parametric form is hard.

- Estimate density for weekly sales of each product in an online store.
 - Difficult to know exactly the best fit of distribution of sales for each product.

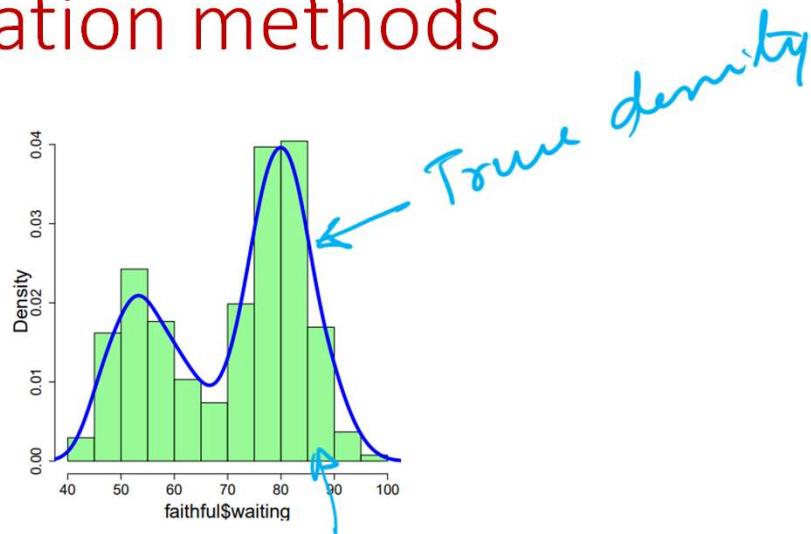


- Estimate distribution of number of steps taken by any arbitrary individual on a phone.

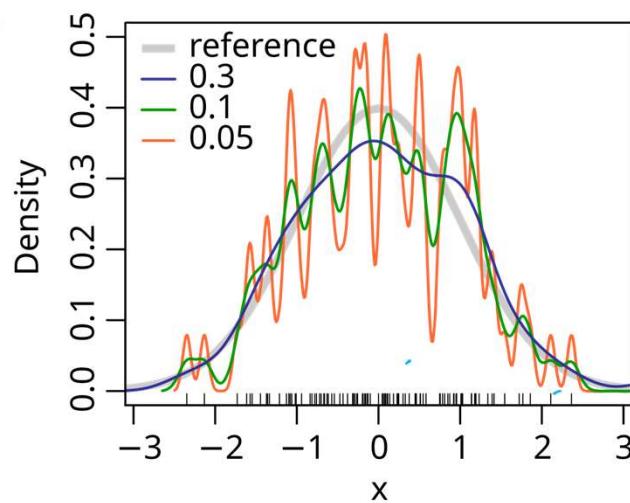


Density estimation methods

- Histogram



- Kernel density

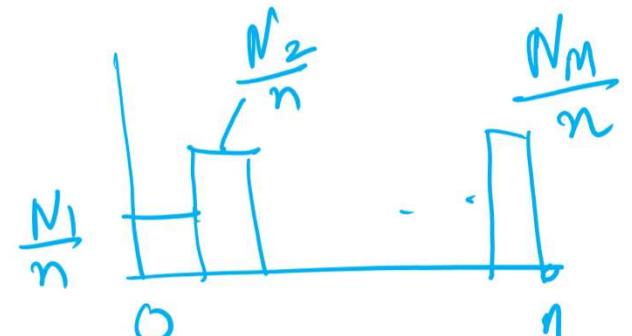


Histogram

- Let $X \in [0,1]$.
- To estimate $\hat{f}_n(x)$ using $D: \{x_1, x_2, \dots, x_n\}$
- Partition the $[0,1]$ range into M equal sized bins.

