

Chapter 1

Density Estimation

The estimation of probability density functions (PDFs) and cumulative distribution functions (CDFs) are cornerstones of applied data analysis in the social sciences. Testing for the equality of two distributions (or moments thereof) is perhaps the most basic test in all of applied data analysis. Economists, for instance, devote a great deal of attention to the study of income distributions and how they vary across regions and over time. Though the PDF and CDF are often the objects of direct interest, their estimation also serves as an important building block for other objects being modeled such as a conditional mean (i.e., a “regression function”), which may be directly modeled using nonparametric or semiparametric methods (a conditional mean is a function of a conditional PDF, which is itself a ratio of unconditional PDFs). After mastering the principles underlying the nonparametric estimation of a PDF, the nonparametric estimation of the workhorse of applied data analysis, the conditional mean function considered in Chapter 2, progresses in a fairly straightforward manner. Careful study of the approaches developed in Chapter 1 will be most helpful for understanding material presented in later chapters.

We begin with the estimation of a univariate PDF in Sections 1.1 through 1.3, turn to the estimation of a univariate CDF in Sections 1.4 and 1.5, and then move on to the more general multivariate setting in Sections 1.6 through 1.8. Asymptotic normality, uniform rates of convergence, and bias reduction methods appear in Sections 1.9 through 1.12. Numerous illustrative applications appear in Section 1.13, while theoretical and applied exercises can be found in Section 1.14

We now proceed with a discussion of how to estimate the PDF

$f_X(x)$ of a random variable X . For notational simplicity we drop the subscript X and simply use $f(x)$ to denote the PDF of X . Some of the treatments of the kernel estimation of a PDF discussed in this chapter are drawn from the two excellent monographs by Silverman (1986) and Scott (1992).

1.1 Univariate Density Estimation

To best appreciate why one might consider using nonparametric methods to estimate a PDF, we begin with an illustrative example, the parametric estimation of a PDF.

Example 1.1. Suppose X_1, X_2, \dots, X_n represent independent and identically distributed (i.i.d.) draws from a normal distribution with mean μ and variance σ^2 . We wish to estimate the normal PDF $f(x)$.

By assumption, $f(x)$ has a known parametric functional form (i.e., univariate normal) given by $f(x) = (2\pi\sigma^2)^{-1/2} \exp[-\frac{1}{2}(x - \mu)^2/\sigma^2]$, where the mean $\mu = E(X)$ and variance $\sigma^2 = E[(X - E(X))^2] = \text{var}(X)$ are the only unknown parameters to be estimated. One could estimate μ and σ^2 by the method of maximum likelihood as follows. Under the i.i.d. assumption, the joint PDF of (X_1, \dots, X_n) is simply the product of the univariate PDFs, which may be written as

$$f(X_1, \dots, X_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i - \mu)^2}{2\sigma^2}} = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2}.$$

Conditional upon the observed sample and taking the logarithm, this gives us the log-likelihood function

$$\begin{aligned} \mathcal{L}(\mu, \sigma^2) &\equiv \ln f(X_1, \dots, X_n; \mu, \sigma^2) \\ &= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2. \end{aligned}$$

The method of maximum likelihood proceeds by choosing those parameters that make it most likely that we observed the sample at hand given our distributional assumption. Thus, the likelihood function (or a monotonic transformation thereof, e.g., \ln) expresses the plausibility of different values of μ and σ^2 given the observed sample. We then maximize the likelihood function with respect to these two unknown parameters.

The necessary first order conditions for a maximization of the log-likelihood function are $\partial\mathcal{L}(\mu, \sigma^2)/\partial\mu = 0$ and $\partial\mathcal{L}(\mu, \sigma^2)/\partial\sigma^2 = 0$. Solving these first order conditions for the two unknown parameters μ and σ^2 yields

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{and} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2.$$

$\hat{\mu}$ and $\hat{\sigma}^2$ above are the maximum likelihood estimators of μ and σ^2 , respectively, and the resulting estimator of $f(x)$ is

$$\hat{f}(x) = \frac{1}{\sqrt{2\pi\hat{\sigma}^2}} \exp \left[-\frac{1}{2} \left(\frac{x - \hat{\mu}}{\hat{\sigma}} \right)^2 \right].$$

The “Achilles heel” of any parametric approach is of course the requirement that, prior to estimation, the analyst must specify the exact parametric functional form for the object being estimated. Upon reflection, the parametric approach is somewhat circular since we initially set out to estimate an unknown density but must first assume that the density is in fact known (up to a handful of unknown parameters, of course). Having based our estimate on the assumption that the density is a member of a known parametric family, we must then naturally confront the possibility that the parametric model is “mis-specified,” i.e., not consistent with the population from which the data was drawn. For instance, by assuming that X is drawn from a normally distributed population in the above example, we in fact impose a number of potentially quite restrictive assumptions: symmetry, unimodality, monotonically decreasing away from the mode and so on. If the true density were in fact asymmetric or possessed multiple modes, or was nonmonotonic away from the mode, then the presumption of distributional normality may provide a misleading characterization of the true density and could thereby produce erroneous estimates and lead to unsound inference.

At this juncture many readers will no doubt be pointing out that, having estimated a parametric PDF, one can always test whether the underlying distributional assumption is valid. We are, of course, completely sympathetic toward such arguments. Often, however, the rejection of a distributional assumption fails to provide any clear alternative. That is, we can reject the assumption of normality, but this rejection leaves us where we started, perhaps having ruled out but one of a large

number of candidate distributions. Against this backdrop, researchers might instead consider nonparametric approaches.

Nonparametric methods circumvent problems arising from the need to specify parametric functional forms prior to estimation. Rather than presume one knows the exact functional form of the object being estimated, one instead presumes that it satisfies some regularity conditions such as smoothness and differentiability. This does not, however, come without cost. By imposing less structure on the functional form of the PDF than do parametric methods, nonparametric methods require more data to achieve the same degree of precision as a *correctly specified* parametric model. Our primary focus in this text is on a class of estimators known as “nonparametric kernel estimators” (a “kernel function” is simply a weighting function), though in Chapters 14 and 15 we provide a treatment of alternative nonparametric methodologies including nearest neighbor and series methods.

Before proceeding to a formal theoretical analysis of nonparametric density estimation methods, we first consider a popular example of estimating the probability of a head on a toss of a coin which is closely related to the nonparametric estimation of a CDF. This in turn will lead us to the nonparametric estimation of a PDF.

Example 1.2. *Suppose we have a coin (perhaps an unfair one) and we want to estimate the probability of flipping the coin and having it land heads up. Let $p = P(H)$ denote the (unknown) population probability of obtaining a head. Taking a relative frequency approach, we would flip the coin n times, count the frequency of heads in n trials, and compute the relative frequency given by*

$$\hat{p} = \frac{1}{n} \{ \# \text{ of heads } \}, \quad (1.1)$$

which provides an estimate of p . The \hat{p} defined in (1.1) is often referred to as a “frequency estimator” of p , and it is also the maximum likelihood estimator of p (see Exercise 1.2). The estimator \hat{p} is, of course, fully nonparametric. Intuitively, one would expect that, if n is large, then \hat{p} should be “close” to p . Indeed, one can easily show that the mean squared error (MSE) of \hat{p} is given by (see Exercise 1.3)

$$\text{MSE}(\hat{p}) \stackrel{\text{def}}{=} E[(\hat{p} - p)^2] = \frac{p(1-p)}{n},$$

so $\text{MSE}(\hat{p}) \rightarrow 0$ as $n \rightarrow \infty$, which is termed as \hat{p} converges to p in mean square error; see Appendix A for the definitions of various modes of convergence.

We now discuss how to obtain an estimator of the CDF of X , which we denote by $F(x)$. The CDF is defined as

$$F(x) = P[X \leq x].$$

With i.i.d. data X_1, \dots, X_n (i.e., random draws from the distribution $F(\cdot)$), one can estimate $F(x)$ by

$$F_n(x) = \frac{1}{n} \{ \# \text{ of } X_i\text{'s} \leq x \}. \quad (1.2)$$

Equation (1.2) has a nice intuitive interpretation. Going back to our coin-flip example, if a coin is such that the probability of obtaining a head when we flip it equals $F(x)$ ($F(x)$ is unknown), and if we treat the collection of data X_1, \dots, X_n as flipping a coin n times and we say that a head occurs on the i^{th} trial if $X_i \leq x$, then $P(H) = P(X_i \leq x) = F(x)$. The familiar frequency estimator of $P(H)$ is equal to the number of heads divided by the number of trials:

$$\hat{P}(H) = \frac{\# \text{ of heads}}{n} = \frac{1}{n} \{ \# \text{ of } X_i\text{'s} \leq x \} \equiv F_n(x). \quad (1.3)$$

Therefore, we call (1.2) a frequency estimator of $F(x)$. Just as before when estimating $P(H)$, we expect intuitively that as n gets large, $\hat{P}(H)$ should yield a more accurate estimate of $P(H)$. By the same reasoning, one would expect that as $n \rightarrow \infty$, $F_n(x)$ yields a more accurate estimate of $F(x)$. Indeed, one can easily show that $F_n(x) \rightarrow F(x)$ in MSE, which implies that $F_n(x)$ converges to $F(x)$ in probability and also in distribution as $n \rightarrow \infty$. In Appendix A we introduce the concepts of convergence in mean square error, convergence in probability, convergence in distribution, and almost sure convergence. It is well established that $F_n(x)$ indeed converges to $F(x)$ in each of these various senses. These concepts of convergence are necessary as it is easy to show that the ordinary limit of $F_n(x)$ does not exist, i.e., $\lim_{n \rightarrow \infty} F_n(x)$ does not exist (see Exercise 1.3, where the definition of an ordinary limit is provided). This example highlights the necessity of introducing new concepts of convergence modes such as convergence in mean square error and convergence in probability.

Now we take up the question of how to estimate a PDF $f(x)$ without making parametric presumptions about its functional form. From the

definition of $f(x)$ we have¹

$$f(x) = \frac{d}{dx}F(x). \quad (1.4)$$

From (1.2) and (1.4), an obvious estimator of $f(x)$ is²

$$\hat{f}(x) = \frac{F_n(x+h) - F_n(x-h)}{2h}, \quad (1.5)$$

where h is a small positive increment.

By substituting (1.2) into (1.5), we obtain

$$\hat{f}(x) = \frac{1}{2nh} \{ \# \text{ of } X_1, \dots, X_n \text{ falling in the interval } [x-h, x+h] \}. \quad (1.6)$$

If we define a uniform kernel function given by

$$k(z) = \begin{cases} 1/2 & \text{if } |z| \leq 1 \\ 0 & \text{otherwise,} \end{cases} \quad (1.7)$$

then it is easy to see that $\hat{f}(x)$ given by (1.5) can also be expressed as

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{X_i - x}{h}\right). \quad (1.8)$$

Equation (1.8) is called a uniform kernel estimator because the kernel function $k(\cdot)$ defined in (1.7) corresponds to a uniform PDF. In general, we refer to $k(\cdot)$ as a kernel function and to h as a smoothing parameter (or, alternatively, a bandwidth or window width). Equation (1.8) is sometimes referred to as a “naïve” kernel estimator.

In fact one might use many other possible choices for the kernel function $k(\cdot)$ in this context. For example, one could use a standard normal kernel given by

$$k(v) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}v^2}, \quad -\infty < v < \infty. \quad (1.9)$$

This class of estimators can be found in the first published paper on kernel density estimation by Rosenblatt (1956), while Parzen (1962) established a number of properties associated with this class of estimators

¹We only consider the continuous X case in this chapter. We deal with the discrete X case in Chapters 3 and 4.

²Recall that the definition of the derivative of a function $g(x)$ is given by $dg(x)/dx = \lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h}$, or, equivalently, $dg(x)/dx = \lim_{h \rightarrow 0} \frac{g(x+h) - g(x-h)}{2h}$.

and relaxed the nonnegativity assumption in order to obtain estimators which are more efficient. For this reason, this approach is sometimes referred to as “Rosenblatt-Parzen kernel density estimation.”

We will prove shortly that the kernel estimator $\hat{f}(x)$ defined in (1.8) constructed from any general nonnegative bounded kernel function $k(\cdot)$ that satisfies

$$\begin{aligned} (i) \quad & \int k(v) dv = 1 \\ (ii) \quad & k(v) = k(-v) \\ (iii) \quad & \int v^2 k(v) dv = \kappa_2 > 0 \end{aligned} \tag{1.10}$$

is a consistent estimator of $f(x)$. Note that the symmetry condition (ii) implies that $\int vk(v) dv = 0$. By consistency, we mean that $\hat{f}(x) \rightarrow f(x)$ in probability (convergence in probability is defined in Appendix A). Note that $k(\cdot)$ defined in (1.10) is a (symmetric) PDF. For recent work on kernel methods with asymmetric kernels, see Abadir and Lawford (2004).

To define various modes of convergence, we first introduce the concept of the “Euclidean norm” (“Euclidean length”) of a vector. Given a $q \times 1$ vector $x = (x_1, x_2, \dots, x_q)' \in \mathbb{R}^q$, we use $\|x\|$ to denote the Euclidean length of x , which is defined by

$$\|x\| = [x'x]^{1/2} \equiv \sqrt{x_1^2 + x_2^2 + \dots + x_q^2}.$$

When $q = 1$ (a scalar), $\|x\|$ is simply the absolute value of x .

In the appendix we discuss the notation $O(\cdot)$ (“big Oh”) and $o(\cdot)$ (“small Oh”). Let a_n be a nonstochastic sequence. We say that $a_n = O(n^\alpha)$ if $|a_n| \leq Cn^\alpha$ for all n sufficiently large, where α and $C (> 0)$ are constants. Similarly, we say that $a_n = o(n^\alpha)$ if $a_n/n^\alpha \rightarrow 0$ as $n \rightarrow \infty$. We are now ready to prove the MSE consistency of $\hat{f}(x)$.

Theorem 1.1. *Let X_1, \dots, X_n denote i.i.d. observations having a three-times differentiable PDF $f(x)$, and let $f^{(s)}(x)$ denote the s th order derivative of $f(x)$ ($s = 1, 2, 3$). Let x be an interior point in the support of X , and let $\hat{f}(x)$ be that defined in (1.8). Assume that the kernel function $k(\cdot)$ is bounded and satisfies (1.10). Also, as $n \rightarrow \infty$, $h \rightarrow 0$ and $nh \rightarrow \infty$, then*

$$\begin{aligned} \text{MSE}(\hat{f}(x)) &= \frac{h^4}{4} [\kappa_2 f^{(2)}(x)]^2 + \frac{\kappa f(x)}{nh} + o(h^4 + (nh)^{-1}) \\ &= O(h^4 + (nh)^{-1}), \end{aligned} \tag{1.11}$$

where $\kappa_2 = \int v^2 k(v) dv$ and $\kappa = \int k^2(v) dv$.

Proof of Theorem 1.1.

$$\begin{aligned} \text{MSE}(\hat{f}(x)) &\equiv \mathbb{E} \left\{ \left[\hat{f}(x) - f(x) \right]^2 \right\} \\ &= \text{var}(\hat{f}(x)) + \left[\mathbb{E}(\hat{f}(x)) - f(x) \right]^2 \\ &\equiv \text{var}(\hat{f}(x)) + \left[\text{bias}(\hat{f}(x)) \right]^2. \end{aligned}$$

We will evaluate the $\text{bias}(\hat{f}(x))$ and $\text{var}(\hat{f}(x))$ terms separately.

