

# Lecture 12

## Nonparametric Regression

1

### Non Parametric Regression: Introduction

- The goal of a regression analysis is to produce a reasonable analysis to the unknown response function  $f$ , where for  $N$  data points  $(X_i, Y_i)$ , the relationship can be modeled as

$$y_i = m(x_i) + \varepsilon_i, \quad i = 1, \dots, N$$

- Note:  $m(\cdot) = E[y|x]$       if  $E[\varepsilon|x]=0$  –i.e.,  $\varepsilon \perp x$

- We have different ways to model the conditional expectation function (CEF),  $m(\cdot)$ :
  - Parametric approach
  - Nonparametric approach
  - Semi-parametric approach.

2

## Non Parametric Regression: Introduction

- Parametric approach:  $m(\cdot)$  is known and smooth. It is fully described by a finite set of parameters, to be estimated. Easy interpretation. For example, a linear model:

$$y_i = x_i' \beta + \varepsilon_i, \quad i = 1, \dots, N$$

- Nonparametric approach:  $m(\cdot)$  is smooth, flexible, but unknown. Let the data determine the shape of  $m(\cdot)$ . Difficult interpretation.

$$y_i = m(x_i) + \varepsilon_i, \quad i = 1, \dots, N$$

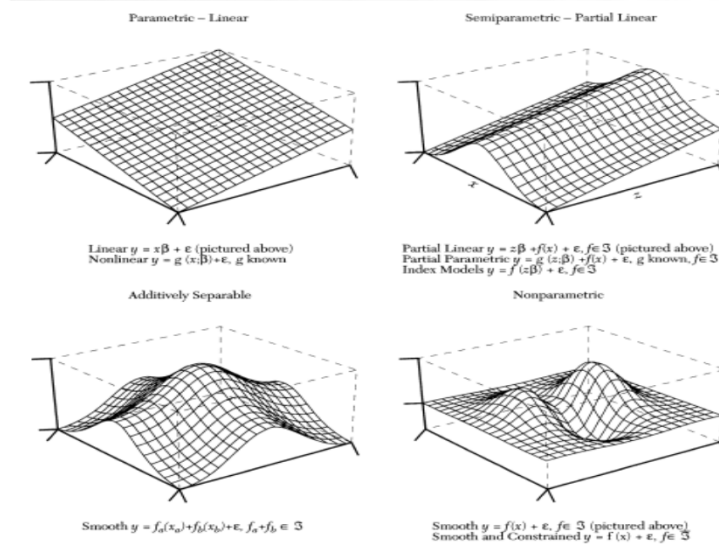
- Semi-parametric approach:  $m(\cdot)$  have some parameters -to be estimated-, but some parts are determined by the data.

$$y_i = x_i' \beta + m_z(z_i) + \varepsilon_i, \quad i = 1, \dots, N$$

3

## Non Parametric Regression: Introduction

Figure 2. Categorization of Regression Functions



$\mathfrak{F}$  is a smooth family of functions.  $\mathfrak{F}$  is a smooth family with additional constraints such as monotonicity, concavity, symmetry or other constraints.

4

## Regression: Smoothing

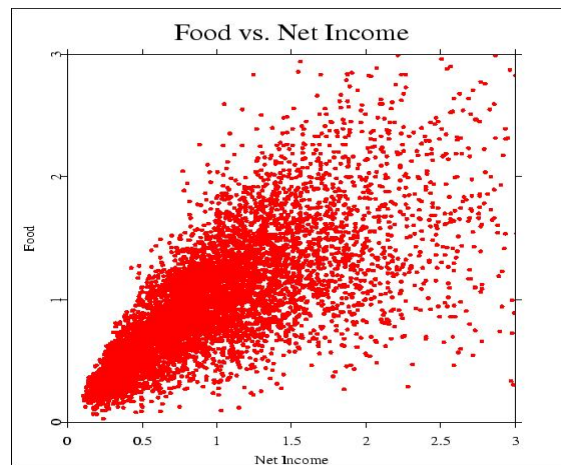
- We want to relate  $y$  with  $x$ , without assuming any functional form. First, we consider the one regressor case:

$$y_i = m(x_i) + \varepsilon_i, \quad i = 1, \dots, N$$

- In the CLM, a linear functional form is assumed:  $m(x_i) = x_i'\beta$ .
- In many cases, it is not clear that the relation is linear.
- Non-parametric models attempt to discover the (approximate) relation between  $y_i$  and  $x_i$ . Very flexible approach, but we need to make some assumptions.

5

## Regression: Smoothing



- The functional form between income and food is not clear from the scatter plot. From Hardle (1990).

## Regression: Smoothing

- A reasonable approximation to the regression curve  $m(x_i)$  will be the mean of response variables near a point  $x_i$ . This *local averaging* procedure can be defined as

$$\hat{m}(x) = N^{-1} \sum_{i=1}^N W_{N,h,i}(x) y_i$$

- The averaging will *smooth* the data. The weights depend on the value of  $x$  and on a  $h$ . Recall that as  $h$  gets smaller,  $\hat{m}(x)$  is less biased but also has greater variance.

Note: Every smoothing method to be described follows this form. Ideally, we give smaller weights for  $x$ 's that are farther from  $x$ .

- It is common to call the regression estimator  $\hat{m}(x)$  a *smoother* and the outcome of the smoothing procedure is called the *smooth*.

7

## Regression: Smoothing – Example 1

- From Hansen (2013). To illustrate the concept, suppose we use the naive histogram estimator as the basis for the weight function,  $w_i$ :

$$W_{N,h,i}(x_0) = \frac{I[|x_i - x_0| \leq h]}{\sum_{i=1}^n I[|x_i - x_0| \leq h]}$$

- Let  $x_0=2$ ,  $h=0.5$ . The estimator  $\hat{m}(x)$  at  $x=2$  is the average of the  $y_i$  for the observations such that  $x_i$  falls in the interval  $[1.5 \leq x_i \leq 2.5]$ .

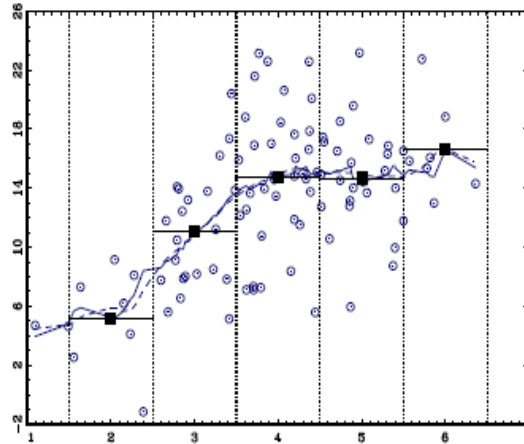
- Hansen simulates observations (see next Figure) and calculate  $\hat{m}(x)$  at  $x=2, 3, 4, 5$  &  $6$ . For example,  $\hat{m}(x=2) = 5.16$ , shown in the Figure as the first solid square.

- This process is equivalent to partitioning the support of  $x_i$  into the regions  $[1.5,2.5]$ ;  $[2.5,3.5]$ ;  $[3.5,4.5]$ ;  $[4.5,5.5]$ ; &  $[5.5,6.5]$ . It produces a step function. Reasonable behavior in the bins, but unrealistic jumps.

8

## Regression: Smoothing – Example 1

- Figure 11.1 - Simulated data and  $\hat{m}(x)$  from Hansen (2013).



- Obviously, we can calculate  $\hat{m}(x)$  at a finer grid for  $x$ . It will track the data better. But, the unrealistic jumps (discontinuities) will remain. 9

## Regression: Smoothing – Example 1

- The source of the discontinuity is the weights  $w_i$  are constructed from indicator functions, which are themselves discontinuous.
- If instead the weights are constructed from continuous functions,  $K(\cdot)$ ,  $\hat{m}(x)$  will also be continuous in  $x$ . It will produce a true *smooth*! For example,

$$W_{N,h,i}(x_0) = \frac{K\left(\frac{x_i - x_0}{h}\right)}{\sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right)}$$

- The bandwidth  $h$  determines the degree of smoothing. A large  $h$  increases the width of the bins, increasing the smoothness of  $\hat{m}(x)$ . A small  $h$  decreases the width of the bins, producing a less smooth  $\hat{m}(x)$ .