$$B_1 = \left[0, \frac{1}{M}\right], B_2 = \left[\frac{1}{M}, \frac{2}{M}\right], \dots, B_{M-1} = \left[\frac{M-2}{M}, \frac{M-1}{M}\right], B_M = \left[\frac{M-1}{M}, 1\right].$$

$$\frac{1}{M}$$



- Estimated density

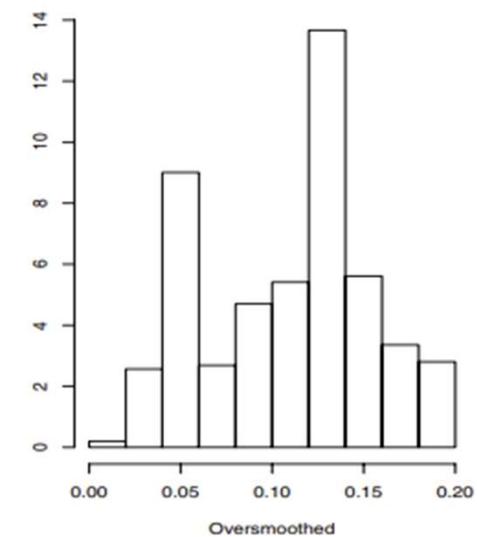
$$\begin{aligned} \hat{f}_n(x) &= (\text{Fraction of instances in } D \text{ in } B_j) / \text{Width of bin} \\ &= \frac{\sum_{i=1}^n I(x_i \in B_j)}{n} \cdot \frac{1}{M} \end{aligned}$$

\downarrow

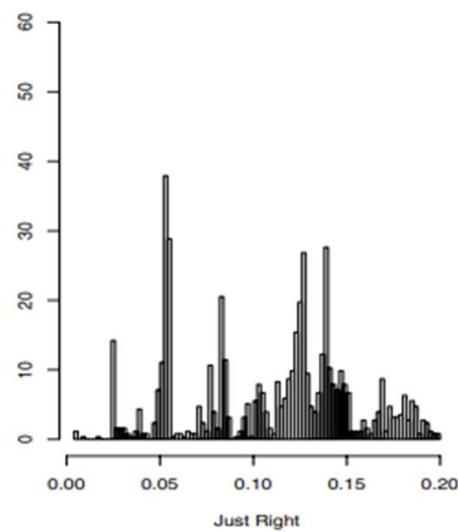
$x \in B_j \Rightarrow \frac{j-1}{M} \leq x \leq \frac{j}{M}$

$N_j = \# \text{ of } x_i \text{ in } B_j$
 $\text{in } D$

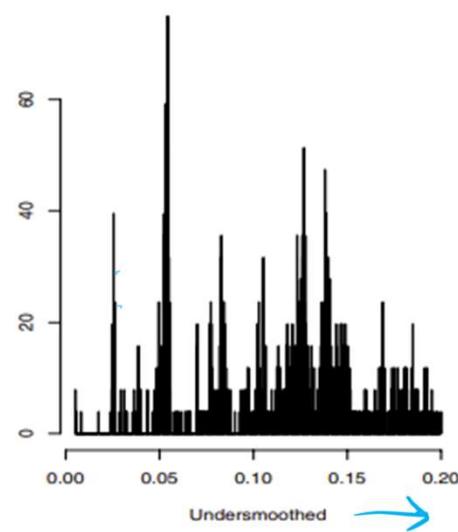
20.2 Example. Figure 20.3 shows three different histograms based on $n = 1,266$ data points from an astronomical sky survey. Each data point represents the distance from us to a galaxy. The galaxies lie on a “pencilbeam” pointing directly from the Earth out into space. Because of the finite speed of light, looking at galaxies farther and farther away corresponds to looking back in time. Choosing the right number of bins involves finding a good tradeoff between bias and variance. We shall see later that the top left histogram has too few bins resulting in oversmoothing and too much bias. The bottom left histogram has too many bins resulting in undersmoothing and too few bins. The top right histogram is just right. The histogram reveals the presence of clusters of galaxies. Seeing how the size and number of galaxy clusters varies with time, helps cosmologists understand the evolution of the universe. ■