For the bias calculation we will need to use the Taylor expansion formula. For a univariate function $g(x)$ that is m times differentiable, we have

$$\begin{aligned} g(x) = & g(x_0) + g^{(1)}(x_0)(x - x_0) + \frac{1}{2!}g^{(2)}(x_0)(x - x_0)^2 + \\ & \cdots + \frac{1}{(m-1)!}g^{(m-1)}(x_0)(x - x_0)^{m-1} + \frac{1}{m!}g^{(m)}(\xi)(x - x_0)^m, \end{aligned}$$

where $g^{(s)}(x_0) = \frac{\partial^s g(x)}{\partial x^s} \big|_{x=x_0}$, and ξ lies between x and x_0 .

The bias term is given by

$$\begin{aligned}
 \text{bias}(\hat{f}(x)) &= E \left\{ \frac{1}{nh} \sum_{i=1}^n k \left(\frac{X_i - x}{h} \right) \right\} - f(x) \\
 &= h^{-1} E \left[k \left(\frac{X_1 - x}{h} \right) \right] - f(x) \\
 &\quad (\text{by identical distribution}) \\
 &= h^{-1} \int f(x_1) k \left(\frac{x_1 - x}{h} \right) dx_1 - f(x) \\
 &= h^{-1} \int f(x + hv) k(v) h dv - f(x) \\
 &\quad (\text{change of variable, } x_1 - x = hv) \\
 &= \int \left\{ f(x) + f^{(1)}(x)hv + \frac{1}{2}f^{(2)}(x)h^2v^2 + O(h^3) \right\} k(v) dv \\
 &\quad - f(x) \\
 &= \left\{ f(x) + 0 + \frac{h^2}{2}f^{(2)}(x) \int v^2 k(v) dv + O(h^3) \right\} - f(x) \\
 &\quad \text{by (1.10)} \\
 &= \frac{h^2}{2}f^{(2)}(x) \int v^2 k(v) dv + O(h^3), \tag{1.12}
 \end{aligned}$$

where the $O(h^3)$ term comes from

$$(1/3!)h^3 \left| \int f^{(3)}(\tilde{x})v^3 k(v) dv \right| \leq Ch^3 \int |v^3 k(v)| dv = O(h^3),$$

where C is a positive constant, and where \tilde{x} lies between x and $x + hv$.

Note that in the above derivation we assume that $f(x)$ is three-times differentiable. We can weaken this condition to $f(x)$ being twice differentiable, resulting in $O(h^3)$ becomes $o(h^2)$, see Exercise 1.5)

$$\begin{aligned}
 \text{bias}(\hat{f}(x)) &= E(\hat{f}(x)) - f(x) \\
 &= \frac{h^2}{2}f^{(2)}(x) \int v^2 k(v) dv + o(h^2). \tag{1.13}
 \end{aligned}$$

Next we consider the variance term. Observe that

$$\begin{aligned}
 \text{var} \left(\hat{f}(x) \right) &= \text{var} \left[\frac{1}{nh} \sum_{i=1}^n k \left(\frac{X_i - x}{h} \right) \right] \\
 &= \frac{1}{n^2 h^2} \left\{ \sum_{i=1}^n \text{var} \left[k \left(\frac{X_i - x}{h} \right) \right] + 0 \right\} \\
 &\quad (\text{by independence}) \\
 &= \frac{1}{nh^2} \text{var} \left(k \left(\frac{X_1 - x}{h} \right) \right) \\
 &\quad (\text{by identical distribution}) \\
 &= \frac{1}{nh^2} \left\{ \text{E} \left[k^2 \left(\frac{X_1 - x}{h} \right) \right] - \left[\text{E} \left(k \left(\frac{X_1 - x}{h} \right) \right) \right]^2 \right\} \\
 &= \frac{1}{nh^2} \left\{ \int f(x_1) k^2 \left(\frac{x_1 - x}{h} \right) dx_1 \right. \\
 &\quad \left. - \left[\int f(x_1) k \left(\frac{x_1 - x}{h} \right) dx_1 \right]^2 \right\} \\
 &= \frac{1}{nh^2} \left\{ h \int f(x + hv) k^2(v) dv \right. \\
 &\quad \left. - \left[h \int f(x + hv) k(v) dv \right]^2 \right\} \\
 &= \frac{1}{nh^2} \left\{ h \int \left[f(x) + f^{(1)}(\xi) hv \right] k^2(v) dv - O(h^2) \right\} \\
 &= \frac{1}{nh} \left\{ f(x) \int k^2(v) dv + O \left(h \int |v| k^2(v) dv \right) - O(h) \right\} \\
 &= \frac{1}{nh} \{ \kappa f(x) + O(h) \}, \tag{1.14}
 \end{aligned}$$

where $\kappa = \int k^2(v) dv$.

Equations (1.12) and (1.14) complete the proof of Theorem 1.1. \square

Theorem 1.1 implies that (by Theorem A.7 of Appendix A)

$$\hat{f}(x) - f(x) = O_p \left(h^2 + (nh)^{-1/2} \right) = o_p(1).$$

By choosing $h = cn^{-1/\alpha}$ for some $c > 0$ and $\alpha > 1$, the conditions required for consistent estimation of $f(x)$, $h \rightarrow 0$ and $nh \rightarrow \infty$,

are clearly satisfied. The overriding question is what values of c and α should be used in practice. As can be seen, for a given sample size n , if h is small, the resulting estimator will have a small bias but a large variance. On the other hand, if h is large, then the resulting estimator will have a small variance but a large bias. To minimize $\text{MSE}(\hat{f}(x))$, one should balance the squared bias and the variance terms. The optimal choice of h (in the sense that $\text{MSE}(\hat{f}(x))$ is minimized) should satisfy $d\text{MSE}(\hat{f}(x))/dh = 0$. By using (1.11), it is easy to show that the optimal h that minimizes the leading term of $\text{MSE}(\hat{f}(x))$ is given by

$$h_{\text{opt}} = c(x)n^{-1/5}, \quad (1.15)$$

where $c(x) = \{\kappa f(x)/[\kappa_2 f^{(2)}(x)]^2\}^{1/5}$.

$\text{MSE}(\hat{f}(x))$ is clearly a “pointwise” property, and by using this as the basis for bandwidth selection we are obtaining a bandwidth that is optimal when estimating a density *at a point* x . Examining $c(x)$ in (1.15), we can see that a bandwidth which is optimal for estimation at a point x located in the tail of a distribution will differ from that which is optimal for estimation at a point located at, say, the mode. Suppose that we are interested not in tailoring the bandwidth to the pointwise estimation of $f(x)$ but instead in tailoring the bandwidth globally *for all points* x , that is, for all x in the support of $f(\cdot)$ (the support of x is defined as the set of points of x for which $f(x) > 0$, i.e., $\{x : f(x) > 0\}$). In this case we can choose h optimally by minimizing the “integrated MSE” (IMSE) of $\hat{f}(x)$. Using (1.11) we have

$$\begin{aligned} \text{IMSE}(\hat{f}) &\stackrel{\text{def}}{=} \int \mathbb{E} [\hat{f}(x) - f(x)]^2 dx = \frac{1}{4} h^4 \kappa_2^2 \int [f^{(2)}(x)]^2 dx \\ &\quad + \frac{\kappa}{nh} + o(h^4 + (nh)^{-1}). \end{aligned} \quad (1.16)$$

Again letting h_{opt} denote the optimal smoothing parameter that minimizes the leading terms of (1.16), we use simple calculus to get

$$h_{\text{opt}} = c_0 n^{-1/5}, \quad (1.17)$$

where $c_0 = \kappa_2^{-2/5} \kappa^{1/5} \left\{ \int [f^{(2)}(x)]^2 dx \right\}^{-1/5} > 0$ is a positive constant. Note that if $f^{(2)}(x) = 0$ for (almost) all x , then c_0 is not well defined. For example, if X is, say, uniformly distributed over its support, then $f^{(s)}(x) = 0$ for all x and for all $s \geq 1$, and (1.17) is not defined in this case. It can be shown that in this case (i.e., when X is uniformly

distributed), h_{opt} will have a different rate of convergence equal to $n^{-1/3}$; see the related discussion in Section 1.3.1 and Exercise 1.16.

An interesting extension of the above results can be found in Zinde-Walsh (2005), who examines the asymptotic process for the kernel density estimator by means of generalized functions and generalized random processes and presents novel results for characterizing the behavior of kernel density estimators when the density does not exist, i.e., when the density does not exist as a locally summable function.

1.2 Univariate Bandwidth Selection: Rule-of-Thumb and Plug-In Methods

Equation (1.17) reveals that the optimal smoothing parameter depends on the integrated second derivative of the unknown density through c_0 . In practice, one might choose an initial “pilot value” of h to estimate $\int [f^{(2)}(x)]^2 dx$ nonparametrically, and then use this value to obtain h_{opt} using (1.17). Such approaches are known as “plug-in methods” for obvious reasons. One popular way of choosing the initial h , suggested by Silverman (1986), is to assume that $f(x)$ belongs to a parametric family of distributions, and then to compute h using (1.17). For example, if $f(x)$ is a normal PDF with variance σ^2 , then $\int [f^{(2)}(x)]^2 dx = 3/[8\pi^{1/2}\sigma^5]$. If a standard normal kernel is used, using (1.17), we get the pilot estimate

$$h_{\text{pilot}} = (4\pi)^{-1/10} \left[(3/8)\pi^{-1/2} \right]^{-1/5} \sigma n^{-1/5} \approx 1.06\sigma n^{-1/5}, \quad (1.18)$$

which is then plugged into $\int [\hat{f}^{(2)}(x)]^2 dx$, which then may be used to obtain h_{opt} using (1.17). A clearly undesirable property of the plug-in method is that it is not fully automatic because one needs to choose an initial value of h to estimate $\int [f^{(2)}(x)]^2 dx$ (see Marron, Jones and Sheather (1996) and also Loader (1999) for further discussion).

Often, practitioners will use (1.18) itself for the bandwidth. This is known as the “normal reference rule-of-thumb” approach since it is the optimal bandwidth for a particular family of distributions, in this case the normal family. Should the underlying distribution be “close” to a normal distribution, then this will provide good results, and for exploratory purposes it is certainly computationally attractive. In practice, σ is replaced by the sample standard deviation of $\{X_i\}_{i=1}^n$, while Silverman (1986, p. 47) advocates using a more robust measure

of spread which replaces σ with A , an “adaptive” measure of spread given by

$$A = \min(\text{standard deviation}, \text{interquartile range}/1.34).$$

We now turn our attention to a discussion of a number of fully automatic or “data-driven” methods for selecting h that are tailored to the sample at hand.

1.3 Univariate Bandwidth Selection: Cross-Validation Methods

In both theoretical and practical settings, nonparametric kernel estimation has been established as relatively insensitive to choice of kernel function. However, the same cannot be said for bandwidth selection. Different bandwidths can generate radically differing impressions of the underlying distribution. If kernel methods are used simply for “exploratory” purposes, then one might undersmooth the density by choosing a small value of h and let the eye do any remaining smoothing. Alternatively, one might choose a range of values for h and plot the resulting estimates. However, for sound analysis and inference, a principle having some known optimality properties must be adopted. One can think of choosing the bandwidth as being analogous to choosing the number of terms in a series approximation; the more terms one includes in the approximation, the more flexible the resulting model becomes, while the smaller the bandwidth of a kernel estimator, the more flexible it becomes. However, increasing flexibility (reducing potential bias) necessarily leads to increased variability (increasing potential variance). Seen in this light, one naturally appreciates how a number of methods discussed below are motivated by the need to balance the squared bias and variance of the resulting estimate.

1.3.1 Least Squares Cross-Validation

Least squares cross-validation is a fully automatic data-driven method of selecting the smoothing parameter h , originally proposed by Rudemo (1982), Stone (1984) and Bowman (1984) (see also Silverman (1986, pp. 48-51)). This method is based on the principle of selecting a bandwidth that minimizes the integrated squared error of the resulting estimate, that is, it provides an optimal bandwidth tailored to *all* x in the support of $f(x)$.

The integrated squared difference between \hat{f} and f is

$$\int [\hat{f}(x) - f(x)]^2 dx = \int \hat{f}(x)^2 dx - 2 \int \hat{f}(x)f(x) dx + \int f(x)^2 dx. \quad (1.19)$$

As the third term on the right-hand side of (1.19) is unrelated to h , choosing h to minimize (1.19) is therefore equivalent to minimizing

$$\int \hat{f}(x)^2 dx - 2 \int \hat{f}(x)f(x) dx \quad (1.20)$$

with respect to h . In the second term, $\int \hat{f}(x)f(x) dx$ can be written as $E_X[\hat{f}(X)]$, where $E_X(\cdot)$ denotes expectation with respect to X and not with respect to the random observations $\{X_j\}_{j=1}^n$ used for computing $\hat{f}(\cdot)$. Therefore, we may estimate $E_X[\hat{f}(X)]$ by $n^{-1} \sum_{i=1}^n \hat{f}_{-i}(X_i)$ (i.e., replacing E_X by its sample mean), where

$$\hat{f}_{-i}(X_i) = \frac{1}{(n-1)h} \sum_{j=1, j \neq i}^n k\left(\frac{X_i - X_j}{h}\right) \quad (1.21)$$

is the leave-one-out kernel estimator of $f(X_i)$.³ Finally, we estimate the first term $\int \hat{f}(x)^2 dx$ by

$$\begin{aligned} \int \hat{f}(x)^2 dx &= \frac{1}{n^2 h^2} \sum_{i=1}^n \sum_{j=1}^n \int k\left(\frac{X_i - x}{h}\right) k\left(\frac{X_j - x}{h}\right) dx \\ &= \frac{1}{n^2 h} \sum_{i=1}^n \sum_{j=1}^n \bar{k}\left(\frac{X_i - X_j}{h}\right), \end{aligned} \quad (1.22)$$

where $\bar{k}(v) = \int k(u)k(v-u) du$ is the twofold convolution kernel derived from $k(\cdot)$. If $k(v) = \exp(-v^2/2)/\sqrt{2\pi}$, a standard normal kernel, then $\bar{k}(v) = \exp(-v^2/4)/\sqrt{4\pi}$, a normal kernel (i.e., normal PDF) with mean zero and variance two, which follows since two independent $N(0, 1)$ random variables sum to a $N(0, 2)$ random variable.

³Here we emphasize that it is important to use the leave-one-out kernel estimator for computing $E_X(\cdot)$ above. This is because the expectations operator presumes that the X and the X_j 's are independent of one another. Without using the leave-one-out estimator, the cross-validation method will break down; see Exercise 1.6 (iii).

Least squares cross-validation therefore chooses h to minimize

$$CV_f(h) = \frac{1}{n^2 h} \sum_{i=1}^n \sum_{j=1}^n \bar{k} \left(\frac{X_i - X_j}{h} \right) - \frac{2}{n(n-1)h} \sum_{i=1}^n \sum_{j \neq i, j=1}^n k \left(\frac{X_i - X_j}{h} \right), \quad (1.23)$$

which is typically undertaken using numerical search algorithms.

It can be shown that the leading term of $CV_f(h)$ is CV_{f0} given by (ignoring a term unrelated to h ; see Exercise 1.6)

$$CV_{f0}(h) = B_1 h^4 + \frac{\kappa}{nh}, \quad (1.24)$$

where $B_1 = (\kappa_2^2/4) \int [f^{(2)}(x)]^2 dx$ ($\kappa_2 = \int v^2 k(v) dv$, $\kappa = \int k^2(v) dv$). Thus, as long as $f^{(2)}(x)$ does not vanish for (almost) all x , we have $B_1 > 0$.