10

## Regression: Smoothing – Example 2

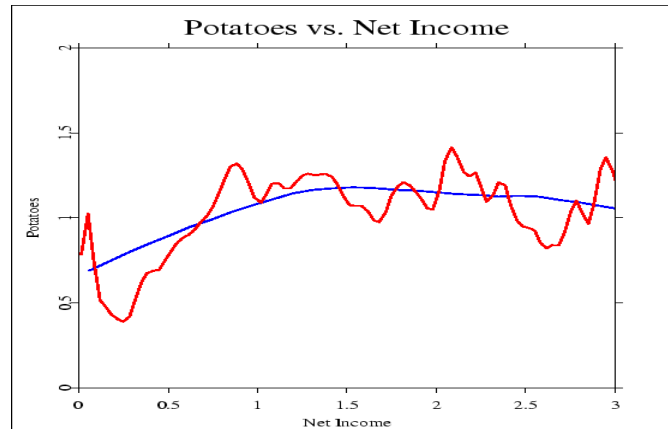


Figure 1. Expenditure of potatoes as a function of net income.

$h = 0.1, 1.0, N = 7125$ , year = 1973. Blue line is the *smooth*. From Hardle (1990).

## Regression: Smoothing - Interpretation

- Suppose the weights add up to 1 for all  $x_i$ . The  $\hat{m}(x)$  is a least squares estimates at  $x$  since we can write  $\hat{m}(x)$  as a solution to

$$\min_{\theta} N^{-1} \sum_{i=1}^N W_{N,h,i}(x)(y_i - \theta)^2$$

That is, a kernel regression estimator is a *local constant regression*, since it sets  $m(x)$  equal to a constant,  $\theta$ , in the very small neighborhood of  $x_0$ :

$$\min_{\theta} N^{-1} \sum_{i=1}^N W_{N,h,i}(x)(y_i - \theta)^2 = N^{-1} \sum_{i=1}^N W_{N,h,i}(x)(y_i - \hat{m}(x))^2$$

Note: The residuals are weighted quadratically => weighted LS!

- Since we are in a LS world, outliers can create problems. Robust techniques can be better.

## Regression: Smoothing - Issues

- Q: What does smoothing do to the data?
  - (1) Since averaging is done over neighboring observations, an estimate of  $m(\cdot)$  at peaks or bottoms will flatten them. This finite sample bias depends on the local curvature of  $m(\cdot)$ . *Solution:* Shrink neighborhood!
  - (2) At the boundary points, half the weights are not defined. This also creates a bias.
  - (3) When there are regions of sparse data, weights can be undefined – no observations to average. *Solution:* Define weights with variable span.

- Computational efficiency is important.

A naive way to calculate the smooth  $\hat{m}(x)$  consists in calculating the  $w_i(x_j)$ 's for  $j=1,\dots,N$ . This results in  $O(N^2)$  operations. If we use an iterative algorithm, calculations can take very long.

## Kernel Regression

- Kernel regressions are weighted average estimators that use kernel functions as weights.
- Recall that the *kernel*  $K$  is a continuous, bounded and symmetric real function which integrates to 1. The weight is defined by

$$W_{hi}(x) = K_h(x - X_i) / \hat{f}_h(x)$$

where  $\hat{f}_h(x) = N^{-1} \sum_{i=1}^N K_h(x - X_i)$ , and  $K_b(u) = b^{-1} K(u/b)$ ;

- The functional form of the kernel virtually always implies that the weights are much larger for the observations where  $x_i$  is close to  $x_0$ . This makes sense!

## Kernel Regression

- Standard statistical formulas allow us to calculate  $E[y | x]$ :

$$E[y | x] = m(x) = \int y f_C(y | x) dy$$

where  $f_C$  is the distribution of  $y$  conditional on  $x$ . As always, we can express this conditional distribution in several ways. In particular:

$$E[y|x] = m(x) = \frac{\int_{-\infty}^{\infty} y f_J(y, x) dy}{f_M(x)} = \frac{\int_{-\infty}^{\infty} y f_J(y, x) dy}{\int_{-\infty}^{\infty} f_J(y, x) dy}$$

where the subscripts M and J refer to the marginal and the joint distributions, respectively.

- Q: How can we estimate  $m(x)$  using these formulas?
- First, consider first  $f_M(x)$ . This is just the density of  $x$ . Estimate this using the density estimation results. For a given value of  $x$  (say,  $x_0$ ) as:

$$\hat{f}_M(x_0) = \hat{f}(x_0) = (Nh)^{-1} \sum_{i=1}^N K\left(\frac{x_i - x_0}{h}\right)$$

## Kernel Regression : *Nadaraya-Watson estimator*

- First, consider first  $f_M(x)$ :  $\hat{f}_M(x) = (Nh)^{-1} \sum_{i=1}^N K\left(\frac{x_i - x_0}{h}\right)$

- Second, consider  $\int f_J(y, x_0) dy = (Nh)^{-1} \sum_{i=1}^N K\left(\frac{x_i - x_0}{h}\right)$

which suggests  $\int y f_J(y, x_0) dy = (Nh)^{-1} \sum_{i=1}^N y_i K\left(\frac{x_i - x_0}{h}\right)$

- Plugging these two kernel estimates of the terms in the numerator and the denominator of the expression for  $m(x)$  gives the *Nadaraya-Watson (NW) kernel estimator*:

$$\hat{m}(x_0) = \frac{\sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right) y_i}{\sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right)}$$



### Kernel Regression: NW estimator - Different $K(\cdot)$

- The shape of the kernel weights is determined by  $K$  and the size of the weights is parameterized by  $h$  ( $h$  plays the usual smoothing role).
- The normalization of the weights  $\hat{f}_h(x) = N^{-1} \sum_{i=1}^N K_h(x - X_i)$  is called the *Rosenblatt-Parzen kernel density estimator*. It makes sure that the weights add up to 1.
- Two important constants associated with a kernel function  $K(\cdot)$  are its variance  $\sigma_K^2 = d_K$  and roughness  $c_K$ , (also denoted  $R_K$ ), which are defined as:

$$d_K = \int z^2 K(z) du$$

$$c_K = \int K^2(z) dz$$

### Kernel Regression: NW estimator - Different $K(\cdot)$

- Many  $K(\cdot)$  are possible. Practical and theoretical considerations limit the choices. Usual choices: Epanechnikov, Gaussian, Quartic (biweight), and Tricube (triweight).

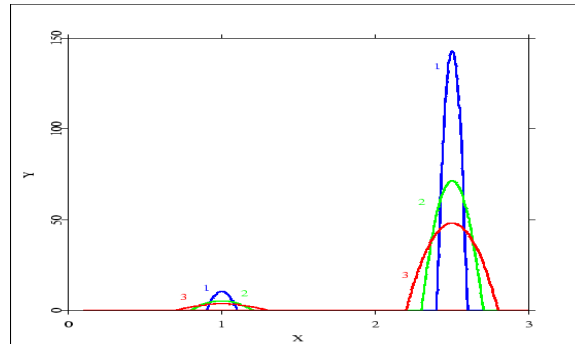
Kernel	Equation	$R_k$	$\sigma_k^2$
Uniform	$k_0(u) = \frac{1}{2} 1( u  \leq 1)$	1/2	1/3
Epanechnikov	$k_1(u) = \frac{3}{4} (1 - u^2) 1( u  \leq 1)$	3/5	1/5
Biweight	$k_2(u) = \frac{15}{16} (1 - u^2)^2 1( u  \leq 1)$	5/7	1/7
Triweight	$k_3(u) = \frac{35}{32} (1 - u^2)^3 1( u  \leq 1)$	350/429	1/9
Gaussian	$k_\phi(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right)$	$1/(2\sqrt{\pi})$	1

- Figure 11.1 shows the NW estimator with Epanechnikov kernel and  $h=0.5$  with the dashed line. (The full line uses a uniform kernel.)

- Recall that the Epanechnikov kernel enjoys optimal properties.

### Kernel Regression: Epanechnikov kernel.

Figure 3. The effective kernel weights for the food/ income data: At  $x=1$  and  $x=2.5$  for  $h = 0.1$  (label 1, blue),  $h = 0.2$  (label 2, green),  $h = 0.3$  (label 3, red) with *Epanechnikov kernel*. From Hardle (1990).



- The smaller  $h$ , the more concentrated the  $w_{i's}$ . In sparse regions, say  $x=2.5$  (low marginal pdf), it gives more weight to observations around  $x$ .