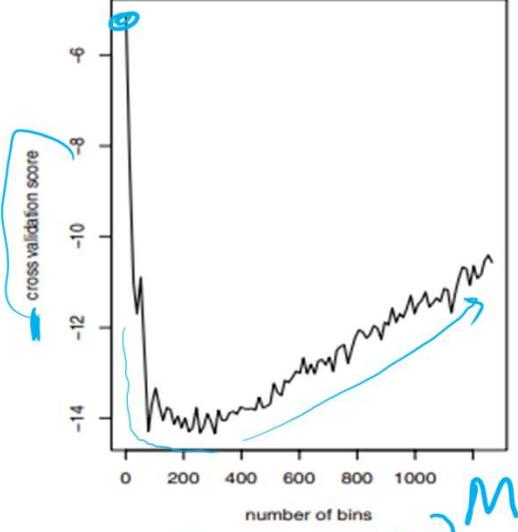
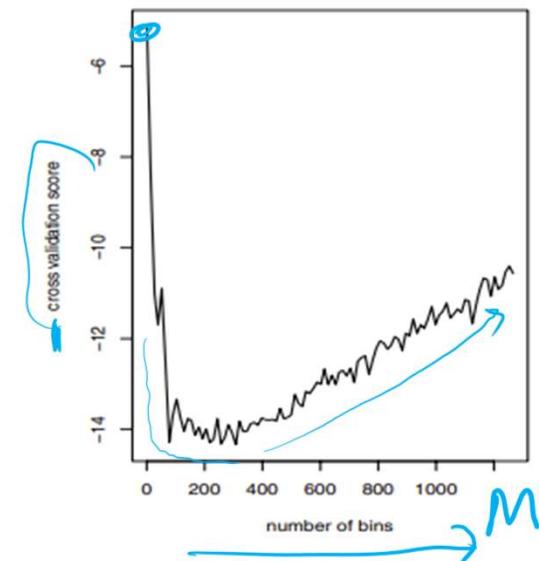
Oversmoothed



Just Right



Undersmoothed



Undersmoothed

FIGURE 20.3. Three versions of a histogram for the astronomy data. The top left histogram has too few bins. The bottom left histogram has too many bins. The top right histogram is just right. The lower, right plot shows the estimated risk versus the number of bins.

$n = 1266$
 $x \sim$ distance
 observed
 from earth
 of an
 galaxy

Bias, Variance, Risk

$$E[\hat{f}_n(x)] - f(x) \equiv \text{Bias}$$

$$E[\hat{f}_n(x)] = E\left[\frac{1}{n} \sum_{i=1}^n I(x_i \in B_j)\right] \text{ if } x \in B_j$$

$$= \frac{M}{n} \cdot n E_x(I(x \in B_j))$$

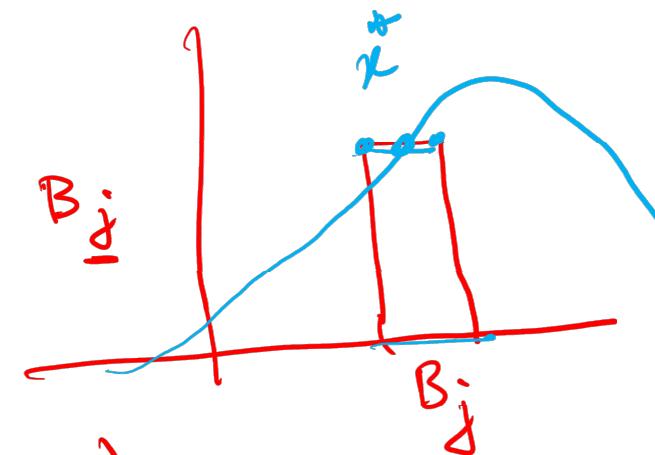
$$= \frac{M}{n} \cdot n \int_{x \in B_j} f(x) dx = M P(x \in B_j)$$

$$= M [F\left(\frac{j}{M}\right) - F\left(\frac{(j-1)}{M}\right)]$$

$$= F\left(\frac{j}{M}\right) - F\left(\frac{(j-1)}{M}\right) = F'(x^*) = f(x^*)$$