Let h^0 denote the value of h that minimizes CV_{f0} . Simple calculus shows that $h^0 = c_0 n^{-1/5}$ where

$$c_0 = [\kappa/(4B_1)]^{1/5} = \kappa^{1/5} \kappa_2^{-2/5} \left\{ \left[\int f^{(2)}(x) \right]^2 dx \right\}^{-1/5}.$$

A comparison of h^0 with h_{opt} in (1.17) reveals that the two are identical, i.e., $h^0 \equiv h_{\text{opt}}$. This arises because h_{opt} minimizes $\int E[\hat{f}(x) - f(x)]^2 dx$, while h^0 minimizes $E[CV_f(h)]$, the leading term of $CV_f(h)$. It can be easily seen that $E[CV_f(h)] + \int f(x)^2 dx$ is an alternative version of $\int E[\hat{f}(x) - f(x)]^2 dx$; hence, $E[CV_f(h)] + \int f(x)^2 dx$ also estimates $\int E[\hat{f}(x) - f(x)]^2 dx$. Given that $\int f(x)^2 dx$ is unrelated to h , one would expect that h^0 and h_{opt} should be the same.

Let \hat{h} denote the value of h that minimizes $CV_f(h)$. Given that $CV_f(h) = CV_{f0} + (s.o.)$, where $(s.o.)$ denotes smaller order terms (than CV_{f0}) and terms unrelated to h , it can be shown that $\hat{h} = h^0 + o_p(h^0)$, or, equivalently, that

$$\frac{\hat{h} - h^0}{h^0} \equiv \frac{\hat{h}}{h^0} - 1 \rightarrow 0 \text{ in probability.} \quad (1.25)$$

Intuitively, (1.25) is easy to understand because $CV_f(h) = CV_{f0}(h) + (s.o.)$, thus asymptotically an h that minimizes $CV_f(h)$ should be

close to an h that minimizes $CV_{f_0}(h)$; therefore, we expect that \hat{h} and h^0 will be close to each other in the sense of (1.25). Härdle, Hall and Marron (1988) showed that $(\hat{h} - h^0)/h^0 = O_p(n^{-1/10})$, which indeed converges to zero (in probability) but at an extremely slow rate.

We again underscore the need to use the leave-one-out kernel estimator when constructing CV_f as given in (1.23). If instead one were to use the standard kernel estimator, least squares cross-validation will break down, yielding $\hat{h} = 0$. Exercise 1.6 shows that if one does not use the leave-one-out kernel estimator when estimating $f(X_i)$, then $h = 0$ minimizes the objective function, which of course violates the consistency condition that $nh \rightarrow \infty$ as $n \rightarrow \infty$.

Here we implicitly impose the restriction that $f^{(2)}(x)$ is not a zero function, which rules out the case for which $f(x)$ is a uniform PDF. In fact this condition can be relaxed. Stone (1984) showed that, as long as $f(x)$ is bounded, then the least squares cross-validation method will select h optimally in the sense that

$$\frac{\int [\hat{f}(x, \hat{h}) - f(x)]^2 dx}{\inf_h \int [\hat{f}(x, h) - f(x)]^2 dx} \rightarrow 1 \text{ almost surely,} \quad (1.26)$$

where $\hat{f}(x, \hat{h})$ denotes the kernel estimator of $f(x)$ with cross-validation selected \hat{h} , and $\hat{f}(x, h)$ is the kernel estimator with a generic h . Obviously, the ratio defined in (1.26) should be greater than or equal to one for any n . Therefore, Stone's (1984) result states that, asymptotically, cross-validated smoothing parameter selection is optimal in the sense of minimizing the estimation integrated square error. In Exercise 1.16 we further discuss the intuition underlying why $\hat{h} \rightarrow 0$ even when $f(x)$ is a uniform PDF.

1.3.2 Likelihood Cross-Validation

Likelihood cross-validation is another automatic data-driven method for selecting the smoothing parameter h . This approach yields a density estimate which has an entropy theoretic interpretation, since the estimate will be close to the actual density in a Kullback-Leibler sense. This approach was proposed by Duin (1976).

Likelihood cross-validation chooses h to maximize the (leave-one-out) log likelihood function given by

$$\mathcal{L} = \ln L = \sum_{i=1}^n \ln \hat{f}_{-i}(X_i),$$

where $\hat{f}_{-i}(X_i)$ is the leave-one-out kernel estimator of $f(X_i)$ defined in (1.21). The main problem with likelihood cross-validation is that it is severely affected by the tail behavior of $f(x)$ and can lead to inconsistent results for fat tailed distributions when using popular kernel functions (see Hall (1987*a*, 1987*b*)). For this reason the likelihood cross-validation method has elicited little interest in the statistical literature.

However, the likelihood cross-validation method may work well for a range of standard distributions (i.e., thin tailed). We consider the performance of likelihood cross-validation in Section 1.3.3, when we compare the impact of different bandwidth selection methods on the resulting density estimate, and in Section 1.13, where we consider empirical applications.

1.3.3 An Illustration of Data-Driven Bandwidth Selection

Figure 1.1 presents kernel estimates constructed from $n = 500$ observations drawn from a simulated bimodal distribution. The second order Gaussian (normal) kernel was used throughout, and least squares cross-validation was used to select the bandwidth for the estimate appearing in the upper left plot of the figure, with $h_{\text{lscv}} = 0.19$. We also plot the estimate based on the normal reference rule-of-thumb ($h_{\text{ref}} = 0.34$) along with an undersmoothed estimate ($1/5 \times h_{\text{lscv}}$) and an oversmoothed estimate ($5 \times h_{\text{lscv}}$).⁴

Figure 1.1 reveals that least squares cross-validation appears to yield a reasonable density estimate for this data, while the reference rule-of-thumb is inappropriate as it oversmooths somewhat. Extreme oversmoothing can lead to a unimodal estimate which completely obscures the true bimodal nature of the underlying distribution. Also, undersmoothing leads to too many false modes. See Exercise 1.17 for an empirical application that investigates the effects of under- and over-smoothing on the resulting density estimate.

1.4 Univariate CDF Estimation

In Section 1.1 we introduced the empirical CDF estimator $F_n(x)$ given in (1.2), while Exercise 1.4 shows that it is a \sqrt{n} -consistent estimator

⁴Likelihood cross-validation yielded a bandwidth of $h_{\text{mlcv}} = 0.15$, which results in a density estimate virtually identical to that based upon least squares cross-validation for this dataset.

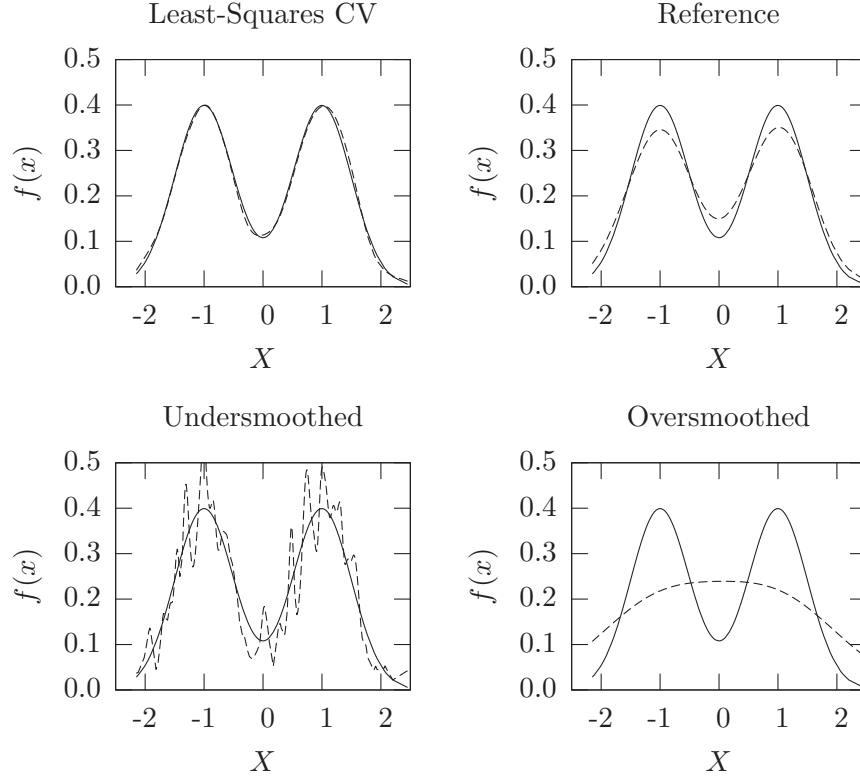


Figure 1.1: Univariate kernel estimates of a mixture of normals using least squares cross-validation, the normal reference rule-of-thumb, undersmoothing, and oversmoothing ($n = 500$). The correct parametric data generating process appears as the solid line, the kernel estimate as the dashed line.

of $F(x)$. However, this empirical CDF $F_n(x)$ is not smooth as it jumps by $1/n$ at each sample realization point. One can, however, obtain a smoothed estimate of $F(x)$ by integrating $\hat{f}(x)$. Define

$$\hat{F}(x) = \int_{-\infty}^x \hat{f}(v) dv = \frac{1}{n} \sum_{i=1}^n G\left(\frac{x - X_i}{h}\right), \quad (1.27)$$

where $G(x) = \int_{-\infty}^x k(v) dv$ is a CDF (which follows directly because $k(\cdot)$ is a PDF; see (1.10)). The next theorem provides the MSE of $\hat{F}(x)$.

Theorem 1.2. *Under conditions given in Bowman, Hall and Prvan (1998), in particular, assuming that $F(x)$ is twice continuously differentiable, $k(v) = dG(v)/dv$ is bounded, symmetric, and compactly supported, and that $d^2F(x)/dx^2$ is Hölder-continuous, $0 \leq h \leq Cn^{-\epsilon}$ for some $0 < \epsilon < \frac{1}{8}$, then as $n \rightarrow \infty$,*

$$\begin{aligned} \text{MSE}(\hat{F}) &= \mathbb{E} \left[\hat{F}(x) - F(x) \right]^2 \\ &= c_0(x)n^{-1} - c_1(x)hn^{-1} + c_2(x)h^4 + o(h^4 + hn^{-1}), \end{aligned}$$

where $c_0 = F(x)(1 - F(x))$, $c_1(x) = \alpha_0 f(x)$, $\alpha_0 = 2 \int vG(v)k(v) dv$, $f(x) = dF(x)/dx$, $c_2(x) = [(\kappa_2/2)F^{(2)}(x)]^2$, $\kappa_2 = \int v^2k(v) dv$, and where $F^{(s)}(x) = d^sF(x)/dx^s$ is the s th derivative of $F(x)$.

Proof. Note that $\mathbb{E} \left[\hat{F}(x) \right] = \mathbb{E} \left[G \left(\frac{x - X_i}{h} \right) \right]$. Then we have ($\int = \int_{-\infty}^{\infty}$)

$$\begin{aligned} \mathbb{E} \left[G \left(\frac{x - X_i}{h} \right) \right] &= \int G \left(\frac{x - z}{h} \right) f(z) dz \\ &= h \int G(v) f(x - hv) dv = - \int G(v) dF(x - hv) \\ &= - [G(v)F(x - hv)] \Big|_{v=-\infty}^{v=\infty} + \int k(v) F(x - hv) dv \\ &= \int k(v) \left[F(x) - F^{(1)}(x)hv + (1/2)h^2F^{(2)}(x)v^2 \right] dv \\ &\quad + o(h^2) \\ &= F(x) + (1/2)\kappa_2h^2F^{(2)}(x) + o(h^2), \end{aligned} \tag{1.28}$$

where at the second equality above we used

$$- \int_{-\infty}^{\infty} [\dots] dv = \int_{-\infty}^{\infty} [\dots] dv.$$

Also note that we did not use a Taylor expansion in $\int G(v)F(x - hv) dv$ since $\int v^m G(v) dv = +\infty$ for any $m \geq 0$. We first used integration by parts to get $k(v)$, and then used the Taylor expansion since $\int v^m k(v) dv$ is usually finite. For example, if $k(v)$ has bounded support or $k(v)$ is a standard normal kernel function, then $\int v^m k(v) dv$ is finite for any $m \geq 0$.

Similarly,

$$\begin{aligned}
 \mathbb{E} \left[G^2 \left(\frac{x - X_i}{h} \right) \right] &= \int G^2 \left(\frac{x - z}{h} \right) f(z) dz = h \int G^2(v) f(x - hv) dv \\
 &= - \int G^2(v) dF(x - hv) \\
 &= 2 \int G(v) k(v) F(x - hv) dv \\
 &= 2 \int G(v) k(v) [F(x) - F^{(1)}(x)hv] dv + O(h^2) \\
 &= F(x) - \alpha_0 h f(x) + O(h^2),
 \end{aligned} \tag{1.29}$$

where $\alpha_0 = 2 \int v G(v) k(v) dv$, and where we have used the fact that

$$2 \int_{-\infty}^{\infty} G(v) k(v) dv = \int_{-\infty}^{\infty} dG^2(v) = G^2(\infty) - G^2(-\infty) = 1,$$

because $G(\cdot)$ is a (user-specified) CDF kernel function.

From (1.28) we have $\text{bias}[\hat{F}(x)] = (1/2)\kappa_2 h^2 F^{(2)}(x) + o(h^2)$, and from (1.28) and (1.29) we have

$$\begin{aligned}
 \text{var} [\hat{F}(x)] &= n^{-1} \text{var} \left[G \left(\frac{x - X_i}{h} \right) \right] \\
 &= n^{-1} \left\{ \mathbb{E} \left[G^2 \left(\frac{x - X_i}{h} \right) \right] - \left[\mathbb{E} G \left(\frac{x - X_i}{h} \right) \right]^2 \right\} \\
 &= n^{-1} F(x) [1 - F(x)] - \alpha_0 f(x) h n^{-1} + o(h/n).
 \end{aligned}$$

Hence,

$$\begin{aligned}
 \mathbb{E} \left(\hat{F}(x) - F(x) \right)^2 &= \left[\text{bias} \left(\hat{F}(x) \right) \right]^2 + \text{var} [\hat{F}(x)] \\
 &= n^{-1} F(x) [1 - F(x)] + h^4 (\kappa_2/2)^2 \left[F^{(2)}(x) \right]^2 \\
 &\quad - \alpha_0 f(x) \frac{h}{n} + o(h^4 + n^{-1}h).
 \end{aligned} \tag{1.30}$$

This completes the proof of Theorem 1.2. \square

From Theorem 1.2 we immediately obtain the following result on the IMSE of \hat{F} :

$$\begin{aligned} \text{IMSE}(\hat{F}) &= \int \mathbb{E} \left[\hat{F}(x) - F(x) \right]^2 dx \\ &= C_0 n^{-1} - C_1 h n^{-1} + C_2 h^4 + o(h^4 + h n^{-1}), \end{aligned} \quad (1.31)$$

where $C_j = \int c_j(x) dx$ ($j = 0, 1, 2$). Letting h_0 denote the value of h that minimizes the leading term of IMSE, we obtain

$$h_0 = a_0 n^{-1/3},$$

where $a_0 = [C_1/(4C_2)]^{1/3}$, hence the optimal smoothing parameter for estimating univariate a CDF has a faster rate of convergence than the optimal smoothing parameter for estimating a univariate PDF ($n^{-1/3}$ versus $n^{-1/5}$). With $h \sim n^{-1/3}$, we have $h^2 = O(n^{-2/3}) = o(n^{-1/2})$. Hence, $\sqrt{n}[\hat{F}(x) - F(x)] \rightarrow N(0, F(x)[1 - F(x)])$ in distribution by the Liapunov central limit theorem (CLT); see Theorem A.5 in Appendix A for this and a range of other useful CLTs.

As is the case for nonparametric PDF estimation, nonparametric CDF estimation has widespread potential application though it is not nearly as widely used. For instance, it can be used to test stochastic dominance without imposing parametric assumptions on the underlying CDFs; see, e.g., Barrett and Donald (2003) and Linton, Whang and Maasoumi (2005).

1.5 Univariate CDF Bandwidth Selection: Cross-Validation Methods

Bowman et al. (1998) suggest choosing h for $\hat{F}(x)$ by minimizing the following cross-validation function:

$$CV_F(h) = \frac{1}{n} \sum_{i=1}^n \int \left\{ \mathbf{1}(X_i \leq x) - \hat{F}_{-i}(x) \right\}^2 dx, \quad (1.32)$$

where $\hat{F}_{-i}(x) = (n-1)^{-1} \sum_{j \neq i}^n G\left(\frac{x - X_j}{h}\right)$ is the leave-one-out estimator of $F(x)$.