### Kernel Regression: NW estimator - Properties

- The *NW estimator* is defined by

$$\hat{m}_h(x) = \frac{\sum_{i=1}^N y_i K_h(x - X_i)}{\sum_{i=1}^N K_h(x - X_i)} = \sum_{i=1}^N w_{N,h,i}(x) y_i$$

- Similar situation as in KDE: No finite sample distribution theory for  $\hat{m}(x)$ . All statistical properties are based on asymptotic theory.

- Details. One regressor ( $d=1$ ), but straightforward to generalize.

Fix  $x$ . Note that  $y_i = m(x_i) + \varepsilon_i = m(x) + (m(x_i) - m(x)) + \varepsilon_i$

Then,

$$\begin{aligned} \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x_i - x}{h}\right) y_i &= \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x_i - x}{h}\right) [m(x) + (m(x_i) - m(x)) + \varepsilon_i] \\ &= \hat{f}(x) m(x) + \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x_i - x}{h}\right) (m(x_i) - m(x)) + \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x_i - x}{h}\right) \varepsilon_i \\ &= \hat{f}(x) m(x) + \hat{m}_1(x) + \hat{m}_2(x) \end{aligned}$$

## Kernel Regression: NW estimator - Properties

- It follows that

$$\hat{m}(x) = m(x) + \hat{m}_1(x) / \hat{f}(x) + \hat{m}_2(x) / \hat{f}(x)$$

(1)  $\hat{m}_2(x)$ .

- Mean. Since  $E[\varepsilon_i | x_i] = 0 \Rightarrow E[\hat{m}_2(x)] = 0$ .

- Variance.

$$\text{var}[\hat{m}_2(x)] = \frac{1}{Nh^2} E\left[K\left(\frac{x_i - x}{h}\right)\varepsilon_i\right]^2 = \frac{1}{Nh^2} E\left[K\left(\frac{x_i - x}{h}\right)^2 \sigma^2(x_i)\right]$$

(by conditioning), and then

$$\text{var}[\hat{m}_2(x)] = \frac{1}{Nh^2} \int K\left(\frac{z-x}{h}\right)^2 \sigma^2(z) f(z) dz$$

Change of variables,  $(z-x)/h = u$ , and assume  $\sigma^2(x)$  and  $f(x)$  are smooth:

$$\begin{aligned} \text{var}[\hat{m}_2(x)] &= \frac{1}{Nh^2} \int K(u)^2 \sigma^2(x+hu) f(x+hu) (hdu) \\ &= \frac{1}{Nh} \int K(u)^2 \sigma^2(x) f(x) (du) + o\left(\frac{1}{Nh}\right) = \frac{\sigma^2(x) f(x)}{Nh} c_k + o\left(\frac{1}{Nh}\right) \end{aligned}$$

## Kernel Regression: NW estimator - Properties

- We can apply the CLT to obtain that as  $h \rightarrow 0$ , and  $Nh \rightarrow \infty$

$$\sqrt{Nh} \hat{m}_2(x) \xrightarrow{d} N(0, \sigma^2(x) f(x) c_k)$$

(1)  $\hat{m}_1(x)$ .

- Mean.

$$\begin{aligned} E[\hat{m}_1(x)] &= \frac{1}{h} E\left[K\left(\frac{x_i - x}{h}\right)(m(x_i) - m(x))\right] = \frac{1}{h} \int K\left(\frac{z-x}{h}\right)(m(z) - m(x)) f(z) dz \\ &= \int K(u)(m(x+hu) - m(x)) f(x+hu) du \end{aligned}$$

Expand  $m(x+hu)$  and  $f(x+hu)$  into (2nd- and 1st-order, respectively)

Taylor expansions around  $x$ : Up to  $o(h^2)$  we get:

$$\begin{aligned} E[\hat{m}_1(x)] &= \int K(u)(m(x+hu) - m(x)) f(x+hu) du \\ &= \int K(u)\left(m(x) + hu m'(x) + \frac{h^2 u^2}{2} m''(x) - m(x)\right) (f(x) + hu f'(x)) du \end{aligned}$$

### Kernel Regression: NW estimator - Properties

- Then, we get

$$\begin{aligned}
 E[\hat{m}_1(x)] &= \int K(u)(hum'(x) + \frac{h^2u^2}{2}m''(x))(f(x) + huf'(x))du \\
 &= hm'(x)f(x)\int K(u)udu + h^2[\frac{1}{2}m''(x)f(x) + m'(x)f'(x)]\int K(u)u^2du + o(h^3) \\
 &\approx h(m'(x)f(x))\kappa_1 + h^2[\frac{1}{2}m''(x)f(x) + m'(x)f'(x)]\kappa_2 \\
 &= h^2\kappa_2B(x)f(x)
 \end{aligned}$$

where  $B(x) = m''(x)/2 + m'(x)f'(x)/f(x)$

- Variance. A similar expansion shows that  $\text{var}[\hat{m}_1(x)]$  is  $O(h^2/Nb)$ , which is of smaller order than  $O(1/Nb)$ .

- Thus, as  $b \rightarrow 0$ , and  $Nb \rightarrow \infty$ ,  $\sqrt{Nh}[\hat{m}_1(x) - h^2d_K B(x)f(x)] \xrightarrow{p} 0$   
and since  $\hat{f}(x) \xrightarrow{p} f(x) \Rightarrow \sqrt{Nh}[\hat{m}_1(x)/\hat{f}(x) - h^2d_K B(x)f(x)] \xrightarrow{p} 0$

### Kernel Regression: NW estimator - Properties

- This bias is of size  $O(b^2)$ . Intuitively, the bias is larger the “curvier”  $m(x_0)$  is -i.e., the larger  $m'(x_0)$  and  $m''(x_0)$  are.

- The kernel regression estimator,  $\hat{m}(x)$ , is *consistent*. But, convergence is at the rate  $\text{sqrt}(Nb)$ , not the usual  $\text{sqrt}(N)$ .

- Applying the CLT, we get under general assumptions, *asymptotically normality*:

$$\sqrt{Nh}[\hat{m}(x) - m(x) - h^2d_K B(x)] \xrightarrow{d} N(0, \sigma^2(x)c_k / f(x))$$

- The MSE = variance + bias<sup>2</sup>. Given our asymptotic results, we can get the AMSE[ $\hat{m}(x)$ ]:

$$AMSE[\hat{m}(x), h] \approx (Nh)^{-1}\sigma^2(x)c_k / f(x) + h^4d_K^2B(x)^2$$

where  $d_K = \sigma_K^2$  and  $c_k$  is the roughness.

### Kernel Regression: NW estimator - Properties

- Notes about asymptotic distribution:
  - The asymptotic distribution depends on the kernel through  $c_k$  –the roughness- and  $d_k$  –the 2<sup>nd</sup> moment of  $\xi$ .
  - The optimal kernel minimizes  $c_k$ ; the same as for density estimation. Therefore, the Epanechnikov family is optimal for regression.
  - The optimal  $h$  depends on the first and second derivatives of  $m(x)$ , not on  $f(x)$ .
  - Rules of thumb for  $h$  designed for  $f(x)$  have no justification.

### Kernel Regression: NW estimator – C.I.'s

- Given the asymptotic normality, it is easy to construct C.I.'s.  
Usual steps:
  - 1) Compute  $\hat{m}(x)$ , and, using kernel density estimation,  $\hat{f}(x)$ .
  - 2) Estimate  $\sigma^2(x)$ .  $c_k$ , the roughness, can be obtained from Tables.
  - 3) Select  $\alpha\%$  level and use usual formula.
- Note that we are not estimating the bias:
 
$$\text{Bias} [\hat{m}(x)] = h^2 \kappa_2 B(x) f(x)$$
 where  $B(x) = m''(x) / 2 + m'(x) f'(x) / f(x)$
- It is complicated, since it needs estimates of derivatives. In general, it adds noise to the C.I. That is, we do not estimate an asymptotic exact C.I.

### **Kernel Regression: NW estimator - Properties**

- C.I.'s tend to be wider at the boundaries and when the data is sparse.
- Even if we compute the bias, asymptotic C.I.'s are an approximation. A bootstrap may work better.