$$\left(\frac{j}{M} - \frac{(j-1)}{M}\right)$$

$$P(x \in B_j) = \frac{f(x^*)}{M}$$



$x^* \in B_j$
 $f(x)$ is differentiable

Variance of $\hat{f}_n(x)$

$$\begin{aligned}\text{Var}(\hat{f}_n(x)) &= \frac{M^2}{n^2} \text{Var}\left(\sum_{i=1}^n I(x_i \in B_j)\right) \quad \text{if } x \in B_j. \\ &= \frac{M^2}{n} \text{Var}_x(I(x \in B_j)) \\ &= \frac{M^2}{n} P(x \in B_j)(1 - P(x \in B_j)) \\ &= \frac{M^2}{n} \left[\frac{f(x^*)}{M} \right] \left[1 - \frac{f(x^*)}{M} \right] \\ &= \frac{M f(x^*)}{n} - \frac{f(x^*)^2}{n}\end{aligned}$$

Variance increases as M is increased
 $M \equiv \# \text{ of bins:}$

$$\text{Risk}(\hat{f}_n(x)) = \text{Bias}^2 + \text{Var}$$

$$\underbrace{[f(x^*) - f(\bar{x})]^2}_{\text{Bias}} + \frac{\underline{M}f(x^*)}{n} - \frac{\overline{f}(x^*)}{n}^2 \quad x^* \in B_\delta$$

$$= f'(\tilde{x})(\cancel{x^* - \bar{x}}) + \frac{\underline{M}f(x^*)}{n} - \frac{\overline{f}(x^*)}{n}^2$$

MV thm:

$$\leq \underbrace{|f'(\tilde{x})| \left\{ \frac{1}{M} \right\}}_{\text{UB on bias}} + \underbrace{\frac{\underline{M}f(x^*)}{n} - \frac{\overline{f}(x^*)}{n}}_{\text{variance}}$$

Bias drops with increasing M vs Variance
which increases.

Kernel Density Estimation

- One of the most convenient and accurate ways to estimate any density.

- Given data sample

$$D = \{x_1, x_2, \dots, x_n\}.$$

- Assume a special function called a kernel that puts a density mass around each training point.

→ 1. $K(x)$ is symmetric.

2. $\int_K K(x)dx = 1.$

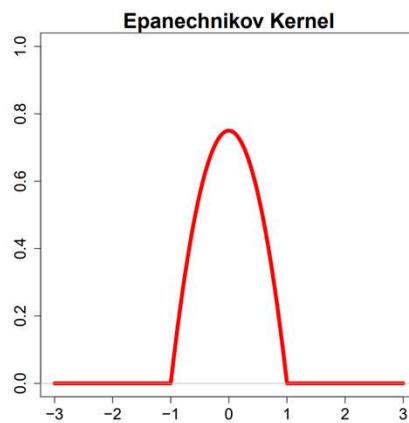
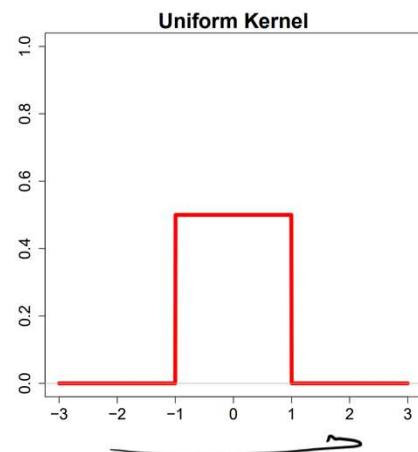
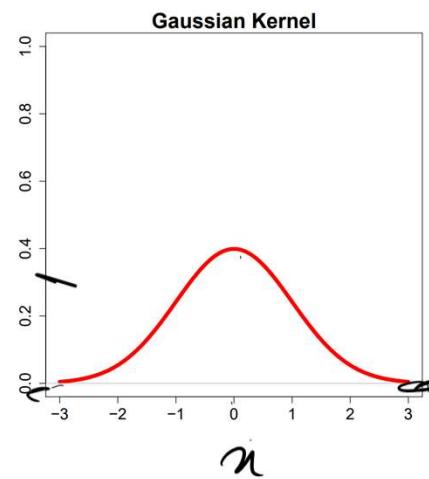
3. $\lim_{x \rightarrow -\infty} K(x) = \lim_{x \rightarrow +\infty} K(x) = 0.$