Bowman et al. (1998) show that $CV_F = \mathbb{E}[CV_F] + (s.o.)$ and that

(see Exercise 1.9)

$$\begin{aligned} E[CV_F(h)] &= \int F(1-F) dx + \frac{1}{n-1} \int F(1-F) dx - C_1 h n^{-1} \\ &\quad + C_2 h^4 + o(h n^{-1} + h^4). \end{aligned} \quad (1.33)$$

We observe that (1.33) has the same leading term as $\text{IMSE}(\hat{F})$ given in (1.31). Thus, asymptotically, selecting h via cross-validation leads to the same asymptotic optimality property for $\hat{F}(x)$ that would arise when using h_0 , the optimal deterministic smoothing parameter. If we let \hat{h} denote the cross-validated smoothing parameter, then it can be shown that $\hat{h}/h_0 \rightarrow 1$ in probability. Note that when using \hat{h} , the asymptotic distribution of $\hat{F}(x, \hat{h})$ is the same as $\hat{F}(x, h_0)$ (by using a stochastic equicontinuity argument as outlined in Appendix A), that is,

$$\sqrt{n} \left(\hat{F}(x) - F(x) \right) \xrightarrow{d} N(0, F(x)(1-F(x))), \quad (1.34)$$

where $\hat{F}(x)$ is defined in (1.27) with h replaced by \hat{h} . Note that no bias term appears in (1.34) since $\text{bias}(\hat{F}(x)) = O(h_0^2) = O(n^{-2/3}) = o(n^{-1/2})$, which was not the case for PDF estimation. Here the squared bias term has order smaller than the leading variance term of $O(n^{-1})$ (i.e., $\text{var}(\hat{F}(x)) = O(n^{-1})$).

We now turn our attention to a generalization of the univariate kernel estimators developed above, namely multivariate kernel estimators. Again, we consider only the continuous case in this chapter; we tackle discrete and mixed continuous and discrete data cases in Chapters 3 and 4.

1.6 Multivariate Density Estimation

Suppose that X_1, \dots, X_n constitute an i.i.d. q -vector ($X_i \in \mathbb{R}^q$, for some $q > 1$) having a common PDF $f(x) = f(x_1, x_2, \dots, x_q)$. Let X_{is} denote the s th component of X_i ($s = 1, \dots, q$). Using a “product kernel function” constructed from the product of univariate kernel functions, we estimate the PDF $f(x)$ by

$$\hat{f}(x) = \frac{1}{n h_1 \dots h_q} \sum_{i=1}^n K \left(\frac{X_i - x}{h} \right), \quad (1.35)$$

where $K\left(\frac{X_i - x}{h}\right) = k\left(\frac{X_{i1} - x_1}{h_1}\right) \times \cdots \times k\left(\frac{X_{iq} - x_q}{h_q}\right)$, and where $k(\cdot)$ is a univariate kernel function satisfying (1.10).

The proof of MSE consistency of $\hat{f}(x)$ is similar to the univariate case. In particular, one can show that

$$\text{bias}\left(\hat{f}(x)\right) = \frac{\kappa_2}{2} \sum_{s=1}^q h_s^2 f_{ss}(x) + O\left(\sum_{s=1}^q h_s^3\right), \quad (1.36)$$

where $f_{ss}(x)$ is the second order derivative of $f(x)$ with respect to x_s , $\kappa_2 = \int v^2 k(v) dv$, and one can also show that

$$\text{var}\left(\hat{f}(x)\right) = \frac{1}{nh_1 \dots h_q} \left[\kappa^q f(x) + O\left(\sum_{s=1}^q h_s^2\right) \right] = O\left(\frac{1}{nh_1 \dots h_q}\right), \quad (1.37)$$

where $\kappa = \int k^2(v) dv$. The proofs of (1.36) and (1.37), which are similar to the univariate X case, are left as an exercise (see Exercise 1.11).

Summarizing, we obtain the result

$$\begin{aligned} \text{MSE}\left(\hat{f}(x)\right) &= \left[\text{bias}\left(\hat{f}(x)\right)\right]^2 + \text{var}\left(\hat{f}(x)\right) \\ &= O\left(\left(\sum_{s=1}^q h_s^2\right)^2 + (nh_1 \dots h_q)^{-1}\right). \end{aligned}$$

Hence, if as $n \rightarrow \infty$, $\max_{1 \leq s \leq q} h_s \rightarrow 0$ and $nh_1 \dots h_q \rightarrow \infty$, then we have $\hat{f}(x) \rightarrow f(x)$ in MSE, which implies that $\hat{f}(x) \rightarrow f(x)$ in probability.

As we saw in the univariate case, the optimal smoothing parameters h_s should balance the squared bias and variance terms, i.e., $h_s^4 = O((nh_1 \dots h_q)^{-1})$ for all s . Thus, we have $h_s = c_s n^{-1/(q+4)}$ for some positive constant c_s ($s = 1, \dots, q$). The cross-validation methods discussed in Section 1.3 can be easily generalized to the multivariate data setting, and we can show that least squares cross-validation can optimally select the h_s 's in the sense outlined in Section 1.3 (see Section 1.8 below).

We briefly remark on the independence assumption invoked for the proofs presented above. Our assumption was that the data is independent across the i index. Note that no restrictions were placed on the s index for each component X_{is} ($s = 1, \dots, q$). The product kernel is used simply for convenience, and it certainly *does not* require that the X_{is} 's

are independent across the s index. In other words, the multivariate kernel density estimator (1.35) is capable of capturing general dependence among the different components of X_i . Furthermore, we shall relax the “independence across observations” assumption in Chapter 18, and will see that all of the results developed above carry over to the weakly dependent data setting.

1.7 Multivariate Bandwidth Selection: Rule-of-Thumb and Plug-In Methods

In Section 1.2 we discussed the use of the so-called normal reference rule-of-thumb and plug-in methods in a univariate setting. The generalization of the univariate normal reference rule-of-thumb to a multivariate setting is straightforward. Letting q be the dimension of X_i , one can choose $h_s = c_s X_{s, sd} n^{-1/(4+q)}$ for $s = 1, \dots, q$, where $X_{s, sd}$ is the sample standard deviation of $\{X_{is}\}_{i=1}^n$ and c_s is a positive constant. In practice one still faces the problem of how to choose c_s . The choice of $c_s = 1.06$ for all $s = 1, \dots, q$ is computationally attractive; however, this selection treats the different X_{is} ’s symmetrically. In practice, should the joint PDF change rapidly in one dimension (say in x_1) but change slowly in another (say in x_2), then one should select a relatively small value of c_1 (hence a small h_1) and a relatively large value for c_2 (h_2). Unlike the cross-validation methods that we will discuss shortly, rule-of-thumb methods do not offer this flexibility.

For plug-in methods, on the other hand, the leading (squared) bias and variance terms of $\hat{f}(x)$ must be estimated, and then h_1, \dots, h_q must be chosen to minimize the leading MSE term of $\hat{f}(x)$. However, the leading MSE term of $\hat{f}(x)$ involves the unknown $f(x)$ and its partial derivative functions, and pilot bandwidths must be selected for *each* variable in order to estimate these unknown functions. How to best select the initial pilot smoothing parameters can be tricky in high-dimensional settings, and the plug-in methods are not widely used in applied settings to the best of our knowledge, nor would we counsel their use other than for exploratory data analysis.

1.8 Multivariate Bandwidth Selection: Cross-Validation Methods

1.8.1 Least Squares Cross-Validation

The univariate least squares cross-validation method discussed in Section 1.3.1 can be readily generalized to the multivariate density estimation setting. Replacing the univariate kernel function in (1.23) by a multivariate product kernel, the cross-validation objective function now becomes

$$CV_f(h_1, \dots, h_q) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \bar{K}_h(X_i, X_j) - \frac{2}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i, j=1}^n K_h(X_i, X_j), \quad (1.38)$$

where

$$K_h(X_i, X_j) = \prod_{s=1}^q h_s^{-1} k\left(\frac{X_{is} - X_{js}}{h_s}\right),$$

$$\bar{K}_h(X_i, X_j) = \prod_{s=1}^q h_s^{-1} \bar{k}\left(\frac{X_{is} - X_{js}}{h_s}\right),$$

and $\bar{k}(v)$ is the twofold convolution kernel based upon $k(\cdot)$, where $k(\cdot)$ is a univariate kernel function satisfying (1.10).

Exercise 1.12 shows that the leading term of $CV_f(h_1, \dots, h_q)$ is given by (ignoring a term unrelated to the h_s 's)

$$CV_{f0}(h_1, \dots, h_q) = \int \left[\sum_{s=1}^q B_s(x) h_s^2 \right]^2 dx + \frac{\kappa^q}{n h_1 \dots h_q}, \quad (1.39)$$

where $B_s(x) = (\kappa_2/2) f_{ss}(x)$.

Defining a_s via $h_s = a_s n^{-1/(q+4)}$ ($s = 1, \dots, q$), we have

$$CV_{f0}(h_1, \dots, h_q) = n^{-4/(q+4)} \chi_f(a_1, \dots, a_q), \quad (1.40)$$

where

$$\chi_f(a_1, \dots, a_q) = \int \left[\sum_{s=1}^q B_s(x) a_s^2 \right]^2 dx + \frac{\kappa^q}{a_1 \dots a_q}. \quad (1.41)$$

Let the a_s^0 's be the values of the a_s 's that minimize $\chi_f(a_1, \dots, a_q)$. Under the same conditions used in the univariate case and, in addition, assuming that $f_{ss}(x)$ is not a zero function for all s , Li and Zhou (2005) show that each a_s^0 is uniquely defined, positive, and finite (see Exercise 1.10). Let h_1^0, \dots, h_q^0 denote the values of h_1, \dots, h_q that minimize CV_{f_0} . Then from (1.40) we know that $h_s^0 = a_s^0 n^{-1/(q+4)} = O(n^{-1/(q+4)})$.

Exercise 1.12 shows that CV_{f_0} is also the leading term of $E[CV_f]$. Therefore, the nonstochastic smoothing parameters h_s^0 can be interpreted as optimal smoothing parameters that minimize the leading term of the IMSE.

Let $\hat{h}_1, \dots, \hat{h}_q$ denote the values of h_1, \dots, h_q that minimize CV_f . Using the fact that $CV_f = CV_{f_0} + (s.o.)$, we can show that $\hat{h}_s = h_s^0 + o_p(h_s^0)$. Thus, we have

$$\frac{\hat{h}_s - h_s^0}{h_s^0} = \frac{\hat{h}_s}{h_s^0} - 1 \rightarrow 0 \quad \text{in probability, for } s = 1, \dots, q. \quad (1.42)$$

Therefore, smoothing parameters selected via cross-validation have the same asymptotic optimality properties as the nonstochastic optimal smoothing parameters.

Note that if $f_{ss}(x) = 0$ almost everywhere (a.e.) for some s , then $B_s = 0$ and the above result does not hold. Stone (1984) shows that the cross-validation method still selects h_1, \dots, h_q optimally in the sense that the integrated estimation square error is minimized; see also Ouyang et al. (2006) for a more detailed discussion of this case.

1.8.2 Likelihood Cross-Validation

Likelihood cross-validation for multivariate models follows directly via (multivariate) maximization of the likelihood function outlined in Section 1.3.2, hence we do not go into further details here. However, we do point out that, though straightforward to implement, it suffers from the same defects outlined for the univariate case in the presence of fat tail distributions (i.e., it has a tendency to oversmooth in such situations).

1.9 Asymptotic Normality of Density Estimators

In this section we show that $\hat{f}(x)$ has an asymptotic normal distribution. The most popular CLT is the Lindeberg-Levy CLT given in

Theorem A.3 of Appendix A, which states that $n^{1/2}[n^{-1} \sum_{i=1}^n Z_i] \rightarrow N(0, \sigma^2)$ in distribution, provided that Z_i is i.i.d. $(0, \sigma^2)$. Though the Lindeberg-Levy CLT can be used to derive the asymptotic distribution of various semiparametric estimators discussed in Chapters 7, 8, and 9, it cannot be used to derive the asymptotic distribution of $\hat{f}(x)$. This is because $\hat{f}(x) = n^{-1} \sum_i Z_{i,n}$, where the summand $Z_{i,n} = K_h(X_i, x)$ depends on n (since $h = h(n)$). We shall make use of the Liapunov CLT given in Theorem A.5 of Appendix A

Theorem 1.3. *Let X_1, \dots, X_n be i.i.d. q -vectors with its PDF $f(\cdot)$ having three-times bounded continuous derivatives. Let x be an interior point of the support of X . If, as $n \rightarrow \infty$, $h_s \rightarrow 0$ for all $s = 1, \dots, q$, $nh_1 \dots h_q \rightarrow \infty$, and $(nh_1 \dots h_q) \sum_{s=1}^q h_s^6 \rightarrow 0$, then*

$$\sqrt{nh_1 \dots h_q} \left[\hat{f}(x) - f(x) - \frac{\kappa_2}{2} \sum_{s=1}^q h_s^2 f_{ss}(x) \right] \xrightarrow{d} N(0, \kappa^q f(x)).$$

Proof. Using (1.36) and (1.37), one can easily show that

$$\sqrt{nh_1 \dots h_q} \left[\hat{f}(x) - f(x) - \frac{\kappa_2}{2} \sum_{s=1}^q h_s^2 f_{ss}(x) \right]$$

has asymptotic mean zero and asymptotic variance $\kappa^q f(x)$, i.e.,

$$\begin{aligned} & \sqrt{nh_1 \dots h_q} \left[\hat{f}(x) - f(x) - \frac{\kappa_2}{2} \sum_{s=1}^q h_s^2 f_{ss}(x) \right] \\ &= \sqrt{nh_1 \dots h_q} \left[\hat{f}(x) - E(\hat{f}(x)) \right] \\ & \quad + \sqrt{nh_1 \dots h_q} \left[E(\hat{f}(x)) - f(x) - \frac{\kappa_2}{2} \sum_{s=1}^q h_s^2 f_{ss}(x) \right] \\ &= \sqrt{nh_1 \dots h_q} \left[\hat{f}(x) - E(\hat{f}(x)) \right] \\ & \quad + O \left(\sqrt{nh_1 \dots h_q} \sum_{s=1}^q h_s^3 \right) \quad (\text{by (1.36)}) \\ &= \sum_{i=1}^n (nh_1 \dots h_q)^{-1/2} \\ & \quad \times \left[K \left(\frac{X_i - x}{h} \right) - E \left(K \left(\frac{X_i - x}{h} \right) \right) \right] + o(1) \\ &\equiv \sum_{i=1}^n Z_{n,i} + o(1) \xrightarrow{d} N(0, \kappa^q f(x)), \end{aligned}$$

by Liapunov's CLT, provided we can verify that Liapunov's CLT condition (A.21) holds, where

$$Z_{n,i} = (nh_1 \dots h_q)^{-1/2} \left[K \left(\frac{X_i - x}{h} \right) - E \left(K \left(\frac{X_i - x}{h} \right) \right) \right]$$

and

$$\sum_{i=1}^n \sigma_{n,i}^2 \stackrel{\text{def}}{=} \sum_{i=1}^n \text{var}(Z_{n,i}) = \kappa^q f(x) + o(1)$$

by (1.37). Pagan and Ullah (1999, p. 40) show that (A.21) holds under the condition given in Theorem 1.3. The condition that $\int k(v)^{2+\delta} dv < \infty$ for some $\delta > 0$ used in Pagan and Ullah is implied by our assumption that $k(v)$ is nonnegative and bounded, and that $\int k(v) dv = 1$, because $\int k(v)^{2+\delta} dv \leq C \int k(v) dv = C$ is finite, where $C = \sup_{v \in \mathbb{R}^q} k(v)^{1+\delta}$. \square

1.10 Uniform Rates of Convergence

Up to now we have demonstrated only the case of pointwise and IMSE consistency (which implies consistency in probability). In this section we generalize pointwise consistency in order to obtain a stronger “uniform consistency” result. We will prove that nonparametric kernel estimators are uniformly almost surely consistent and derive their uniform almost sure rate of convergence. Almost sure convergence implies convergence in probability; however, the converse is not true, i.e., convergence in probability may not imply convergence almost surely; see Serfling (1980) for specific examples.