### **Kernel Regression: NW estimator - Limitations**

- (1) Applied to truly linear data, the NW estimator can be poor.
- Let  $d=1$  and the true conditional mean is linear  $y=\alpha+\beta x$ , with no error. The behavior of the NW estimator depends on the marginal distribution of  $X$ .
  - If they are not spaced at uniform distances, then  $m^{\wedge}(x)\neq m(x)$ . The NW estimator applied to purely linear data yields a nonlinear output.
  - The choice of  $h$  may not help. As  $h$  increases, the estimator becomes a constant, not a linear function.
- (2) Poor behavior at the boundaries of  $X$ . Suppose  $m(x)$  is positively sloped, at the right boundary, the NW estimator will be upward biased. In fact, the estimator is inconsistent at the boundary.
- This restricts application of the NW estimator to interior points.

## Kernel Estimators of Derivatives

- The same principles behind kernel estimation can be used to estimate the derivatives of the regression function. These derivatives can be used to estimate partial effects.
- If the weights are sufficiently smooth and  $h$  is properly chosen, the derivative estimator is consistent.
- Taking the  $k$ -th derivative of  $\hat{m}(x)$ :

$$\hat{m}^{(k)}(x_0) = N^{-1} h^{-(k+1)} \sum_{i=1}^n w^{(k)}\left(\frac{x_i - x_0}{h}\right) y_i$$

- The kernel estimate of the  $k$ -th derivative is also a local average.

## Kernel Regression: Local Linear Estimator

- We motivated the NW estimator at  $x$  as an average of the  $y_i$ 's for observations in a neighborhood of  $x$ : A local constant approximation.
- Instead, we can do OLS in the same neighborhood. If we use a weighting function, this is called the local linear (LL) estimator.
- The idea is to fit the local model  $y_i = \alpha + (x_i - x)' \beta + \varepsilon_i$ .
- We use  $(X_i - x)$  rather than  $X_i$  to have  $m(x) = E[y_i | X_i = x] = \alpha$ .
- We do OLS with observations such that  $|X_i - x| \leq h$ . That is,

$$\min_{\alpha, \beta} N^{-1} \sum_{i=1}^N (y_i - \alpha - (x_i - x)' \beta)^2 I[|x_i - x| \leq h]$$

## Kernel Regression: Local Linear (LL) Estimator

- We have a weighted LS problem, which can be generalized to:

$$\min_{\alpha, \beta} N^{-1} \sum_{i=1}^N (y_i - \alpha - (x_i - x)\beta)^2 K\left(\frac{x_i - x}{h}\right)$$

- Then, setting  $Z_i = [1 \ (X_i - x)]'$  delivers:

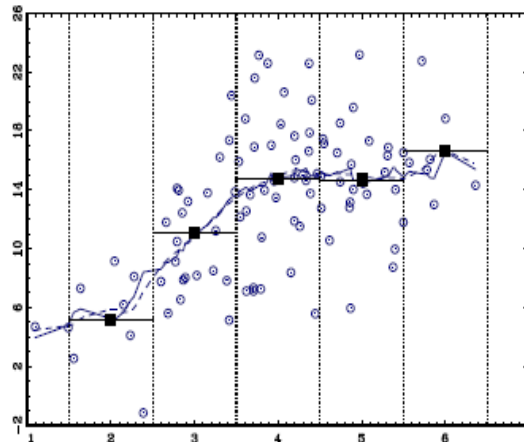
$$\begin{aligned} \begin{pmatrix} \hat{\alpha}(x) \\ \hat{\beta}(x) \end{pmatrix} &= \left( \sum_{i=1}^n 1(|X_i - x| \leq h) Z_i Z_i' \right)^{-1} \left( \sum_{i=1}^n 1(|X_i - x| \leq h) Z_i y_i \right) \\ &= \left( \sum_{i=1}^n K(H^{-1}(X_i - x)) Z_i Z_i' \right)^{-1} \left( \sum_{i=1}^n K(H^{-1}(X_i - x)) Z_i y_i \right) \end{aligned}$$

the second line is valid for any (multivariate) kernel function. This is a (locally) weighted regression of  $y_i$  on  $X_i$ .

- LL estimator preserves linear data and behaves better at the boundaries.

## Regression: LL Smoothing – Example 1

- Figure 11.1 - Simulated data and  $\hat{m}(x)$  from Hansen (2013).



- $\hat{m}(x)$  estimated under NW (dashed line) and LL (points). Overall, very similar smooths.



## Kernel Regression: LL Estimator - LOWESS

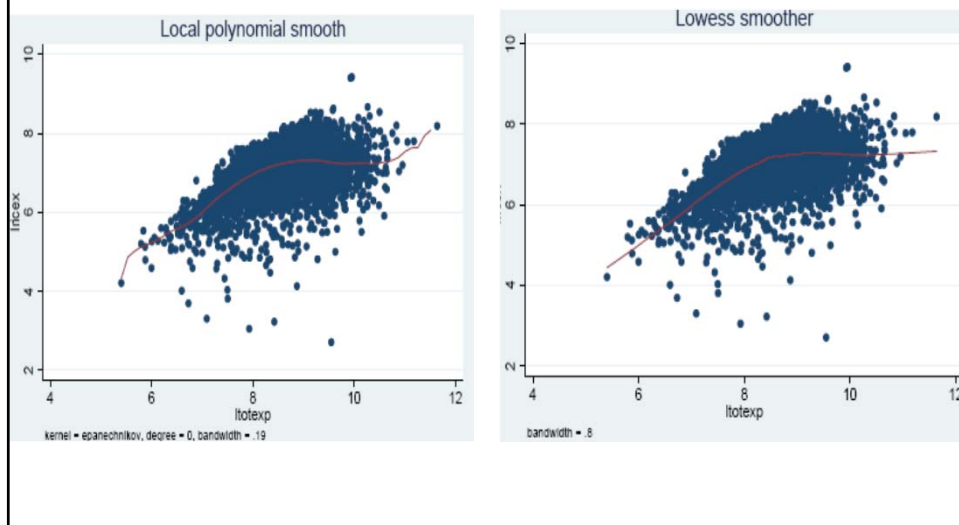
- A popular local regression estimator is *locally weighted scatterplot smoothing* (**lowess**), introduced by Cleveland (1979).
- It uses a variable  $h$ , determined by the distance from  $x_0$  to its  $k$ -th NN, and it uses a *tricubic kernel*:

$$K(\tilde{x}) = (70/81)(1 - |\tilde{x}|^3)^3 \mathbb{I}[|\tilde{x}| < 1].$$

- In principle, we can add higher order polynomial terms, which would make it easier to take higher order derivatives.

## Kernel Regression: LL Estimator - Application

- Rice expenditures as a total of Log total expenditures.



### Kernel Regression: NW or LL Estimator?

- In contrast to the NW estimator, the LL estimator preserves linearity in the data. That is, if the true data is linear, for any sub-sample, a local linear regression fits exactly, so  $\hat{m}(x_0) = m(x_0)$ .
- As  $h \rightarrow \infty$ , the LL estimator collapses to OLS of  $y_i$  on  $X_i$ . That is, we can think of LL as a nonparametric generalization of OLS.
- The asymptotic distribution of the LL estimator is similar to that of the NW estimator. The bias term is simpler, the  $m'(x)$  and  $f'(x)$  disappear. The asymptotic variance is the same.

Q: If LL improves on NW, why not use a local polynomial of order  $p$ ?  
It is possible and doable. In practice, when  $d > 1$ , applying polynomial methods is not easy.

### Kernel Regression: NW or LL Estimator?

- Strictly speaking, we cannot rank the AMSE of the NW versus the LL estimator.
- The AMSE of the LL estimator only depends on  $m''(x)$ ; while that of the NW estimator also depends on  $m'(x)$ . We expect this to translate into reduced bias.
- Since both estimators have the same asymptotic variance, the statistics literature prefers the LL estimator.
- According to Bruce Hansen (2013), caution is warranted. In simple simulations, the LL estimator does not always beat the NW estimator.

## Kernel Regression: NW or LL Estimator?

- Hansen's interesting findings:
  - When the regression function  $m(x)$  is quite flat, the NW estimator does better. When the regression function is steeper and curvier, the LL estimator tends to do better.
  - Intuition from above result: In finite samples the NW estimator tends to have a smaller variance. An advantage in contexts where estimation bias is low (such as when the regression function is flat).

Note: In many economic contexts, it is believed that the regression function may be quite flat with respect to many regressors. In this context it may be better to use NW rather than LL.