- Examples of kernels

→ Gaussian $\underline{K(x)} = \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}},$

Uniform $\underline{K(x)} = \frac{1}{2} I(-1 \leq x \leq 1),$

Epanechnikov $\underline{K(x)} = \frac{3}{4} \cdot \underline{\max\{1 - x^2, 0\}}.$

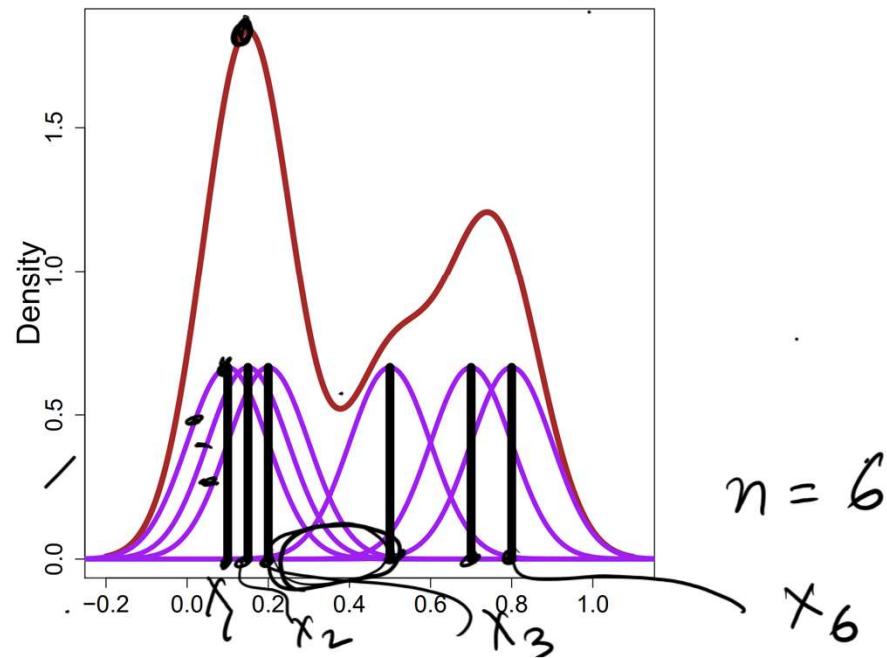


Kernel density estimator

