

Exercise 12- nonparametric regression and density estimation

kernel regression density estimation

$X_1 \dots X_m \stackrel{iid}{\sim} f(x)$ estimate $f(x)$ by $\hat{f}_m(x) = \frac{1}{m h_m} \sum_{i=1}^m k\left(\frac{x-X_i}{h_m}\right)$

for $k: \mathbb{R} \rightarrow [0, \infty)$ a symmetric density, $h_m \xrightarrow{m \rightarrow \infty} 0$, $m h_m \xrightarrow{m \rightarrow \infty} \infty$.
 k kernel, h_m bandwidth
 bias $\hat{f}_m \rightarrow 0$ var $\hat{f}_m \rightarrow 0$

By theorem 17 then $\hat{f}_m(x) \xrightarrow{P} f(x)$ for all x . (under regularity assumptions)

In \mathbb{R} : density function

kernels are taken in the form as in Remark 28. $\tilde{k}(x) = \sqrt{\mu_k} k(\sqrt{\mu_k} x)$

for $\mu_k = \int y^2 k(y) dy = \text{var}(k)$

- for $k = \frac{1}{2} I[t \in [-1, 1]]$ $\mu_k = 1/3 \Rightarrow \tilde{k}(x) = \frac{1}{\sqrt{3} \cdot 2} I[\frac{1}{\sqrt{3}} t \in [-1, 1]]$
 $\text{var } \text{Unif}[-1, 1] = 1/3 \Rightarrow \text{let } \tilde{k} \sim \text{Unif}[-1, 1] \cdot \sqrt{3}$
 $= \frac{1}{2\sqrt{3}} I[t \in [-\sqrt{3}, \sqrt{3}]]$ so that $\mu_{\tilde{k}} = 1$

Bandwidth selection. usually based on asymptotic expansions

nod - normal reference - assume Gaussianity of f Section 9.2.1

$h_m = 1.06 m^{-1/5} \sigma$ estimated by $h_m = 1.06 m^{-1/5} \min\left\{S_m, \frac{IQR_m}{1.34}\right\}$

nod0 - should have a bit better robustness properties **default choice**

$h_m = 0.9 m^{-1/5} \sigma$

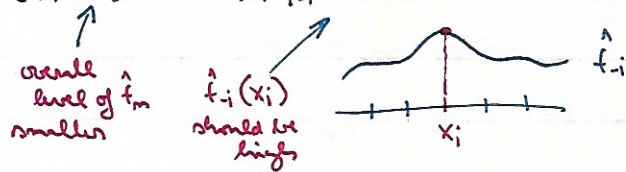
unbiased

ncv - (unbiased CV) choose h that minimizes an estimator of

$MISE(\hat{f}_m) = \int E[\hat{f}_m(x) - f(x)]^2 dx$ given g (Section 9.2.2) - up to a const.

$Z(h) = \int [\hat{f}_m(x)]^2 dx = \frac{2}{m} \sum_{i=1}^m \hat{f}_{-i}(x_i)$ with \hat{f}_{-i} the estimator

without x_i



bcv. - biased CV minimize the asymptotic ^{approximation for} $MISE(\hat{f}_m)$ (sect 9.3)

in that we have $R(f'') = \int [f''(x)]^2 dx$ - the overall curvature of f .

estimator $R(\hat{f}_m'')$ overestimates $R(f'')$ \rightarrow expansion leads to $R(\hat{f}_m'') - \frac{R(f'')}{m h_m^3}$



Boundary effects and

mirror-reflection: if $x_i \geq 0$ a.s. the estimator $\hat{f}_m(x)$ does not take into

account the support of x , i.e. $[0, \infty)$. Define

alternates the $\hat{f}_{MRE,m}(x) := \begin{cases} \hat{f}_m(x) + \hat{f}_m(-x) & \text{if } x \geq 0 \\ \hat{f}_m(x) & \text{if } x < 0 \end{cases}$

as optimal bw (global) - Section 9.2
 $h_m^{(opt)} = m^{-1/5} \left[\frac{R(f)}{R(f'') \mu_k^2} \right]^{1/5}$