## Kernel Regression: Weighted NW estimator

- Hall (1999) proposed a weighted NW estimator is defined by

$$\hat{m}_h(x_0) = \frac{\sum_{i=1}^N p_i(x_0) K_h\left(\frac{x_0 - X_i}{h}\right) y_i}{\sum_{i=1}^N p_i(x_0) K_h\left(\frac{x_0 - X_i}{h}\right)}$$

where  $p_i(x)$  are weights. The weights satisfy:

$$p_i(x) \geq 0$$

$$\sum_i p_i(x) = 1.$$

$$\sum_i p_i(x) K(b^{-1}(x; x_i)) (x_i - x) = 1.$$

- The first two requirements define  $p_i(x)$  as weights. The third equality requires the weights to force the kernel function to satisfy local linearity.

### Kernel Regression: Weighted NW estimator

- The weights are determined by empirical likelihood. Specifically, for each  $x$ ; you maximize  $\sum_i \ln p_i(x)$  s.t. the above constraints.
- The solutions take the form

$$p_i(x) = \frac{1}{n (1 + \lambda' (X_i - x) K (H^{-1} (X_i - x)))}$$

where  $\lambda$  is a LM, found by numerical optimization.

- The estimator  $\hat{m}(x)$  has the same asymptotic distribution as the LL estimator. When  $y_i \geq 0$ ; the standard and weighted NW estimators also satisfy  $\hat{m}(x) \geq 0$ . This is good ( $m(x)$  is non-negative!). On the other side, the LL estimator is not necessarily non-negative.

### Kernel Regression: Weighted NW estimator

- When  $y_i \geq 0$ ; the standard and weighted NW estimators also satisfy  $\hat{m}(x) \geq 0$ . This is good ( $m(x)$  is non-negative!). On the other side, the LL estimator is not necessarily non-negative.
- Disadvantage: More computationally intensive than the LL estimator. The EL weights must be found separately for each  $x_0$  at which  $\hat{m}(x_0)$  is calculated.

## Kernel Regression: Residuals, Fit & CV

- We are used to use the fitted residuals to construct GOF measures. The residuals are defined as usual:

$$e_i = y_i - \hat{m}(x_i), \quad i = 1, \dots, N$$

- Problem: In general, but especially when  $h$  is small, it is hard to view  $e_i$  as a GOF measure. As  $h \rightarrow 0$ ,  $\hat{m}(\cdot) \rightarrow y_i$  (and  $e_i \rightarrow 0$ ). This indicates *overfitting* as the true error is not zero.

- Solution: Measure the fit of the regression at  $x = x_i$  by re-estimating the model excluding the  $i$ -th observation (notation: “- $i$ ,” the  $i$ -th observation excluded). We call this *leave-one-out* estimation For NW regression, we get:

$$\hat{m}_{-i}(x) = \frac{\sum_{j \neq i}^N y_j K_h(x - X_j) y_j}{\sum_{j \neq i}^N K_h(x - X_j)} = \sum_{j \neq i}^N w_{N,h,-i}(x) y_j$$

## Kernel Regression: Residuals, Fit & CV

- Now, the leave-one-out residuals are defined as:

$$e_{-i} = y_i - \hat{m}_{-i}(x_i), \quad i = 1, \dots, N$$

- $e_{-i}$  is not a function of  $y_i$ ; there is no tendency to overfit for small  $h$ :

- The mean squared leave-one-out residual is

$$CV(h) = \frac{1}{N} \sum_{i=1}^N e_{-i}(h)^2,$$

- This function of  $h$  is known as the *cross-validation criterion*. This criterion can be used to select the bandwidth.

- The CV bandwidth  $h_{CV}$  is the value that minimizes  $CV(h)$ . Usually, the restriction  $h_{CV} \geq h_{LB}$  is imposed, where  $h_{LB}$  is a lower bound for  $h_{CV}$ , to make sure the bandwidth is not too small.

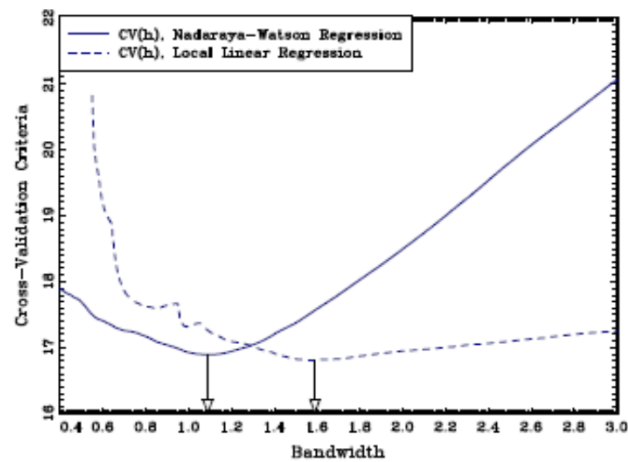
## Kernel Regression: Residuals, Fit & CV

- The CV bandwidth  $h_{CV}$  is calculated numerically.
- A grid search is popular. Plots of  $CV(b)$  against  $b$  are also used.
- It turns out that  $CV(b)$  is an estimator of the mean-squared forecast error (MSFE). That is,

$$E[CV(b)] = MSFE_{N-1}(b) = MISE_{N-1}(b) + \sigma^2$$

## Kernel Regression: Residuals, Fit & CV

- Plots of  $CV(b)$  against  $b$  for Hansen's simulated data for the NW and Local Linear estimators (with Epanechnikov kernel). From Hansen (2013).



### Kernel Regression: NW estimator - Multivariate

- The *NW estimator* is defined by

$$\hat{m}_h(x) = \frac{\sum_{i=1}^N y_i K_h(x - X_i)}{\sum_{i=1}^N K_h(x - X_i)} = \sum_{i=1}^N y_i w_{N,h,i}(x)$$

- The last expression simply shows that this estimator can be thought of as a weighted average of the observations of  $y$ . In matrix notation, we can write  $\hat{Y} = \mathbf{M}(h) \mathbf{Y}$ , with

$$M(h) = \begin{bmatrix} \frac{K(\frac{X_1 - x_1}{h})}{\sum_{i=1}^n K(\frac{X_i - x_1}{h})} & \frac{K(\frac{X_2 - x_1}{h})}{\sum_{i=1}^n K(\frac{X_i - x_1}{h})} & \dots & \frac{K(\frac{X_n - x_1}{h})}{\sum_{i=1}^n K(\frac{X_i - x_1}{h})} \\ \frac{K(\frac{X_1 - x_2}{h})}{\sum_{i=1}^n K(\frac{X_i - x_2}{h})} & \frac{K(\frac{X_2 - x_2}{h})}{\sum_{i=1}^n K(\frac{X_i - x_2}{h})} & \dots & \frac{K(\frac{X_n - x_2}{h})}{\sum_{i=1}^n K(\frac{X_i - x_2}{h})} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{K(\frac{X_1 - x_n}{h})}{\sum_{i=1}^n K(\frac{X_i - x_n}{h})} & \frac{K(\frac{X_2 - x_n}{h})}{\sum_{i=1}^n K(\frac{X_i - x_n}{h})} & \dots & \frac{K(\frac{X_n - x_n}{h})}{\sum_{i=1}^n K(\frac{X_i - x_n}{h})} \end{bmatrix}$$

### Kernel Regression: NW estimator - Multivariate

- Kernel regression predictions:  $\hat{Y} = \mathbf{M}(h) \mathbf{Y}$
- Liner regression predictions:  $\hat{Y} = \mathbf{P}_x \mathbf{Y}$ .
- A multivariate kernel is constructed, row by row, by computing the product of marginal densities for each variable in the matrix of regressors  $X$ . That is,

$$h^{-d} K\left(\frac{X - x_i}{h}\right) = \prod_{j=1}^d h^{-1} K\left(\frac{x_j - x_{ji}}{h}\right)$$

- Usually, we use *leave-one-out* kernels. That is, the current observation is excluded in the kernel construction to avoid overfitting — the principal diagonal in  $M(h)$  is zeroes.