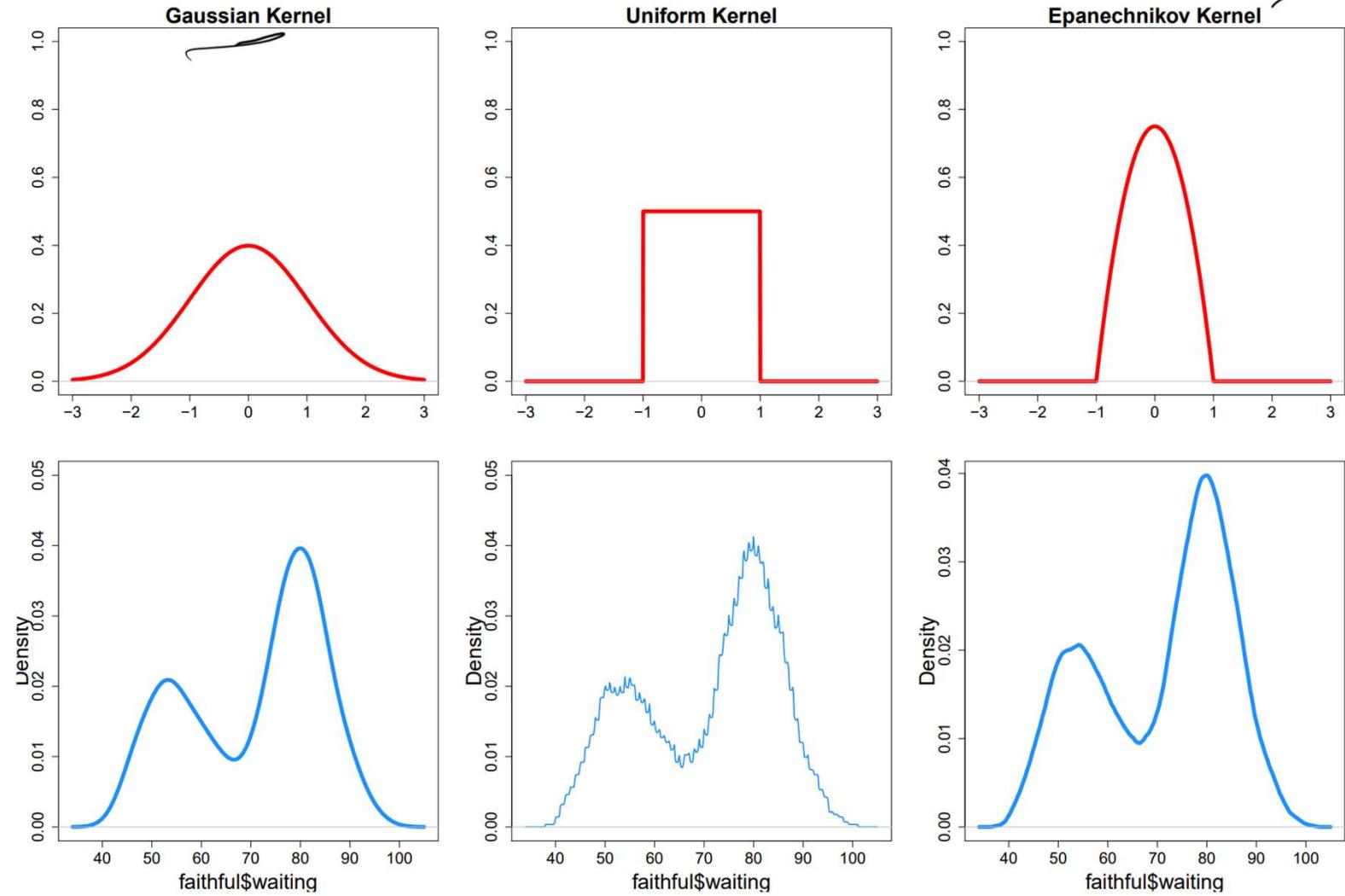
$$h = \frac{1}{M}$$

20.12 Definition. Given a kernel \underline{K} and a positive number \underline{h} , called the bandwidth, the kernel density estimator is defined to be

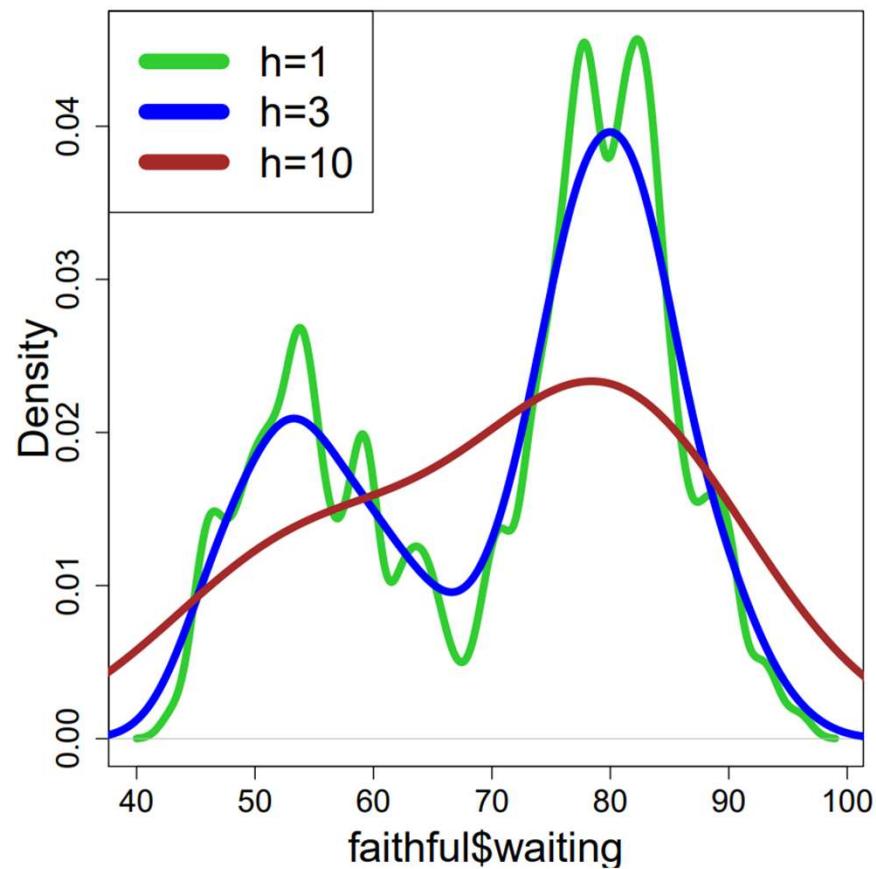
$$\widehat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right). \quad (20.21)$$