We have already established pointwise consistency for an interior point in the support of X . However, it turns out that popular kernel functions such as (1.9) may not lead to consistent estimation of $f(x)$ when x is at the boundary of its support, hence we need to exclude the boundary ranges when considering the uniform convergence rate. This highlights an important aspect of kernel estimation in general, and a number of kernel estimators introduced in later sections are motivated by the desire to mitigate such “boundary effects.” We first show that when x is at (or near) the boundary of its support, $\hat{f}(x)$ may not be a consistent estimator of $f(x)$.

Consider the case where X is univariate having bounded support. For simplicity we assume that $X \in [0, 1]$. The pointwise consistency result $\hat{f}(x) - f(x) = o_p(1)$ obtained earlier requires that x lie in the

interior of its support. Exercise 1.13 shows that, for x at the boundary of its support, $\text{MSE } \hat{f}(x)$ may not be $o(1)$. Therefore, some modifications may be needed to consistently estimate $f(x)$ for x at the boundary of its support. Typical modifications include the use of boundary kernels or data reflection (see Gasser and Müller (1979), Hall and Wehrly (1991), and Scott (1992, pp. 148–149)). By way of example, consider the case where x lies on its lowermost boundary, i.e., $x = 0$, hence $\hat{f}(0) = (nh)^{-1} \sum_{i=1}^n K((X_i - 0)/h)$. Exercise 1.13 shows that for this case, $E[\hat{f}(0)] = f(0)/2 + O(h)$. Therefore, $\text{bias } \hat{f}(0) = E[\hat{f}(0)] - f(0) = -f(0)/2 + O(h)$, which will not converge to zero if $f(0) \neq 0$ (when $f(0) > 0$).

In the literature, various boundary kernels are proposed to overcome the boundary (bias) problem. For example, a simple boundary corrected kernel is given by (assuming that $X \in [0, 1]$)

$$k_h(x, y) = \begin{cases} h^{-1}k\left(\frac{y-x}{h}\right) / \int_{-x/h}^{\infty} k(v) dv & \text{if } x \in [0, h] \\ h^{-1}k\left(\frac{y-x}{h}\right) & \text{if } x \in [h, 1-h] \\ h^{-1}k\left(\frac{y-x}{h}\right) / \int_{-\infty}^{(1-x)/h} k(v) dv & \text{if } x \in (1-h, 1], \end{cases} \quad (1.43)$$

where $k(\cdot)$ is a second order kernel satisfying (1.10). Now, we estimate $f(x)$ by

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n k_h(x, X_i), \quad (1.44)$$

where $k_h(x, X_i)$ is defined in (1.43). Exercise 1.14 shows that the above boundary corrected kernel successfully overcomes the boundary problem.

We now establish the uniform almost sure convergence rate of $\hat{f}(x) - f(x)$ for $x \in \mathcal{S}$, where \mathcal{S} is a bounded set excluding the boundary range of the support of X . In the above example, when the support of x is $[0, 1]$, we can choose $\mathcal{S} = [\epsilon, 1 - \epsilon]$ for arbitrarily *small* positive ϵ ($0 < \epsilon < 1/2$). We assume that $f(x)$ is bounded below by a positive constant on \mathcal{S} .

Theorem 1.4. *Under smoothness conditions on $f(\cdot)$ given in Masry (1996b), and also assuming that $\inf_{x \in \mathcal{S}} f(x) \geq \delta > 0$, we have*

$$\sup_{x \in \mathcal{S}} |\hat{f}(x) - f(x)| = O \left(\frac{(\ln(n))^{1/2}}{(nh_1 \dots h_q)^{1/2}} + \sum_{s=1}^q h_s^2 \right) \text{ almost surely.}$$

A detailed proof of Theorem 1.4 is given in Section 1.12.

Since almost sure convergence implies convergence in probability, the uniform rate also holds in probability, i.e., under the same conditions as in Theorem 1.4, we have

$$\sup_{x \in \mathcal{S}} |\hat{f}(x) - f(x)| = O_p \left(\frac{(\ln(n))^{1/2}}{(nh_1 \dots h_q)^{1/2}} + \sum_{s=1}^q h_s^2 \right).$$

Using the results of (1.36) and (1.37), we can establish the following uniform MSE rate.

Theorem 1.5. *Assuming that $f(x)$ is twice differentiable with bounded second derivatives, then we have*

$$\sup_{x \in \mathcal{S}} \mathbb{E} \left\{ \left[\hat{f}(x) - f(x) \right]^2 \right\} = O \left(\sum_{s=1}^q h_s^4 + (nh_1 \dots h_q)^{-1} \right).$$

Proof. This follows from (1.36) and (1.37), by noting that $\sup_{x \in \mathcal{S}} f(x)$ and $\sup_{x \in \mathcal{S}} |f_{ss}(x)|$ are both finite ($s = 1, \dots, q$). \square

Note that although convergence in MSE implies convergence in probability, one cannot derive the uniform convergence rate in probability from Theorem 1.5. This is because

$$\mathbb{E} \left\{ \sup_{x \in \mathcal{S}} \left[\hat{f}(x) - f(x) \right]^2 \right\} \neq \sup_{x \in \mathcal{S}} \mathbb{E} \left[\hat{f}(x) - f(x) \right]^2,$$

and

$$\mathbb{P} \left[\sup_{x \in \mathcal{S}} |\hat{f}(x) - f(x)| > \epsilon \right] \neq \sup_{x \in \mathcal{S}} \mathbb{P} \left[|\hat{f}(x) - f(x)| > \epsilon \right].$$

The sup and the $\mathbb{E}(\cdot)$ or the $\mathbb{P}(\cdot)$ operators do not commute with one another.

Cheng (1997) proposes alternative (local linear) density estimators that achieve automatic boundary corrections and enjoy some typical optimality properties. Cheng also suggests a data-based bandwidth selector (in the spirit of plug-in rules), and demonstrates that the bandwidth selector is very efficient regardless of whether there are non-smooth boundaries in the support of the density.

1.11 Higher Order Kernel Functions

Recall that decreasing the bandwidth h lowers the bias of a kernel estimator but increases its variance. Higher order kernel functions are devices used for bias reduction which are also capable of reducing the MSE of the resulting estimator. Many popular kernel functions such as the one defined in (1.10) are called “second order” kernels. The order of a kernel, ν ($\nu > 0$), is defined as the order of the first nonzero moment. For example, if $\int uk(u) du = 0$, but $\int u^2k(u) du \neq 0$, then $k(\cdot)$ is said to be a second order kernel ($\nu = 2$). A general ν th order kernel ($\nu \geq 2$ is an integer) must therefore satisfy the following conditions:

$$\begin{aligned} (i) \quad & \int k(u) du = 1, \\ (ii) \quad & \int u^l k(u) du = 0, \quad (l = 1, \dots, \nu - 1), \\ (iii) \quad & \int u^\nu k(u) du = \kappa_\nu \neq 0. \end{aligned} \tag{1.45}$$

Obviously, when $\nu = 2$, (1.45) collapses to (1.10).

If one replaces the second order kernel in $\hat{f}(x)$ of (1.35) by a ν th order kernel function, then as was the case when using a second order kernel, under the assumption that $f(x)$ is ν th order differentiable, and assuming that the h_s 's all have the same order of magnitude, one can show that

$$\text{bias}(\hat{f}(x)) = O\left(\sum_{s=1}^q h_s^\nu\right) \tag{1.46}$$

and

$$\text{var}(\hat{f}(x)) = O((nh_1 \dots h_q)^{-1}) \tag{1.47}$$

(see Exercise 1.15). Hence, we have

$$\text{MSE}(\hat{f}(x)) = O\left(\sum_{s=1}^q h_s^{2\nu} + (nh_1 \dots h_q)^{-1}\right) \tag{1.48}$$

and

$$\hat{f}(x) - f(x) = O_p\left(\sum_{s=1}^q h_s^\nu + (nh_1 \dots h_q)^{-1/2}\right).$$

Thus, by using a ν th higher order kernel function ($\nu > 2$), one can reduce the order of the bias of $\hat{f}(x)$ from $O(\sum_{s=1}^q h_s^2)$ to $O(\sum_{s=1}^q h_s^\nu)$,

and the optimal value of h_s may once again be obtained by balancing the squared bias and the variance, giving $h_s = O(n^{-1/(2\nu+q)})$, while the rate of convergence is now $\hat{f}(x) - f(x) = O_p(n^{-\nu/(2\nu+q)})$. Assuming that $f(x)$ is differentiable up to any finite order, then one can choose ν to be sufficiently large, and the resulting rate can be made arbitrarily close to $O_p(n^{-1/2})$. Note, however, that for $\nu > 2$, no nonnegative kernel exists that satisfies (1.45). This means that, necessarily, we have to assign negative weights to some range of the data which implies that one may get *negative* density estimates, clearly an undesirable side-effect. Furthermore, in finite-sample applications nonnegative second order kernels have often been found to yield more stable estimation results than their higher order counterparts. Therefore, higher order kernel functions are mainly used for theoretical purposes; for example, to achieve a \sqrt{n} -rate of convergence for some finite dimensional parameter in a semiparametric model, one often has to use high order kernel functions (see Chapter 7 for such an example).

Higher order kernel functions are quite easy to construct. Assuming that $k(u)$ is symmetric around zero,⁵ i.e., $k(u) = k(-u)$, then $\int u^{2m+1}k(u) du = 0$ for all positive integers m . By way of example, in order to construct a simple fourth order kernel (i.e., $\nu = 4$), one could begin with, say, a second order kernel such as the standard normal kernel, set up a polynomial in its argument, and solve for the roots of the polynomial subject to the desired moment constraints. For example, letting $\Phi(u) = (2\pi)^{-1/2} \exp(-u^2/2)$ be a second order Gaussian kernel, we could begin with the polynomial

$$k(u) = (a + bu^2)\Phi(u), \quad (1.49)$$

where a and b are two constants which must satisfy the requirements of a fourth order kernel. Letting $k(u)$ satisfy (1.45) with $\nu = 4$ ($\int u^l k(u) du = 0$ for $l = 1, 3$ because $k(u)$ is an even function), we therefore only require $\int k(u) du = 1$ and $\int u^2 k(u) du = 0$. From these two restrictions, one can easily obtain the result $a = 3/2$ and $b = -1/2$. For readers requiring some higher order kernel functions, we provide a few examples based on the second order Gaussian and Epanechnikov kernels, perhaps the two most popular kernels in applied nonparametric estimation. As noted, the fourth order univariate Gaussian kernel

⁵Typically, only symmetric kernel functions are used in practice, though see Abadir and Lawford (2004) for recent work involving optimal asymmetric kernels.

is given by the formula

$$k(u) = \left(\frac{3}{2} - \frac{1}{2}u^2 \right) \frac{\exp(-u^2/2)}{\sqrt{2\pi}},$$

while the sixth order univariate Gaussian kernel is given by

$$k(u) = \left(\frac{15}{8} - \frac{5}{4}u^2 + \frac{1}{8}u^4 \right) \frac{\exp(-u^2/2)}{\sqrt{2\pi}}.$$

The second order univariate Epanechnikov kernel is the *optimal* kernel based on a calculus of variations solution to minimizing the IMSE of the kernel estimator (see Serfling (1980, pp. 40–43)). The univariate second order Epanechnikov kernel is given by the formula

$$k(u) = \begin{cases} \frac{3}{4\sqrt{5}} \left(1 - \frac{1}{5}u^2 \right) & \text{if } u^2 < 5.0 \\ 0 & \text{otherwise,} \end{cases}$$

the fourth order univariate Epanechnikov kernel by

$$k(u) = \begin{cases} \frac{3}{4\sqrt{5}} \left(\frac{15}{8} - \frac{7}{8}u^2 \right) \left(1 - \frac{1}{5}u^2 \right) & \text{if } u^2 < 5.0 \\ 0 & \text{otherwise,} \end{cases}$$

while the sixth order univariate Epanechnikov kernel is given by

$$k(u) = \begin{cases} \frac{3}{4\sqrt{5}} \left(\frac{175}{64} - \frac{105}{32}u^2 + \frac{231}{320}u^4 \right) \left(1 - \frac{1}{5}u^2 \right) & \text{if } u^2 < 5.0 \\ 0 & \text{otherwise.} \end{cases}$$

Figure 1.2 plots the second, fourth, and sixth order Epanechnikov kernels defined above. Clearly, for $\nu > 2$, the kernels indeed assign negative weights which can result in negative density estimates, not a desirable feature.

For related work involving exact mean integrated squared error for higher order kernels in the context of univariate kernel density estimation, see Hansen (2005). Also, for related work using iterative methods to estimate transformation-kernel densities, see Yang and Marron (1999) and Yang (2000).

1.12 Proof of Theorem 1.4 (Uniform Almost Sure Convergence)

The proof below is based on the arguments presented in Masry (1996b), who establishes uniform almost sure rates for local polynomial regression with weakly dependent (α -mixing) data; see Chapter 18 for further details on weakly dependent processes. Since the bias of the kernel

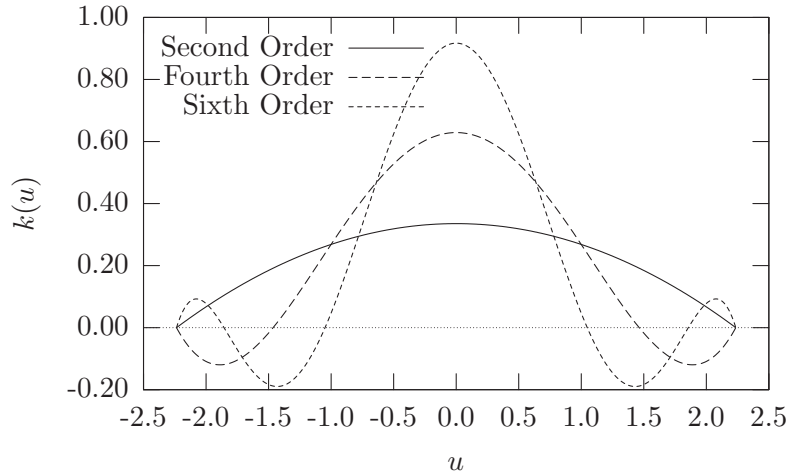


Figure 1.2: Epanechnikov kernels of varying order.

density estimator is of order $O(\sum_{s=1}^q h_s^2)$ and the variance is of order $O((nh_1 \dots h_q)^{-1})$, it is easy to show that the optimal rate of convergence requires that all h_s should be of the same order of magnitude. Therefore, for expositional simplicity, we will make the simplifying assumption that

$$h_1 = \dots = h_q = h.$$

This will not affect the optimal rate of convergence, but it simplifies the derivation tremendously. We emphasize that, in practice, one should always allow h_s ($s = 1, \dots, s$) to differ from each other, which is of course always permitted when using fully data-driven methods of bandwidth selection such as cross-validation. Only for the theoretical analysis that immediately follows do we assume that all smoothing parameters are the same.