## Comparison: Mean vs Kernel Smoother

- Mean (uniform) smoother

$$\hat{m}(x) = \frac{\sum_{i=1}^N w\left(\frac{x-x_i}{h}\right) x_i}{\sum_{i=1}^N w\left(\frac{x-x_i}{h}\right)}$$

where

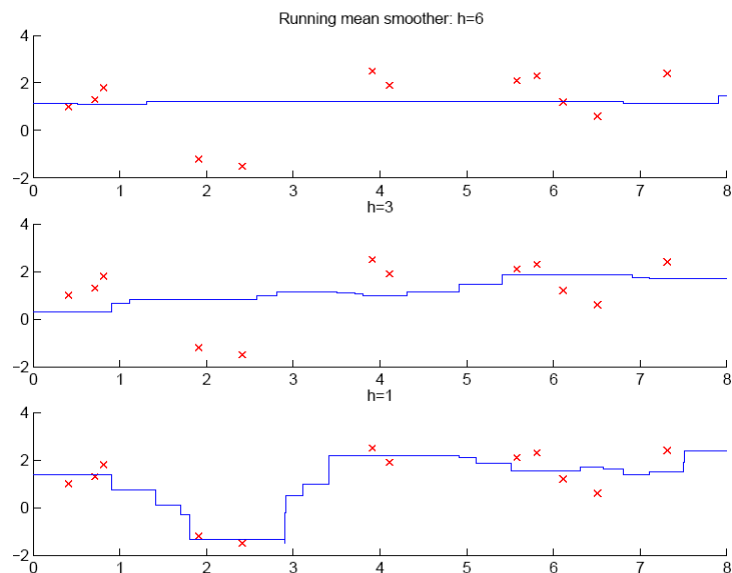
$$w(u) = \begin{cases} 1 & \text{if } |u| < 1 \\ 0 & \text{otherwise} \end{cases}$$

- Kernel smoother

$$\hat{m}(x) = \frac{\sum_{i=1}^N K\left(\frac{x-x_i}{h}\right) x_i}{\sum_{i=1}^N K\left(\frac{x-x_i}{h}\right)}$$

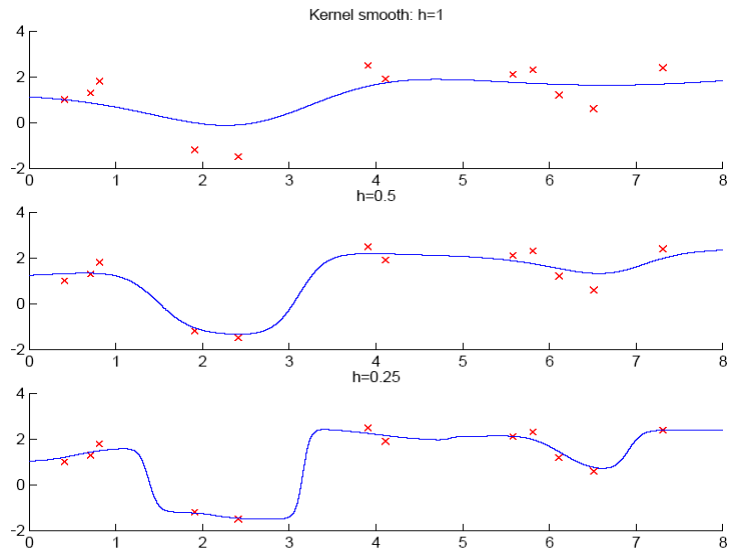
where  $K(\cdot)$  is Gaussian.

## Comparison: Mean vs Kernel Smoother





## Comparison: Mean vs Kernel Smoother



## $k$ -Nearest Neighbor Estimates

- $k$ -NN methods are more commonly used for regression than for density estimation. The classic  $k$ -NN smoother is defined as

$$\hat{m}_k(x_0) = k^{-1} \sum_{i=1}^k I(\|x_0 - x_i\| \leq d_k(x_0)) Y_i$$

This is the average value of  $y_i$  among the observations which are the  $k$  nearest neighbors of  $x_0$ . ( $d_k$  is the distance between  $x$  and  $x_0$ .)

- A smooth  $k$ -NN estimator is:

$$\hat{m}_h(x_0) = \frac{\sum_{i=1}^N w_k(\|x_0 - x_i\| \leq d_k) y_i}{\sum_{i=1}^N w_k(\|x_0 - x_i\| \leq d_k)} = \sum_{i=1}^N W_{N,k,i}(x_0) y_i$$

a weighted average of the  $k$  nearest neighbors.

### ***k*-Nearest Neighbor Estimates**

• Example:

Suppose we have the  $\{X,Y\} = \{(1,5),(7,12),(3,1),(4,0),(5,4)\}$ . Set  $k=3$ . We want to calculate  $\hat{m}(x=4)$  for the classic  $k$ -NN estimator, using Euclidian distance. Then,  $Neighborhood_{x=4} = \{3,4,5\}$ .

The weights are

$$\{W_{k=3,i}(x=4)\} = \{0, 0, 1/3, 1/3, 1/3\}$$

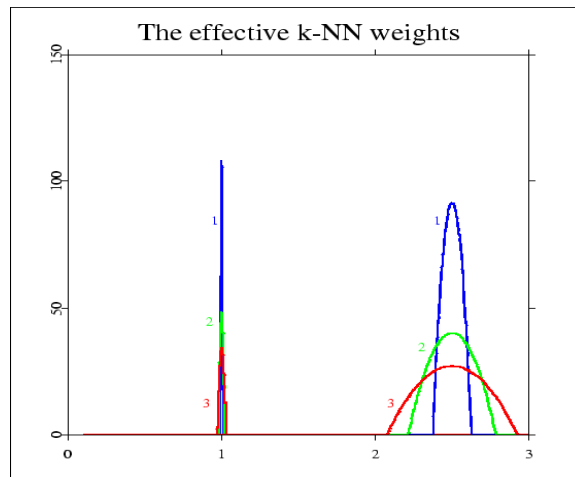
$$\hat{m}_k(x=4) = k^{-1} \sum_{i=1}^N I(\|4 - x_i\| \leq d_{k=3}(4)) Y_i = (1+0+4)/3 = 5/3$$

Note: If the X-variable is chosen from an equidistant grid, the  $k$ -NN weight are equivalent to kernel weights.

- If Epanechnikov weights are used, when observations get thin, the  $k$ -NN weights spread out more. See the food/income example, when  $x=2.5$ . (Very different weights from previous (fixed) case.)

### ***k*-Nearest Neighbor Estimates**

Figure 4. The effective  $k$ -NN weights for the food versus net income data set. At  $x=1$  and  $x=2.5$  for  $k = 100$  (label 1),  $k = 200$  (label 2),  $k = 300$  (label 3) with Epanechnikov kernel. From Hardle (1990).



## ***k*-Nearest Neighbor Estimates**

- The smoothing parameter  $k$  regulates the degree of smoothness of the estimated curve. It plays a role similar to  $h$  for kernel smoothers.
- The influence of varying  $k$  on qualitative features of the estimated curve is similar to that observed for kernel estimation with a uniform kernel.
- When  $k > N$ , the  $k$ -NN smoother is equal to the average of the response variables. When  $k = 1$ , the observations are reproduced at  $X_p$ , and for an  $x$  between two adjacent predictor variables a step function is obtained with a jump in the middle between the two observations.
- When  $\mathbf{X}$  is a vector, scaling matters. Then, always scale  $\mathbf{X}$ .

## ***k*-Nearest Neighbor Estimates**

- For the one regressor case, we have similar asymptotic results as in the univariate density case.
- Let  $N \rightarrow \infty$ ,  $k \rightarrow 0$ , and  $Nk \rightarrow \infty$ . Bias and variance of the  $k$ -NN estimate with *uniform* weights are given by

$$E[\hat{m}_k(x) - m(x)] \approx \frac{1}{24 f(x)^3} [(m'' f + 2m' f')(x)] (k/n)^2$$

$$\text{var}\{\hat{m}_k(x)\} \approx \sigma^2(x) / k$$

Note: The optimal trade-off between bias<sup>2</sup> and variance is thus achieved in an asymptotic sense by setting  $k \sim N^{4/(4+q)}$ , ( $q$ =dimension of  $\mathbf{X}$ ).  $\Rightarrow$  when  $q=1$ ,  $k \sim N^{4/5}$ .

- If  $k=2Nb f(x)$  we have exactly the same MSE at  $x$  for both kernel and  $k$ -NN estimators.

### ***k*-Nearest Neighbor Estimates**

- For the multivariate case, the asymptotic analysis is the same as for density estimation.
- Conditional on  $d_k(x)$ ; the bias and variance are approximately as for NW regression. The conditional bias is proportional to  $d_k(x)$  and the variance to  $1/[N d_k(x)^q]$  ( $q$ =dimension of vector  $X$ ).
- The optimal  $k \sim N^{4/(4+q)}$  and the optimal convergence rate is the same as for NW estimation.