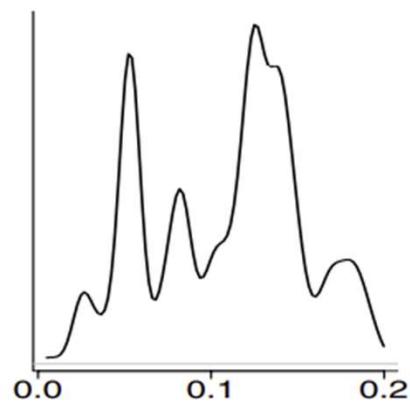


KDE with different kernels



Effect of kernel width





large h .

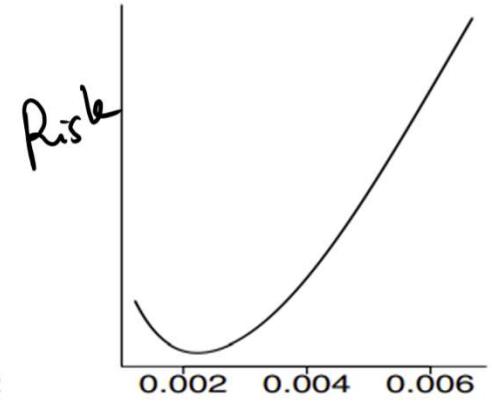
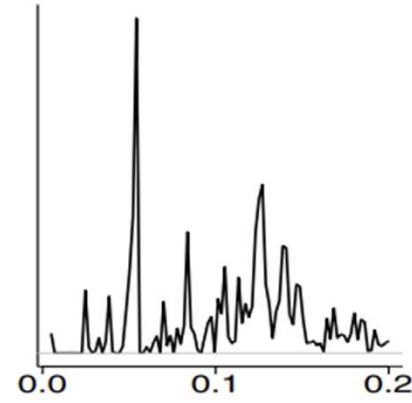
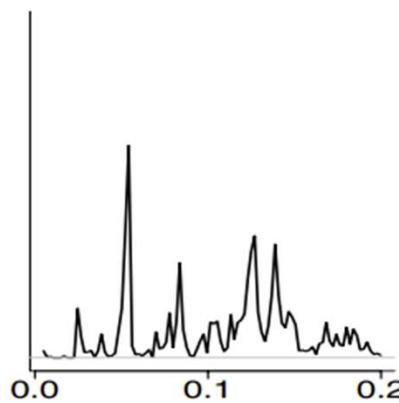


FIGURE 20.6. Kernel density estimators and estimated risk for the astronomy data.
Top left: oversmoothed. Top right: just right (bandwidth chosen by cross-validation).
Bottom left: undersmoothed. Bottom right: cross-validation curve as a function of
bandwidth h . The bandwidth was chosen to be the value of h where the curve is a
minimum.