Proof. Let $W_n = W_n(x) = |\hat{f}(x) - f(x)|$. To prove that the random variable W_n is of order (η) almost surely (a.s.), we can show that $\sum_{n=1}^{\infty} P(|W_n/\eta| > 1)$ is finite (for some $\eta > 0$). Then by the Borel-Cantelli lemma (see Lemma A.7 in Appendix A), we know that $W_n = O(\eta)$ a.s. Here, the supremum operator complicates the proof because \mathcal{S} is an uncountable set. Letting L_n denote a countable set,

then we have

$$P\left(\max_{x \in L_n} W_n(x) > \eta\right) \leq (\# \text{ of } L_n) \max_{x \in L_n} P(W_n(x) > \eta). \quad (1.50)$$

But in our case, $x \in \mathcal{S}$ is uncountable and we cannot simply use an inequality like (1.50) to bound $P(\sup_{x \in \mathcal{S}} W_n(x) > \eta)$.

However, since \mathcal{S} is a bounded set we can partition \mathcal{S} into countable subsets with the volume of each subset being as small as necessary. Then $P(\sup_{x \in \mathcal{S}} |W_n(x)| > \eta)$ can be transformed into a problem like $P(\max_{x \in L_n} |W_n(x)| > \eta)$ and the inequality of (1.50) can be used to handle this term. Using this idea we prove Theorem 1.4 below.

We write

$$\begin{aligned} |\hat{f}(x) - f(x)| &= \left| \hat{f}(x) - E(\hat{f}(x)) + E(\hat{f}(x)) - f(x) \right| \\ &\leq \left| \hat{f}(x) - E(\hat{f}(x)) \right| + \left| E(\hat{f}(x)) - f(x) \right|. \end{aligned}$$

We prove Theorem 1.4 by showing that

$$\sup_{x \in \mathcal{S}} |E(\hat{f}(x)) - f(x)| = O(h^2), \quad (1.51)$$

and that

$$\sup_{x \in \mathcal{S}} |\hat{f}(x) - E(\hat{f}(x))| = O\left(\frac{(\ln(n))^{1/2}}{(nh^q)^{1/2}}\right) \text{ almost surely.} \quad (1.52)$$

We first prove (1.51). Because the compact set \mathcal{S} is in the interior of its support, by a change-of-variables argument, we have, for $x \in \mathcal{S}$,

$$\begin{aligned} E(\hat{f}(x)) - f(x) &= \int f(x + hv) K(v) dv - f(x) \\ &= h^2 \int v' f^{(2)}(\tilde{x}) v K(v) dv \\ &\leq C_1 h^2 \int v' v K(v) dv \leq Ch^2 = O(h^2) \end{aligned}$$

uniformly in $x \in \mathcal{S}$. Thus, we have proved (1.51).

We now turn to the proof of (1.52). Since \mathcal{S} is compact (closed and bounded), it can be covered by a finite number $L_n = L(n)$ of (q -dimensional) cubes $I_k = I_{k,n}$, with centers $x_{k,n}$ and length l_n ($k =$

$1, \dots, L(n)$). We know that $L_n = \text{constant}/(l_n)^q$ because \mathcal{S} is compact, which gives $l_n = \text{constant}/L_n^{1/q}$. We write

$$\begin{aligned} \sup_{x \in \mathcal{S}} \left| \hat{f}(x) - \mathbb{E} \left(\hat{f}(x) \right) \right| &= \max_{1 \leq k \leq L(n)} \sup_{x \in \mathcal{S} \cap I_k} \left| \hat{f}(x) - \mathbb{E} \left(\hat{f}(x) \right) \right| \\ &\leq \max_{1 \leq k \leq L(n)} \sup_{x \in \mathcal{S} \cap I_k} \left| \hat{f}(x) - \hat{f}(x_{k,n}) \right| \\ &\quad + \max_{1 \leq k \leq L(n)} \left| \hat{f}(x_{k,n}) - \mathbb{E} \left(\hat{f}(x_{k,n}) \right) \right| \\ &\quad + \max_{1 \leq k \leq L(n)} \sup_{x \in \mathcal{S} \cap I_k} \left| \mathbb{E} \left(\hat{f}(x_{k,n}) \right) - \mathbb{E} \left(\hat{f}(x) \right) \right| \\ &\equiv Q_1 + Q_2 + Q_3. \end{aligned}$$

Note that Q_2 does not depend on x , so $\sup_{x \in \mathcal{S} \cap I_k}$ does not appear in the definition of Q_2 .

We first consider Q_2 . Write $W_n(x) = \hat{f}(x) - \mathbb{E} \left(\hat{f}(x) \right) = \sum_i Z_{n,i}$, where $Z_{n,i} = (nh^q)^{-1} \{K((X_i - x)/h) - \mathbb{E}[K((X_i - x)/h)]\}$. For any $\eta > 0$, we have

$$\begin{aligned} \mathbb{P}[Q_2 > \eta] &= \mathbb{P} \left[\max_{1 \leq k \leq L(n)} |W_n(x_{k,n})| > \eta \right] \\ &\leq \mathbb{P}[W_n(x_{1,n}) > \eta \text{ or } W_n(x_{2,n}) > \eta, \dots, \text{ or } W_n(x_{L(n),n}) > \eta] \\ &\leq \mathbb{P}(W_n(x_{1,n}) > \eta) + \mathbb{P}(W_n(x_{2,n}) > \eta) + \dots \\ &\quad + \mathbb{P}(W_n(x_{L(n),n}) > \eta) \\ &\leq L(n) \sup_{x \in \mathcal{S}} \mathbb{P}[|W_n(x)| > \eta]. \end{aligned} \tag{1.53}$$

Since $K(\cdot)$ is bounded, and letting $A_1 = \sup_x |K(x)|$, we have $|Z_{n,i}| \leq 2A_1/(nh^q)$ for all $i = 1, \dots, n$. Define $\lambda_n = (nh^q \ln(n))^{1/2}$. Then $\lambda_n |Z_{n,i}| \leq 2A_1 [\ln(n)/(nh^q)]^{1/2} \leq 1/2$ for all $i = 1, \dots, n$ for n sufficiently large.⁶ Using the inequality $\exp(x) \leq 1 + x + x^2$ for $|x| \leq 1/2$, we have $\exp(\pm \lambda_n Z_{n,i}) \leq 1 + \lambda_n Z_{n,i} + \lambda_n^2 Z_{n,i}^2$. Hence,

$$\mathbb{E}[\exp(\pm \lambda_n Z_{n,i})] \leq 1 + \lambda_n^2 \mathbb{E}[Z_{n,i}^2] \leq \exp[\mathbb{E}(\lambda_n^2 Z_{n,i}^2)], \tag{1.54}$$

where we used $\mathbb{E}(Z_{n,i}) = 0$ while for the second inequality we used $1 + v \leq \exp(v)$ for $v \geq 0$ ($v = \mathbb{E}[\lambda_n^2 Z_{n,i}^2]$).

⁶For now, any choice of $\lambda_n \leq (nh^q)/(4A_1)$ will lead to $|\lambda_n Z_{n,i}| \leq 1/2$. Later on we will show that, in order to obtain the optimal rate for Q_2 , one needs to choose $\lambda_n = (nh^q \ln(n))^{1/2}$.

By the Markov inequality (see Lemma A.23 with $\phi(x) = \exp(ax)$) we know that

$$P[X > c] \leq \frac{E[\exp(Xa)]}{\exp(ac)}, \quad (a > 0). \quad (1.55)$$

Using (1.55) we have

$$\begin{aligned} P[|W_n(x)| > \eta] &= P\left[\left|\sum_{i=1}^n Z_{n,i}\right| > \eta\right] \\ &= P\left[\sum_{i=1}^n Z_{n,i} > \eta\right] + P\left[\sum_{i=1}^n Z_{n,i} < -\eta\right] \\ &\leq P\left[\sum_{i=1}^n Z_{n,i} > \eta\right] + P\left[-\sum_{i=1}^n Z_{n,i} > \eta\right] \\ &\leq \frac{E[\exp(\lambda_n \sum_{i=1}^n Z_{n,i})] + E[\exp(-\lambda_n \sum_{i=1}^n Z_{n,i})]}{\exp(\lambda_n \eta_n)} \\ &\quad (\text{by (1.55), } a = \lambda_n, c = \eta) \\ &\leq 2 \exp(-\lambda_n \eta) \left[\exp\left(\lambda_n^2 \sum_{i=1}^n E(Z_{n,i}^2)\right) \right] \\ &\quad (\text{by (1.54)}) \\ &\leq 2 \exp(-\lambda_n \eta) \left[\exp(A_2 \lambda_n^2 / (nh^q)) \right], \end{aligned} \quad (1.56)$$

where we used

$$E[Z_{n,i}^2] \leq (nh^q)^{-2} E[K^2((X_i - x)/h)] \leq A_2(n^2 h^q)^{-1} [1 + o(1)].$$

Because the last bound in (1.56) is independent of x , it is also the uniform bound, i.e.,

$$\sup_{x \in \mathcal{S}} P[|W_n(x)| > \eta] \leq 2 \exp\left(-\lambda_n \eta + \frac{A_2 \lambda_n^2}{nh^q}\right). \quad (1.57)$$

We want to have $\eta \rightarrow 0$ as fast as possible, and at the same time we need $\lambda_n \eta \rightarrow \infty$ at a rate which ensures that (1.57) is summable.⁷ We can choose $\lambda_n \eta = C_4 \ln(n)$, or $\lambda_n = C_4 \ln(n)/\eta$. Finding the fastest rate for which $\eta \rightarrow 0$ is equivalent to finding the fastest rate for which $\lambda_n \rightarrow \infty$. We also need the order of $\lambda_n \eta \geq \lambda_n^2 / (nh^d)$, or $\ln(n) \geq \lambda_n^2 / (nh^d)$.

⁷A sequence $\{a_n\}_{n=1}^\infty$ is said to be summable if $|\sum_{j=1}^\infty a_j| < \infty$.

Thus, we simply need to maximize the order of $\lambda_n \rightarrow \infty$ subject to $\lambda_n^2 \leq (nh^q) \ln(n)$. Doing so, we get

$$\lambda_n = [(nh^q) \ln(n)]^{1/2} \quad \text{and} \quad \eta = C_4 \ln(n) / \lambda_n = C_4 [\ln(n) / (nh^q)]^{1/2}. \quad (1.58)$$

Using (1.58) we get

$$\begin{aligned} -\lambda_n \eta / 2 + A_2 \lambda_n^2 / (nh^q) &= -C_4 \ln(n) + A_2 \ln(n) \\ &= -\alpha \ln(n), \end{aligned}$$

where $\alpha = C_4 - A_2$. Substituting this into (1.57) and then into (1.53) gives us

$$P[Q_2 > \eta_n] \leq 2L(n)/n^\alpha. \quad (1.59)$$

By choosing C_4 sufficiently large, we can obtain the result that $L(n)/n^\alpha$ is summable by properly choosing the order of $L(n)$, i.e., $\sum_{n=1}^{\infty} P(|Q_2/\eta_n| > 1) \leq 4 \sum_{n=1}^{\infty} L(n)/n^\alpha < \infty$. Therefore, by the Borel-Cantelli lemma we know that

$$Q_2 = O(\eta_n) = O\left((\ln(n))^{1/2} / (nh^q)^{1/2}\right) \quad \text{almost surely.} \quad (1.60)$$

We now consider Q_1 and Q_3 . Recall that $\|\cdot\|$ denotes the usual Euclidean norm of a vector. By the Lipschitz condition on $K(\cdot)$, we know that

$$\begin{aligned} \sup_{x \in S \cap I_k} |K((X_i - x)/h) - K((X_i - x_{k,n})/h)| &\leq C_1 h^{-1} \sup_{x \in S \cap I_k} \|x - x_{k,n}\| \\ &\leq C_2 h^{-1} l_n. \end{aligned}$$

Therefore, by choosing $l_n = (\ln(n))^{1/2} h^{(q+2)/2} / n^{1/2}$, we have

$$|Q_1| \leq C_2 h^{-(q+1)} l_n = O\left((\ln(n) / (nh^q))^{1/2}\right). \quad (1.61)$$

By exactly the same argument we can show that

$$|Q_3| \leq C_3 h^{-(q+1)} l_n = O\left((\ln(n) / (nh^q))^{1/2}\right). \quad (1.62)$$

Equations (1.60) through (1.62) prove (1.52), and this completes the proof of Theorem 1.4. \square

1.13 Applications

We now consider a number of applications of univariate and multivariate density estimation that illustrate the flexibility and power of the kernel approach.

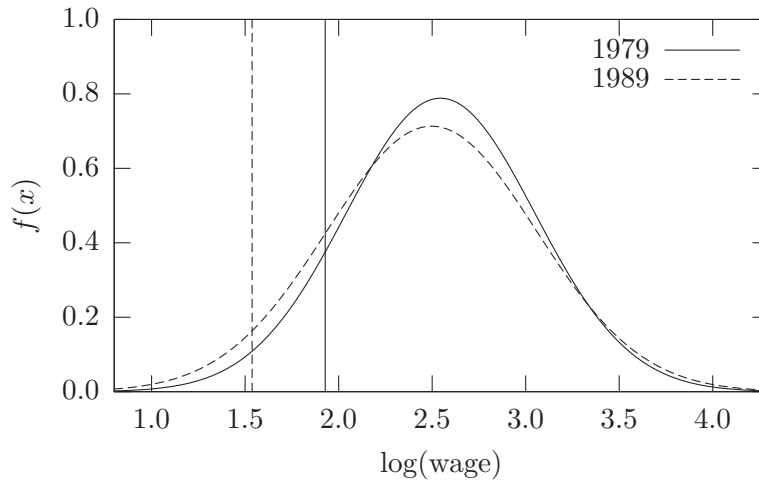


Figure 1.3: Parametric density estimate (vertical lines represent (log) minimum wages in 1979 and 1989).

1.13.1 Female Wage Inequality

DiNardo and Tobias (2001, p. 12) used nonparametric kernel methods to investigate the phenomenon of female wage inequality which grew from 1979 to 1989. Sometimes the scale of a parametric distribution is used as a crude measure of inequality, and the standard deviation of log wages increased 25% from 0.41 to 0.50 over this period.⁸ One might think that common culprits underlying such changes would include international trade, technical change, or perhaps organizational change. As we will see below, DiNardo and Tobias show that the kernel estimator can help reveal who the true culprit is.

If one used a parametric model and assumed, say, a normal distribution for log wages, one would arrive at the description of the data presented in Figure 1.3.

Use of nonparametric kernel methods and a simple “normal refer-

⁸The minimum wages in 1979 and 1989 were \$2.90/hour and \$3.35/hour, while the CPI was 72.6, 124.0, and 172.2 in 1979, 1989, and 2000 respectively. Wages were taken from the Current Population Survey (CPS). There were 140,284 and 167,863 observations in the 1979 and 1989 samples respectively. The Gaussian kernel was used, and the normal reference rule-of-thumb bandwidths were 0.050 and 0.053 for the 1979 and 1989 samples respectively. Wage values appearing in Figures 1.3 and 1.4 are in current (2000) dollars.

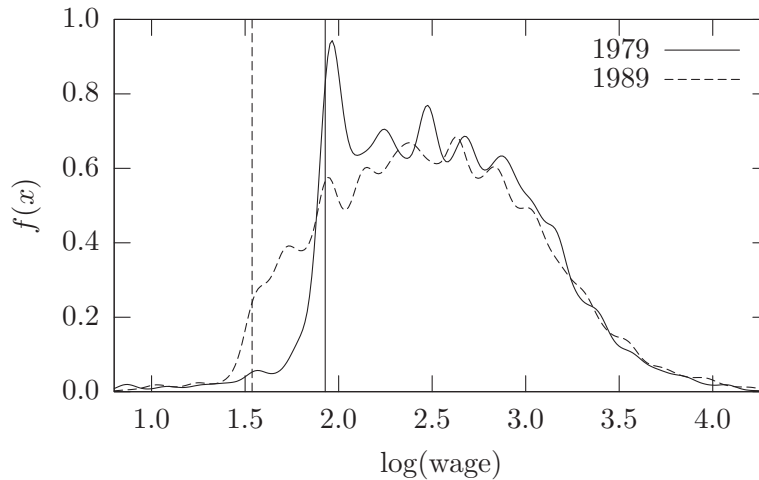


Figure 1.4: Nonparametric density estimate (vertical lines represent (log) minimum wages in 1979 and 1989).

ence rule-of-thumb” ($h = 1.06\sigma n^{-1/5}$) bandwidth along with a second order Gaussian kernel yields the estimates plotted in Figure 1.4.