R script

Multivariate kernel density estimation. $(\begin{smallmatrix} x_1 \\ y_1 \end{smallmatrix}) \dots (\begin{smallmatrix} x_m \\ y_m \end{smallmatrix}) \sim f_{X,Y}$

$$\hat{f}_m(x,y) := \frac{1}{m} \sum_{i=1}^m \frac{1}{|H|^{d/2}} k\left(H^{-1/2} \begin{bmatrix} x \\ y \end{bmatrix} - (\begin{smallmatrix} x_i \\ y_i \end{smallmatrix})\right)$$

for k a bivariate symmetric density and H a symmetric, pos. def. matrix

Bias reduction and higher order kernels. In the usual setting

$$E \hat{f}_m(x) = E \frac{1}{h} k\left(\frac{x-X}{h}\right) = \int_{\mathbb{R}} \frac{1}{h} k\left(\frac{x-y}{h}\right) f(y) dy = \int_{\mathbb{R}} L(x-y) f(y) dy$$

$L(\cdot) = \frac{1}{h} k\left(\frac{\cdot}{h}\right)$ is a density of $h \cdot Z$ for $Z \sim k$

convolution, density of $h \cdot Z + X$

for h fixed and m growing, we estimate the density of $hZ + X$

$$\hat{f}_m(x) = \frac{1}{m} \sum_{i=1}^m \frac{1}{h} k\left(\frac{x-X_i}{h}\right)$$

bias introducing kernel

Bias: $E \hat{f}_m(x) = \int_{\mathbb{R}} k(t) f(x-th) dt = f(x) \int k(t) dt - f'(x)h \int t k(t) dt + f(x-th) = f(x) + f'(x)(-ht) + \frac{f''(x)(-ht)^2}{2!} + \dots + \frac{f^{(p)}(x)(-ht)^p}{p!} + R_m$

0 for a sym kernel

$$\dots + \frac{f^{(p)}(x)(-h)^p}{p!} \int (t)^p k(t) dt \quad \text{typically bias } \hat{f}_m(x) = O(h^2)$$

But if also $\int t^j k(t) dt = 0 \quad \forall j=1,2,\dots,p-1$

and $\int t^p k(t) dt \neq 0$, we have bias $\hat{f}_m(x) = O(h^p)$

Though $\int t^2 k(t) dt = 0 \Rightarrow k \equiv 0$ if k is a symmetric density

and k would have to be chosen s.t. $k < 0$ may appear.

$$k(y) = \frac{1}{2} (3-y^2) \frac{1}{\sqrt{2\pi}} e^{-y^2/2}$$

$p=4$ modification of Gaussian kernel.

$$m(x) = E[Y|X=x]$$

Non-parametric regression

$(\begin{smallmatrix} x_1 \\ y_1 \end{smallmatrix}) \dots (\begin{smallmatrix} x_m \\ y_m \end{smallmatrix})$ iid

$$y_i = m(x_i) + \varepsilon_i \quad E[\varepsilon_i|X] = 0$$

estimate m by $\hat{m}(x) = \frac{\sum_{i=1}^m \frac{1}{h_m} k\left(\frac{x-x_i}{h_m}\right) \cdot y_i}{\sum_{i=1}^m \frac{1}{h_m} k\left(\frac{x-x_i}{h_m}\right)} = \sum_{i=1}^m w_i y_i$

$$w_i \geq 0 \\ \sum w_i = 1$$

more generally, estimate $m(x)$ by the intercept in regression

$$y_i \sim \beta_0 + \sum_{j=1}^p \beta_j (x_i - x)^j \quad \text{weighed by } k\left(\frac{x-x_i}{h_m}\right) \frac{1}{h_m}$$

\Rightarrow local polynomial regression

CV: minimize $\frac{1}{m} \sum_{i=1}^m [y_i - \hat{m}_{h,m}(x_i)]^2$



Mathematical script