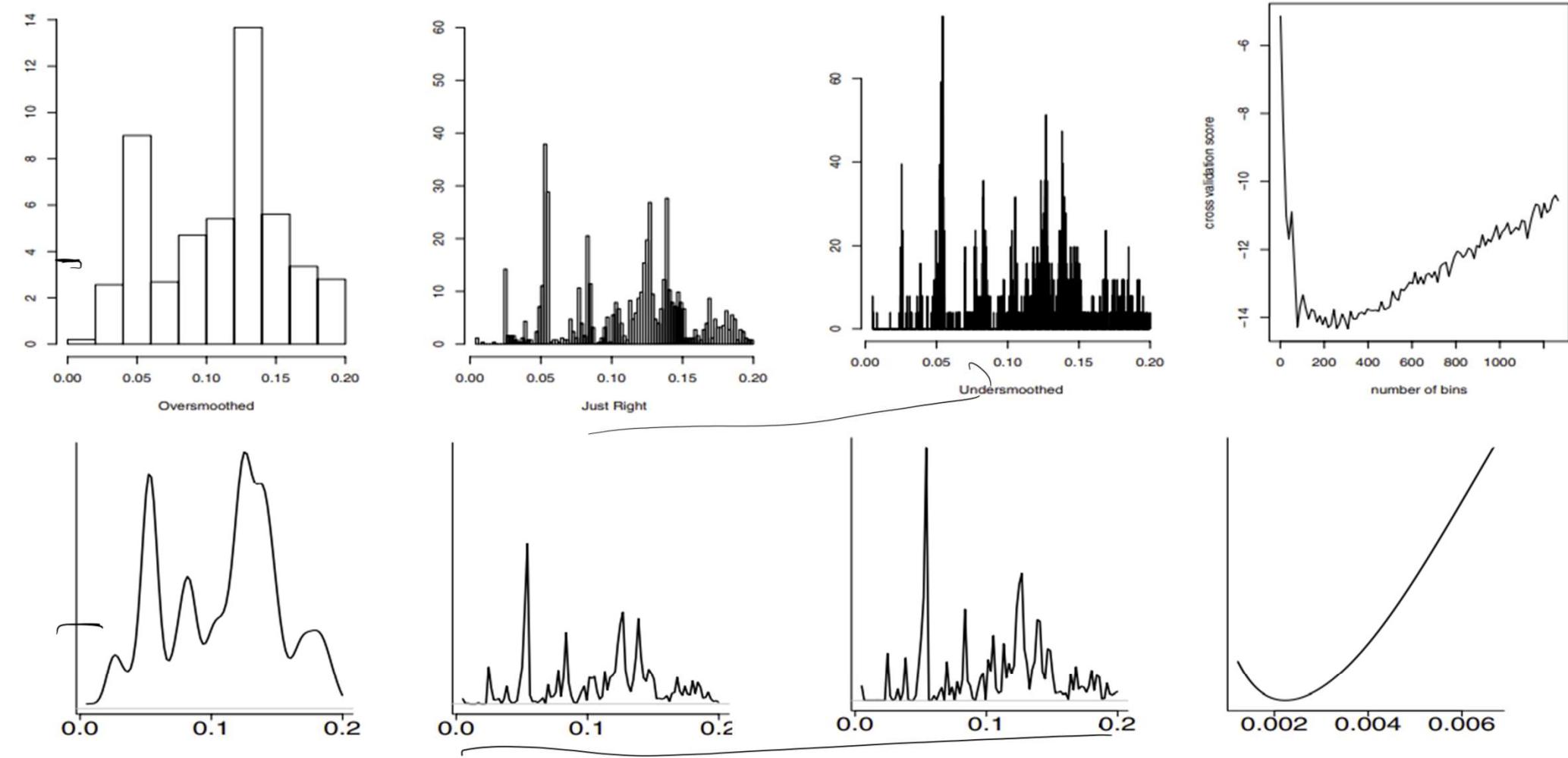


FIGURE 20.6. Kernel density estimators and estimated risk for the astronomy data. Top left: oversmoothed. Top right: just right (bandwidth chosen by cross-validation). Bottom left: undersmoothed. Bottom right: cross-validation curve as a function of bandwidth h . The bandwidth was chosen to be the value of h where the curve is a minimum.

Demo

[https://colab.research.google.com/github/fbeilstein/machine learning/
blob/master/lecture_15_kernel_density_estimation.ipynb#scrollTo=pGU9KIkxe-FY](https://colab.research.google.com/github/fbeilstein/machine_learning/blob/master/lecture_15_kernel_density_estimation.ipynb#scrollTo=pGU9KIkxe-FY)