The two kernel density estimates based on the normal reference rule-of-thumb presented in Figure 1.4 appear to be undersmoothed. However, these estimates clearly reveal a feature not captured by parametric methods: a binding modal minimum wage for 1979 that is no longer binding in 1989 for most women. This finding suggests that the growing wage inequality can be explained by truncation induced by a binding real minimum wage in 1979. That is, in 1979, unlike 1989, employers were paying minimum wage to many employees, which distorts and *reduces* the variance of the wage distribution. The real value of the minimum wage falls over time, becoming nonbinding in 1989. Thus, the nonparametric estimator readily reveals the true reason underlying growing wage inequality, and focuses attention away from other possible explanations, such as international trade, technical change, or possibly organizational change. This example serves simply to underscore the fact that traditional parametric approaches may mask important characteristics present in data.

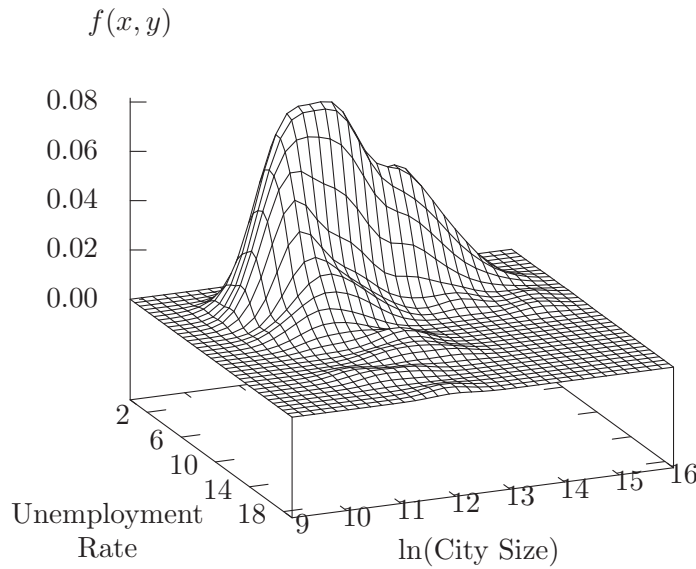


Figure 1.5: Unemployment rate and $\ln(\text{city size})$ joint density estimate.

1.13.2 Unemployment Rates and City Size

For this example we use U.S. data on city population ($\ln(\text{city size})$) and unemployment rates based upon a sample of $n = 295$ cities. Gan and Zhang (2006) present a theory predicting that the larger the city, the smaller the unemployment rate (on average). In Figure 1.5 we plot the estimated joint PDF using least squares cross-validated bandwidth selection and a second order Gaussian kernel. The cross-validated bandwidths are 0.665 and 0.351 for the unemployment rate and population respectively.

The joint density estimate presented in Figure 1.5 is consistent with the hypothesis that large cities tend to have low unemployment rates and vice versa. That is, Figure 1.5 reveals a somewhat “right-angled” distribution having probability mass at low unemployment rates and large city sizes, while as the city size falls we observe the probability mass shifting first toward the origin and then, as city size falls further, the mass shifts toward higher unemployment rates.

1.13.3 Adolescent Growth

Abnormal adolescent growth can provide an early warning that a child has a medical problem. For instance, too rapid growth may indicate the presence of a hydrocephalus (an accumulation of liquid within the cavity of the cranium), a brain tumor, or other conditions that cause macrocephaly (having an unusually large head), while too slow growth may indicate malformations of the brain, early fusion of sutures or other problems. Insufficient gain in weight, height or a combination may indicate failure-to-thrive, chronic illness, neglect or other problems.

We consider data from the population of healthy U.S. children obtained from the Center for Disease Control and Prevention's (CDC) National Health and Nutrition Examination Survey. We combine data and use two recent cross-sectional nationally representative health examination surveys for the years 1999/2000 and 2001/2002. For each cross section, two separate datasets must be linked (a body measurement dataset and a demographic variable dataset). The combined linked datasets contains 8,399 complete observations for children and youths ages 2-20 years of age. We model the joint distribution of height and weight by sex.

Figures 1.6 and 1.7 reveal that the joint distribution of height and weight is similar for males and females; however, that for males contains greater probability mass at higher values of both weight and height. That is, one is more likely to observe both taller and heavier boys than girls. Such data lays the foundation for the construction of adolescent growth charts, for instance, weight for stature charts.⁹ See also Wei and He (2006) for related work on conditional growth charts.

1.13.4 Old Faithful Geyser Data

The Old Faithful Geyser is a tourist attraction located in Yellowstone National Park. This famous dataset containing $n = 272$ observations consists of two variables, eruption duration (minutes) and waiting time until the next eruption (minutes). This dataset is used by the park service to model, among other things, expected duration conditional upon the amount of time that has elapsed since the previous eruption. Modeling the joint distribution is, however, of interest in its own right. The underlying bimodal nature of the joint PDF is readily revealed by

⁹See <http://www.cdc.gov/growthcharts> for official growth charts developed by the National Center for Health Statistics.

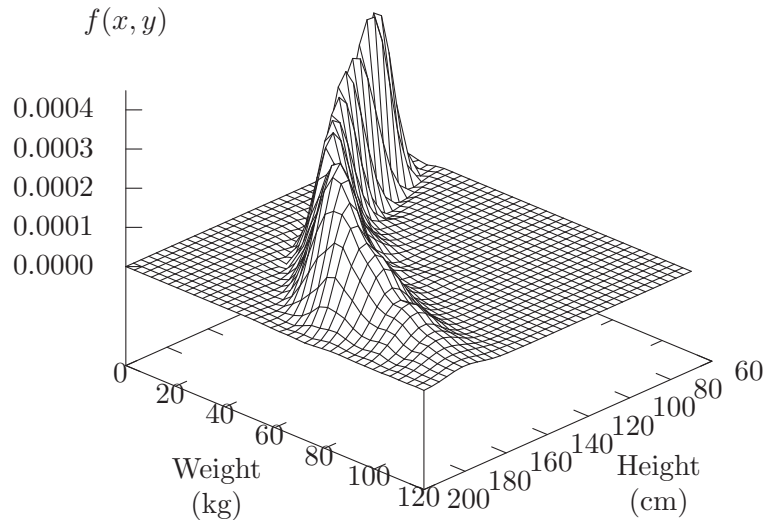


Figure 1.6: Weight and height joint density estimate for males.

the kernel estimator graphed in Figure 1.8 constructed using likelihood cross-validated bandwidths and a second order Gaussian kernel.¹⁰

If one were to instead model this density with a parametric model such as the bivariate normal (being symmetric, unimodal, and monotonically decreasing away from the mode), one would of course fail to uncover the underlying structure readily revealed by the kernel estimate.

1.13.5 Evolution of Real Income Distribution in Italy, 1951–1998

Baiocchi (2006) recently considered the evolution of the distribution of real income in Italy using kernel methods. He considers a series of

¹⁰Likelihood cross-validated bandwidths were computed and were equal to $(h_1, h_2) = (0.368\sigma_1 n^{-1/6}, 0.764\sigma_2 n^{-1/6})$, while least squares cross-validated bandwidths were $(h_1, h_2) = (0.307\sigma_1 n^{-1/6}, 0.733\sigma_2 n^{-1/6})$ where h_1 is that for eruption duration and h_2 that for waiting time.

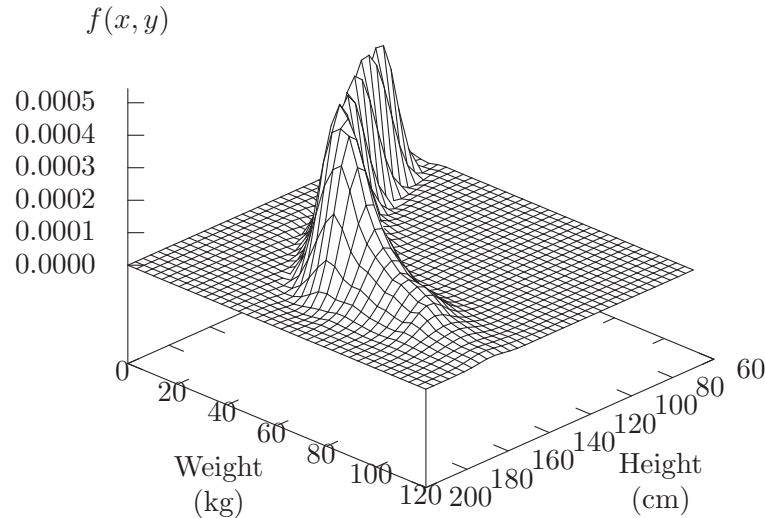


Figure 1.7: Weight and height joint density estimate for females.

“stacked” univariate kernel density estimates of the income distribution for 21 regions and plots the resulting evolution of the univariate kernel density estimates over time. We are indebted to Giovanni Baiocchi for providing the data containing observations for the period 1951–1998 (millions of lire, 1990 = base) used to generate a series of univariate kernel estimates using likelihood cross-validation. Figure 1.9 presents the evolution of real GDP per capita (millions of 1990 lire) by stacking the series of annual (i.e., cross section) univariate kernel estimates in a 3D plot.

Figure 1.9 reveals that the distribution of income has evolved from a unimodal one in the early 1950s to a markedly bimodal one in the 1990s. This result is robust to bandwidth choice, and is observed whether using simple rules-of-thumb or data-driven methods such as likelihood cross-validation. The kernel method readily reveals this evolution which might easily be missed were one to use parametric models of the income distribution (e.g., the lognormal distribution commonly found in

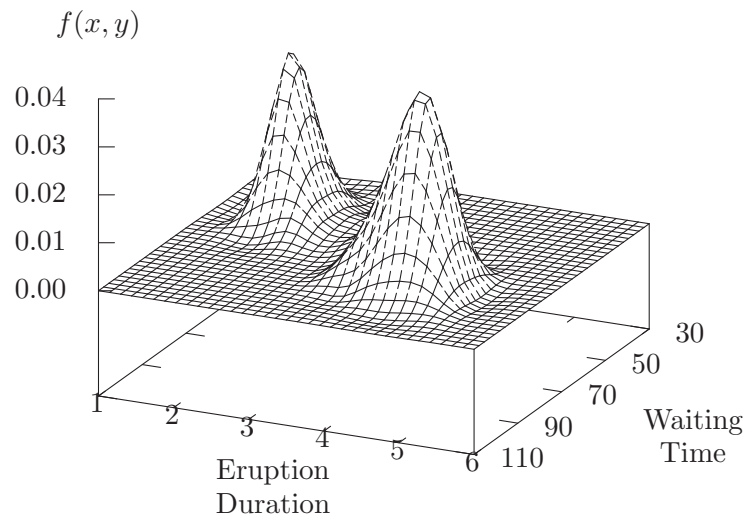


Figure 1.8: Joint density estimate for the Old Faithful data.

applied work).

1.14 Exercises

Exercise 1.1. Consider the following sample of continuous data:

$$\{-0.57, 0.25, -0.08, 1.40, -1.05, -1.00, 0.37, -1.15, 0.73, 1.59\},$$

(e.g., the real seasonally adjusted GDP gap in trillions of dollars).

Recall that the parametric normal density function is given by

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}.$$

- (i) Compute and graph the parametric density function for this data (i.e., compute $\hat{\mu}$ and $\hat{\sigma}^2$) assuming an underlying normal distribution.

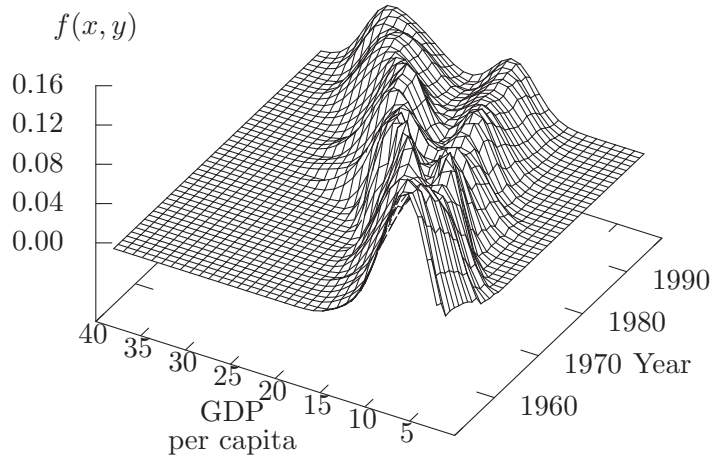


Figure 1.9: Evolution of the income distribution in Italy, 1951–1998 (series of univariate cross section kernel estimates).

- (ii) Compute and graph a histogram for this data using bin widths of 0.5 ranging from -1.5 through 2.0.

Recall that the kernel estimator of a univariate density function for continuous data can be expressed as

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)$$

and that a common (optimal) kernel is the Epanechnikov kernel given by

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

where h is a smoothing parameter restricted to lie in the range $(0, \infty]$.

- (iii) Using the same tiny sample of data, compute the kernel estimator of the density function for every sample realization using the bandwidth $h = 1.5$. Show all steps.
- (iv) Using the same data, compute the kernel estimator of the density function for every sample realization using the bandwidth $h = 0.5$. Show all steps.
- (v) On the same axes, graph your estimates of the density functions using a smooth curve to “connect the dots” for each function.
- (vi) Describe the effect of *increasing* the smoothing parameter on the estimated density function.

Exercise 1.2. Let \hat{p} be defined as in (1.1). Show that \hat{p} is the maximum likelihood estimator of $p = P(H)$.

Hint: Define $X_i = 1$ if the i th trial is H , and $X_i = 0$ if it is T . Then the likelihood function is $\prod_{i=1}^n f(X_i) = \prod_{i=1}^n p^{X_i} (1-p)^{1-X_i}$. The log-likelihood function is $\ln L = (\sum_{i=1}^n X_i) \ln p + [\sum_{i=1}^n (1 - X_i)] \ln(1 - p)$.

Exercise 1.3.

- (i) Show that $\text{MSE}(\hat{p}_n) = p(1-p)/n$, where $p = P(H)$.
- (ii) Show that $\text{plim}_{n \rightarrow \infty} \hat{p} = p$.
- (iii) Supposing that $p = P(H) \in (0, 1)$, show that the ordinary limit $\lim_{n \rightarrow \infty} \hat{p}$ does not exist.

Note that the ordinary limit is defined as follows. Letting a_n be a sequence of real numbers, we write $\lim_{n \rightarrow \infty} a_n = c$ if for all (small) $\epsilon > 0$, there exists a positive integer n_0 such that $|a_n - c| < \epsilon$ for all $n \geq n_0$.

Hint: For (ii) use the result from (i) along with Theorem A.3 of Appendix A.

Hint: For (iii) argue by contradiction.

Exercise 1.4. Let $F(x)$ be defined as in (1.2).

- (i) Show that $\text{MSE}[F_n(x)] = O(n^{-1})$ (note that this implies that $F_n(x) - F(x) = O_p(n^{-1/2})$ by Theorem A.7 of Appendix A.

(ii) Prove that

$$\sqrt{n}(F_n(x) - F(x)) \xrightarrow{d} N(0, F(x)(1 - F(x))).$$

Hint: First show that $E[F_n(x)] = F(x)$ and $\text{var}(F_n(x)) = F(x)(1 - F(x))$. Then use the Lindeberg-Levy CLT.