### ***k*-Nearest Neighbor Estimates - Computations**

- A great advantage of the  $k$ -NN smoother is computational.
- Calculations can be easily updated. The algorithm requires  $O(N)$  operations to compute the smooth at all  $x_i$ 's. Compare this to  $O(N^2h)$  calculations for the kernel estimator.
- Cross-validation is used to set  $k$ , using leave-one-out errors:

$$CV(k) = \frac{1}{N} \sum_{i=1}^N [y_i - \hat{m}_{-i}(x_i)]^2$$

## Nonparametric Variance Estimation

- Suppose we have the following DGP:

$$y_i = m_{\mathbf{x}}(\mathbf{x}_i) + \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i$$

$$E[\varepsilon_i | \mathbf{X}_i, \mathbf{Z}_i] = 0$$

$$\varepsilon_i^2 = \sigma^2(x_i) + \eta_i, \quad E[\eta_i | \mathbf{X}_i] = 0$$

-  $\sigma^2(x)$  is the regression function of  $\varepsilon_i^2$  on  $\mathbf{x}_i$ . We want to estimate it.

- Problem: If  $\varepsilon_i^2$  were observed  $\Rightarrow$  NW or LL regression.
- Solution: Use the nonparametric residual  $e_i$ :  $e_i = y_i - \hat{m}_i(x_i)$ .

- Then, we can use the NW estimator:

$$\hat{\sigma}(x_0) = \frac{\sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right) e_i^2}{\sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right)}$$

## Nonparametric Variance Estimation

- We have a two-step estimator. Similar situation if we use the LL estimator. The bandwidths  $h$  are not the same as for estimation of  $\hat{m}(x)$ ; although we use the same notation.

Note: the LL estimator is not guaranteed to be non-negative, while the NW (or weighted NW) estimator is always non-negative (if non-negative kernels are used).

- Fan and Yao (1998) derive the surprising result that the asymptotic distribution of this two-step estimator is identical to that of the one-step idealized estimator –i.e., using  $e_i$ .

## Series Estimation

- Series estimation is the other nonparametric regression method.
- Series methods approximate an unknown function,  $m(x)$ , with a flexible parametric function, with the number of parameters treated similarly to the bandwidth in kernel regression.
- A series approximation to  $m(x)$  takes the general form:

$$m_k(x) = m_k(x, \beta),$$

where  $m_k(x, \beta)$  is a known parametric family and  $\beta$  is a vector of  $k$  unknowns.

- A linear series approximation takes the form:

$$\hat{m}_k(x) = \sum_{j=1}^K z_{jK}(x) \beta_{jK} = z_K(x)' \beta_K$$

## Series Estimation: Splines

- A linear series approximation takes the form:

$$\hat{m}_k(x) = \sum_{j=1}^K z_{jK}(x) \beta_{jK} = z_K(x)' \beta_K$$

where  $z_{jK}(x)$  are (nonlinear) functions of  $x$ ; known as *basis functions*.

- Several candidates to use for series approximation

(1) Power series. We can use a  $p$ -th order polynomial –i.e.,  $z_{jK}(x) = \mathcal{X}$ .

It works well for low  $p$ 's. But, they tend to be unstable for large  $p$ .

(2) Trigonometric series

It produces bounded functions. It can produce wiggly, wild estimates.

(3) *Splines*. A continuous piecewise polynomial function. Splines can have any polynomial order (linear, quadratic, cubic, etc.). But, it is common to use a *cubic*. It is common to constrain the spline function to have continuous derivatives up to the order of the spline.

## Series Estimation: Splines

- There is more than one way to define a spline series expansion. All are based on the number of *knots* –the points between the segments.

Examples: A piecewise linear function, with 2 segments and a *knot* at  $t$ :

$$m_K(x) = \begin{cases} m_1(x) = \beta_{00} + \beta_{01}(x - t) & x < t \\ m_2(x) = \beta_{10} + \beta_{11}(x - t) & x \geq t \end{cases}$$

The function  $m_K(x)$  is continuous if  $\beta_{00} = \beta_{10}$ . Enforcing this (and transforming the coefficients), we get:

$$\hat{m}_k(x) = \beta_0 + \beta_1 x + \beta_2 (x - t) I[x \geq t]$$

Note: This function has  $K=3$  coefficients --as a quadratic polynomial.

- Following the above process, a piecewise quadratic function, with one knot and a continuous 1st derivative has  $K=4$ .

## Series Estimation: Splines

- Similarly, a piecewise cubic function, with one knot and a continuous 2nd derivative has  $K=5$ . The function  $m_K(x)$  is

$$\hat{m}_k(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 (x - t)^3 I[x \geq t]$$

Note: The polynomial order  $p$  is selected to control the smoothness of the spline, as  $m_K(x)$  has continuous derivatives up to  $p-1$ .

- The approximation improves as the number of knots increases. For example, for a cubic spline with two knots  $t_1$  &  $t_2$  ( $t_1 < t_2$ ). The form is:

$$\hat{m}_k(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 (x - t_1)^3 I[x \geq t_1] + \beta_5 (x - t_2)^3 I[x \geq t_2]$$

- Then, a  $p$ -th-order spline with  $N$  knots at  $t_1, t_2, \dots, t_N$  ( $t_1 < t_2 < \dots < t_N$ ) is

$$\hat{m}_k(x) = \sum_{j=1}^p \beta_j x^j + \sum_{k=1}^N \gamma_k (x - t_k)^p I[x \geq t_k]$$

which has  $K = N + p + 1$  coefficients.

## Series Estimation: Splines

- In spline approximations, the usual approach is to treat  $p$  as fixed, and select  $N$  to determine the complexity of the approximation.
- The  $t_k$ 's are typically treated as fixed. It is common to set evenly spaced  $t_k$ 's. When this happens, the knot sequence is called *uniform*.

Note: For a given set of knots, the function  $m_k(x)$  is linear in the parameters  $\Rightarrow$  LS estimation is possible!

- Another popular class of series approximation are called *B-splines*. (“B” for basis). They are basis functions which are bounded, integrable and density-shaped. They can be constructed from a variety of basic shapes, usually polynomials.

## Series Estimation: B-Splines

- Let  $X \in [0; 1]$  and use uniform knots, that is, equal subintervals, with knots  $t_j = j/J; j = 0; 1, \dots, J$ . We also need knots outside of  $[0; 1]$ . Then, let  $t_j = j/m$  for all integers  $j$ .

- An  $r$ -th order B-spline is a piecewise  $(r-1)$ -order polynomial. For example, a quadratic ( $r=3$ ) B-spline base function is piecewise quadratic over three subintervals:

$$B_{r=3}(x | t_j, t_{j+1}, t_{j+2}, t_{j+r=3}) = (x-t_j)I[x \geq t_j] - 3(x-t_{j+1})I[x > t_{j+1}] + 3(x-t_{j+2})I[x > t_{j+2}] - (x-t_{j+3})I[x > t_{j+3}] = \theta_{K=4} z(x)$$

- The B-spline for an  $r$ -th order:

$$B_r(x | t_j, \dots, t_{j+r}) = \sum_{s=0}^r (-1)^s \binom{r}{s} (x-t_{j+s}) I[x > t_{j+s}]$$



### Series Estimation: B-Splines

- The B-spline is a linear combination of these basis functions:

$$\hat{m}_k(x) = \sum_{j=1-r}^{J-1} \theta_j B_r(x | t_j, \dots, t_{j+r}) = \theta_K' z$$

where  $z=z(x)$  is the vector of the basic functions.

- The number of basis functions,  $K$ , equals sum of the degree of the B-spline basis functions and the number of interior knots plus one.

$$\Rightarrow \text{Dim}(\theta) = K = J + r + 1.$$

- It is not easy to choose the optimal number of knots and their locations, which is an infinite dimensional optimization problem.

### Series Estimation: Uniform Approximations

- A good series approximation  $m_K(x)$  will have the property that it gets close to the true  $m(x)$  as  $K$  increases.

- The *Stone-Weierstrass theorem*, (Weierstrass (1885), Stone (1937, 1948)) states that any continuous function can be arbitrarily uniformly well approximated by a polynomial of sufficiently high order:

$$\sup_{x \in \mathcal{X}} |m_K(x) - m(x)| \leq \varepsilon$$

for any  $\varepsilon > 0$ .

- That is,  $m(x)$  can be arbitrarily well approximated by selecting a suitable polynomial.