Exercise 1.5. Prove (1.13) under the assumption that $f(x)$ has a continuous second order derivative at x .

Hint: Use the dominated convergence theorem given by Lemma A.26 in Appendix A.

Exercise 1.6. Write $k_{ij} = k\left(\frac{X_i - X_j}{h}\right)$ and $\bar{k}_{ij} = \bar{k}\left(\frac{X_i - X_j}{h}\right)$. Using $n^{-2} = (n(n-1))^{-1} - (n^2(n-1))^{-1}$, we obtain from (1.23)

$$\begin{aligned} CV_f(h) &= \frac{1}{n(n-1)h} \sum_{i=1}^n \sum_{j=1}^n \bar{k}_{ij} - \frac{2}{n(n-1)h} \sum_{i=1}^n \sum_{j \neq i}^n k_{ij} \\ &\quad - \frac{1}{n^2(n-1)h} \sum_{i=1}^n \sum_{j=1}^n \bar{k}_{ij} \\ &= \frac{1}{(n-1)h} \bar{k}(0) + \frac{1}{n(n-1)h} \sum_{i=1}^n \sum_{j \neq i}^n [\bar{k}_{ij} - 2k_{ij}] + O_p(n^{-1}) \\ &= \frac{\kappa}{(n-1)h} + J_n + O_p(h(nh)^{-1}), \end{aligned} \tag{1.63}$$

where $J_n = [n(n-1)h]^{-1} \sum_{i=1}^n \sum_{j \neq i}^n [\bar{k}_{ij} - 2k_{ij}]$ and $\kappa = \int k^2(v) dv \equiv \bar{k}(0)$.

- (i) Show that $E(J_n) = B_0 + B_1 h^4 + O(h^5)$, where $B_0 = -\int f(x)^2 dx$, and $B_1 = (\kappa_2^2/4)\{\int [f^{(2)}(x)]^2 dx\}$.
- (ii) Accept the fact that $J_n = E(J_n) + \text{smaller order terms}$. So, asymptotically, minimizing $CV_f(h)$ is equivalent to minimizing $I(h) \stackrel{\text{def}}{=} (nh)^{-1} \kappa + E(J_n)$. Obtain that \hat{h} which minimizes $I(h)$.
- (iii) Assume that $k(0) \geq k(v)$ for all v (which is usually true for kernel estimation). If we do not use the leave-one-out estimator, then we would instead have the objective function $V(h) \stackrel{\text{def}}{=} (nh)^{-1} [\kappa - 2k(0)] + E(J_n)$. Show that $h = 0$ minimizes $V(h)$, which obviously violates the requirement that $nh \rightarrow \infty$ as $n \rightarrow \infty$. This shows that we *must* use the leave-one-out estimator when constructing $CV_f(h)$.

(iv) In deriving (1.63) we used

$$A_n \stackrel{\text{def}}{=} (n^2(n-1)h)^{-1} \sum_i \sum_{j \neq i} \bar{k}((X_i - X_j)/h) = O_p(n^{-1}).$$

Prove this result.

(v) Using the U-statistic H-decomposition given in Lemma A.15 of Appendix A, show that $J_n = E(J_n) + O_p(h^{1/2}(nh)^{-1} + n^{-1/2}h^4) +$ terms unrelated to h . Therefore, we indeed have $J_n = E(J_n) + (s.o.)$ (for a given value of h).

Hints: Note that $\bar{k}(\cdot)$ is also a nonnegative, symmetric PDF, i.e., $\int \bar{k}(v) dv = 1$, $\int v^s \bar{k}(v) dv = 0$ when s is an odd positive integer.

(i)

$$\begin{aligned} E[\bar{k}_{12}] &= h^{-1} E \int k\left(\frac{X_1 - x}{h}\right) k\left(\frac{X_2 - x}{h}\right) dx \\ &= h^{-1} \int \left[Ek\left(\frac{X_1 - x}{h}\right) \right] \left[Ek\left(\frac{X_2 - x}{h}\right) \right] dx \\ &= h^{-1} \int \left[Ek\left(\frac{X_1 - x}{h}\right) \right]^2 dx \\ &= h \int \left[f(x) + 0 + (\kappa_2/2)f^{(2)}(x)h^2 + 0 \right. \\ &\quad \left. + (\kappa_4/4!)f^{(4)}(x)h^4 + O(h^5) \right]^2 dx. \end{aligned}$$

$$\begin{aligned} E[k_{12}] &= Ek\left(\frac{X_1 - X_2}{h}\right) \\ &= \int k\left(\frac{x_1 - x_2}{h}\right) f(x_1)f(x_2) dx_1 dx_2 \\ &= h \int f(x) \left\{ f(x) + 0 + (\kappa_2/2)f^{(2)}(x)h^2 + 0 \right. \\ &\quad \left. + (\kappa_4/4!)f^{(4)}(x)h^4 \right\} dx + O(h^5). \end{aligned}$$

(ii) Note that $\bar{k}(0) = \int k^2(v) dv > 0$.

(iii) Show that $h \rightarrow 0$ produces a value of the objective function $V(h) = -\infty$. Thus, $h = 0$ minimizes $V(h)$.

- (iv) Show that $E[|A_n|] = E(A_n) = O(n^{-1})$, then apply Theorem A.7.
- (v) Using the U-statistic H-decomposition (again, see Appendix A), show that the last two terms in the H-decomposition are of order $O_p(n^{-1/2}h^4)$ (plus terms unrelated to h) and $O_p(h^{1/2}(nh)^{-1})$, respectively.

Exercise 1.7. Derive (1.27), i.e., show that

$$\int_{-\infty}^x \hat{f}(v) dv = n^{-1} \sum_{i=1}^n G\left(\frac{x - X_i}{h}\right).$$

Hint: Use $\hat{f}(v) = (nh)^{-1} \sum_{i=1}^n k\left(\frac{X_i - v}{h}\right)$ and do a change of variable $(x_i - v)/h = t$ and $dx_i = h dv$.

Exercise 1.8.

- (i) Discuss the relationship between the kernel and empirical CDF estimators, i.e., $\hat{F}(x)$ and $F_n(x) = n^{-1} \sum_{i=1}^n \mathbf{1}(X_i \leq x)$.
- (ii) Discuss whether or not one can use $h = 0$ in $\hat{F}(x)$ defined in (1.27), i.e., can one let $h \rightarrow 0$ arbitrarily fast in $\hat{F}(x)$?
- (iii) $\hat{F}(x)$ and $F_n(x)$ have the same asymptotic distribution. What is the advantage of using $\hat{F}(x)$ over $F_n(x)$? Which estimator do you expect to have smaller *finite-sample* MSE? Explain.

Exercise 1.9. Derive (1.33).

Hint: Write $\mathbf{1}_i(x) = \mathbf{1}(X_i \leq x)$ and $G_{x,x_j} = G((x - X_j)/h)$, then

$$\begin{aligned} E[CV_F(h)] &= \frac{1}{n(n-1)^2} \sum_{i=1}^n \sum_{j \neq i}^n \sum_{l \neq i}^n \int E \{ [\mathbf{1}_i(x) - G_{x,x_j}] \\ &\quad \times [\mathbf{1}_i(x) - G_{x,x_l}] \} dx \\ &= \frac{1}{n-1} \int E \{ [\mathbf{1}_i(x) - G_{x,x_j}]^2 \} dx \\ &\quad + \frac{n-2}{n-1} \int E \{ E[\mathbf{1}_i(x) - G_{x,x_j} | X_i] \}^2 dx \\ &= CV_1 + CV_2, \end{aligned}$$

then show that

$$CV_1 = (n-1)^{-1} \left\{ 2 \int F(1-F) dx - C_1 h + O(h^2) \right\} \quad \text{and}$$

$$CV_2 = \left[1 - \frac{1}{n-1} \right] \left\{ \int F(x)(1-F(x)) dx + h^4 \int C_2(x)^2 dx \right\}.$$

Exercise 1.10. Define a $q \times q$ matrix A with its (t, s) th element given by $A_{t,s} = (\kappa_2/2) \int B_t(x)B_s(x) dx$.

- (i) Show that A is positive semidefinite.
- (ii) Show that if A is positive definite, then the a_s^0 's defined in (1.41) are all uniquely determined, positive, and finite.

A necessary condition for A to be positive definite is that $f_{ss}(x)$ is not a zero function for all $s = 1, \dots, q$.

Hint:

- (i) Note that for any $q \times 1$ vector $z = (z_1, \dots, z_q)'$ that $z'Az = \int [\sum_{s=1}^q B_s(x)z_s]^2 dx \geq 0$.
- (ii) Define $z_s = a_s^2$, then $\chi_f = z'Az + \kappa^q / \sqrt{z_1 \dots z_q}$, and let z_s^0 denote values of z_s that minimize χ_f . It is easy to argue that $\infty > \inf_{z_1, \dots, z_q} \chi_f > 0$. This implies that $z_s^0 > 0$ for all s . The fact that A is positive definite implies that $z_s^0 < \infty$ for all s . Finally, z_s^0 is uniquely determined by a result given in Li and Zhou (2005). Therefore, $a_s^0 = \sqrt{z_s^0}$ is uniquely determined, positive and finite for all $s = 1, \dots, q$.

Note that A being positive definite is a sufficient condition. Li and Zhou (2005) provide a weaker necessary and sufficient condition for this result.

Exercise 1.11. Prove (1.36) and (1.37).

Hint: For a multivariate Taylor expansion, we have $f(x_0 + x) = f(x_0) + \sum_{s=1}^q f_s(x_0)(x_s - x_{s0}) + (1/2) \sum_{s=1}^q \sum_{s'=1}^q f_{ss'}(\tilde{x})(x_s - x_{s0})(x_{s'} - x_{s'0})$, \tilde{x} is on the line segment between x and x_0 .

Exercise 1.12. For the multivariate case, we have

$$CV_f(h_1, \dots, h_q) = \frac{\kappa^q}{nh_1 \dots h_q} + J_n + O_p \left((n^2 h_1 \dots h_q)^{-1} \right)$$

where $J_n = [n(n-1)]^{-1} \sum_i \sum_{j \neq i} [\bar{K}_n(X_i, X_j) - 2K_n(X_i, X_j)]$.

- (i) Show that $E(J_n) = \int [\sum_{s=1}^q B_s(x)h_s^2]^2 dx + o(\sum_{s=1}^q h_s^4)$, where the definition of $B_s(x)$ is given in Section 1.8.
- (ii) Use the U-statistic H-decomposition to show that (ignoring the term unrelated to the H_s 's)

$$J_n = E(J_n) + O_p \left(n^{-1/2} \left(\sum_{s=1}^q h_s^2 \right)^2 \right) + O_p \left((h_1 \dots h_q)^{1/2} (nh_1 \dots h_q)^{-1} \right).$$

Note that (i) and (ii) together imply that

$$CV_f = \sum_{s=1}^q B_s h_s^4 + \kappa^q (nh_1 \dots h_q)^{-1} + o_p(\eta_2^2 + \eta_1)$$

where $\eta_2 = \sum_{s=1}^q h_s^2$ and $\eta_1 = (nh_1 \dots h_q)^{-1}$.

Hint: Using H-decomposition, show that the second moments of the second and third terms are of order $O(n^{-1/2}\eta_2^2)$ and $O((h_1 \dots h_q)\eta_1^2)$, respectively.

Exercise 1.13. Assuming that $X \in [0, 1]$ and $f(0) > 0$, show that $E[\hat{f}(0)] = f(0)/2 + O(h)$ so that $\hat{f}(0)$ is a biased estimator of $f(0)$ even asymptotically.

Hint: $\hat{f}(0) = (nh)^{-1} \sum_{i=1}^n k((X_i - 0)/h)$, and

$$\begin{aligned} E[\hat{f}(0)] &= h^{-1} E[k(X_i/h)] = h^{-1} \int_0^1 f(x_1) k(x_1/h) dx_1 \\ &= \int_0^{1/h} f(hv) k(v) dv \\ &\rightarrow f(0) \int_0^\infty k(v) dv = f(0)/2. \end{aligned}$$

Exercise 1.14. With the boundary-corrected kernel defined in (1.43), and with $\hat{f}(x)$ defined in (1.44) and with the support of X being $[0, 1]$, show that for $x \in [0, h]$ at the boundary region, we have $E[\hat{f}(x)] = f(x) + O(h)$. Explicitly state the conditions that you need to derive this result.

Therefore, $\text{bias}[\hat{f}(x)] = O(h) \rightarrow 0$ as $n \rightarrow \infty$, and the boundary-corrected kernel restores the asymptotic unbiasedness for $\hat{f}(x)$ for x at the boundary region.

Hint: One can write $x = \alpha h$ with $0 \leq \alpha \leq 1$. One can assume that $|f(x) - f(z)| \leq C|x - z|$ for all $x, z \in [0, 1]$, where C is a positive constant. Then

$$\begin{aligned} E[\hat{f}(x)] &= h^{-1} \int_{-x/h}^{\infty} \frac{k\left(\frac{x_1-x}{h}\right)}{\int_{-x/h}^{\infty} k(v) dv} f(x_1) dx_1 \\ &= \int_{-\alpha}^{\infty} \frac{k(w)}{\int_{-\alpha}^{\infty} k(v) dv} f(\alpha h + wh) dw \quad (\text{used } x/h = \alpha \text{ and } x = \alpha h) \\ &= f(0) \frac{\int_{-\alpha}^{\infty} k(w) dw}{\left[\int_{-\alpha}^{\infty} k(v) dv\right]} + O(h) \\ &= f(0) + O(h). \end{aligned}$$

Exercise 1.15. With a ν th order kernel, prove (1.46) and (1.47) for the univariate x case (i.e., $q = 1$).

Exercise 1.16. Intuitively, one might think that when $f(x)$ is a uniform density, say on $[0, 1]$, then one can choose a nonshrinking value of h to estimate $f(x)$ for some $x \in [0, 1]$ (i.e., h does not go to zero as $n \rightarrow \infty$). This intuition is correct when x is an interior point of $[0, 1]$. However, at (or near) the boundary of $[0, 1]$, estimation bias will not go to zero even for uniform $f(x)$.

- (i) Show that if h does not converge to 0 as $n \rightarrow \infty$, then $\int_0^1 [\hat{f}(x, h) - f(x)]^2 dx$ will not go to zero, where $f(x)$ is the uniform PDF.
- (ii) Show that if $h \rightarrow 0$ as $n \rightarrow \infty$, then $\int_0^1 [\hat{f}(x, h) - f(x)]^2 dx \rightarrow 0$ as $n \rightarrow \infty$, where $f(x)$ is the uniform PDF.

(i) and (ii) above explain why the cross-validated selection of h , \hat{h} , must converge to zero as $n \rightarrow \infty$, and why one does not need the condition that $f^{(2)}(x)$ is not a zero function. Of course when $f(x)$ is a uniform PDF, \hat{h} will no longer have the usual order ($n^{-1/5}$). Instead it has an order equal to $n^{-1/3}$ since the bias now is of order h rather than h^2 .

Exercise 1.17. Consider the Italian income data from Section 1.13.5. For the two samples of size $n = 21$ for the years 1951 and 1998, compute the density estimates using the reference rule-of-thumb in (1.17) presuming an underlying normal distribution. How many times larger than this would the bandwidth have to be to remove the bimodal feature

present in the 1998 sample? Next, compute the density estimates using least squares cross-validation. Presuming that these bandwidths represent the “optimal” bandwidths, how much larger would the bandwidth for 1998 have to be to produce an apparently unimodal distribution? Finally, compare your least squares cross-validated density estimates with a naïve histogram. Do your estimates appear to be sensible, i.e., do they reflect features you believe are in fact present in the data?