## Series Estimation: Uniform Approximations

- The above result can be strengthened. If the  $s$ -th derivative of  $m(x)$  is continuous, then the uniform approximation error,  $r_{K_i}$ , satisfies

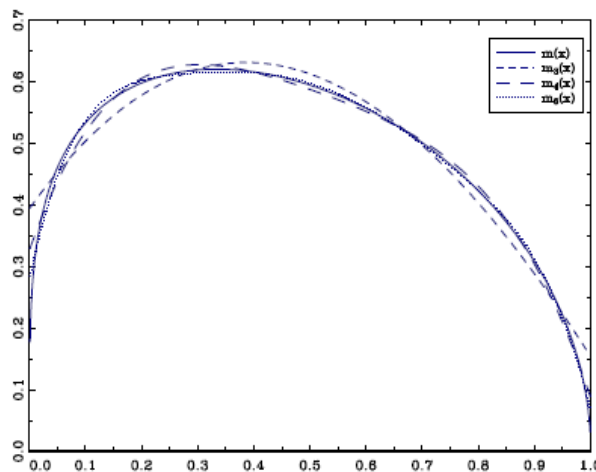
$$\sup_{x \in \mathcal{X}} |r_{K_i} = m_K(x) - m(x)| = O(K^{-\alpha})$$

as  $K \rightarrow \infty$  where  $\alpha = s/d$ . ( $\dim(\mathbf{X}) = N \times d$ )

- Useful result: It gives a rate at which the approximation  $m_K(x)$  approaches  $m(x)$  as  $K$  increases.
- Intuitively, the number of derivatives  $s$  indexes the smoothness of  $m(x)$ . The best rate at which a polynomial (or spline) approximates  $m(x)$  depends on the underlying smoothness of  $m(x)$ .
- Both results hold for spline approximations.

## Series Estimation: Uniform Approximations

- $m(x)$  can be arbitrarily well approximated by selecting a suitable polynomial. We plot approximations of  $m(x) = x^{1/4}(1-x)^{1/2}$  on  $[0,1]$ .



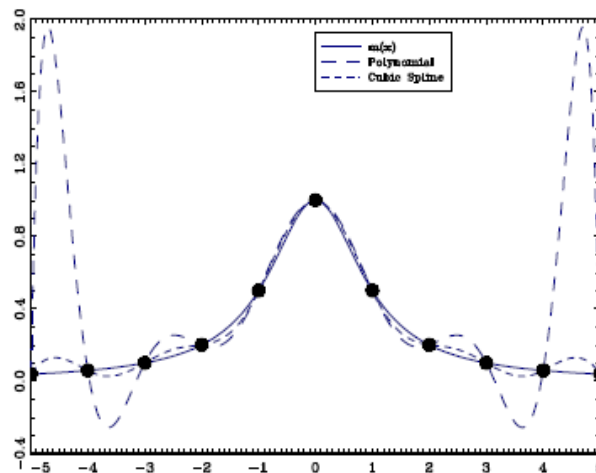
Note: The approximation with  $K = 3$  is fairly crude, but improves with  $K = 4$  and it is very good with  $K = 6$ .

## Series Estimation: Runge's Phenomenon

- Despite the excellent approximation implied by the Stone-Weierstrass theorem, polynomials have the troubling disadvantage that they are very poor at simple interpolation.
- The problem is known as *Runge's phenomenon*.
- In contrast, splines do not show Runge's phenomenon. (See next Figure.) While the fitted spline displays some oscillation relative to  $m(x)$ , but they are relatively small.
- Because of Runge's phenomenon, high-order polynomials are not used for interpolation, and are not popular choices for high-order series approximations. Instead, splines are widely used.

## Series Estimation: Runge's Phenomenon

- We plot approximations of  $m(x) = (1+x^2)^{-1}$  on  $[-5,5]$ , with  $K=11$ . Using a 10-th order polynomial. The discrepancy increases to infinity with  $K$ .



Note: The approximation is not accurate and far from the smoother true  $m(x)$ .

### Series Estimation: Regression

- We have observations on  $(Y, X)$ . Steps:
  - (1) For each  $i$ , construct the regressor vector  $z_{ki} = z_k(x_i)$ , using the series transformations.
  - (2) Stack the observations in the matrices  $\mathbf{y}$  and  $\mathbf{Z}_K$ .
  - (3) Do OLS  $\Rightarrow \mathbf{b} = (\mathbf{Z}_K' \mathbf{Z}_K)^{-1} \mathbf{Z}_K' \mathbf{y}$
  - (4) Compute the LS regression function:  $\hat{m}_k(x) = z_k(x)' \mathbf{b}_k$
  - (5) Compute estimated errors  $e_{ki} = y_{ki} - \hat{m}_k(x_i) = y_{ki} - z_k(x_i)' \mathbf{b}_k$

Note: We estimate one error,  $\varepsilon_{ki}$ , but we have two errors: the usual model error,  $\varepsilon_i$ , and the approximation error,  $r_k(x_i) = r_{ki}$ . That is,

$$\varepsilon_{ki} = r_{ki} + \varepsilon_i$$

- To assess the fit of the regression, we can calculate the  $R^2$  as usual.

### Series Estimation: Regression - $K$

- $\beta_k$  is a function of  $K$ . This reflects the goal to be flexible to incorporate greater complexity when the data are sufficiently informative. That is,  $K$  will typically be increasing with sample size  $N$ .
- $K$  plays the role of  $h$  in kernel estimation. Larger  $K$  implies smaller approximation error but increased estimation variance.
- The number of series terms,  $K$ , can be determined through CV.
- Under certain assumptions (compact set, smoothness of  $m(x)$ , bounded error variance, non-singularities in  $z_K$ , bounded  $E[z_{ki}' z_{ki}]$ ,  $K$  is chosen appropriately – i.e., a function of  $N$  and grows slower than  $N$ , etc.), the LS estimator  $\mathbf{b}_k$  converges to  $\beta_k$  in mean squared distance.

### Series Estimation: Regression - Asymptotics

- Convergence

Under certain assumptions (compact set, smoothness of  $m(x)$ , bounded error variance, non singularities in  $z_K$ , bounded  $E[z_{Ki}'z_{Ki}]$ ,  $K$  is chosen appropriately –i.e, a function of  $N$  and grows slower than  $N$ , etc.), the LS estimator  $\mathbf{b}_k$  converges to  $\beta_k$  in *m.s.* distance. See Newey (1997).

- Asymptotic normality

Even though we are in a situation similar to parametric estimation, the fact that  $K$  can grow and the finite sample bias due to the approximation error, a new theory needs to be developed.

It turns out that under the same assumptions needed for convergence and imposing some mild restrictions on  $K$  and the bias, the estimator is asymptotically normal. See Newey (1997).

### Series Estimation: Regression - Asymptotics

- The estimator has the asymptotic bias component  $r_K(x)$ , due to the finite order series as approximation to the unknown  $m(x)$ . The asymptotic distribution shows that the bias term is negligible if  $K$  diverges fast enough so that  $NK^{-2\alpha} \rightarrow 0$ . (In practical terms, this means that  $K$  is larger than optimal.)

- Asymptotic standard errors for the  $m(x)$  can be estimated with:

$$\hat{s}(x) = \sqrt{\frac{1}{n} z_K(x)' \hat{Q}_K^{-1} \hat{\Omega}_K \hat{Q}_K^{-1} z_K(x)}.$$

where

$$\hat{\Omega}_K = \frac{1}{n} \sum_{i=1}^n z_{Ki} z_{Ki}' \hat{\epsilon}_{iK}^2$$

$$\hat{Q}_K = \frac{1}{n} \sum_{i=1}^n z_{Ki} z_{Ki}'.$$

- See Newey (1997) for details.

## Spline Smoothing

- Determination of  $K$  is not easy. A perfect fit can be achieved by giving a lot of local flexibility to  $\hat{m}(x)$ . The result of this flexibility will be a jerky, difficult to interpret  $\hat{m}(x)$ .
- Spline smoothing quantifies the competition between two goals:
  - producing a good fit to the data –traditionally measured as SSR
  - producing a good curve –i.e., without too much rapid local variation.

- The regression curve  $\hat{m}_\lambda(x)$  is obtained by minimizing the penalized sum of squares

$$S_\lambda(m) = \sum_{i=1}^n \{Y_i - m(X_i)\}^2 + \lambda \int_a^b \{m''(x)\}^2 dx$$

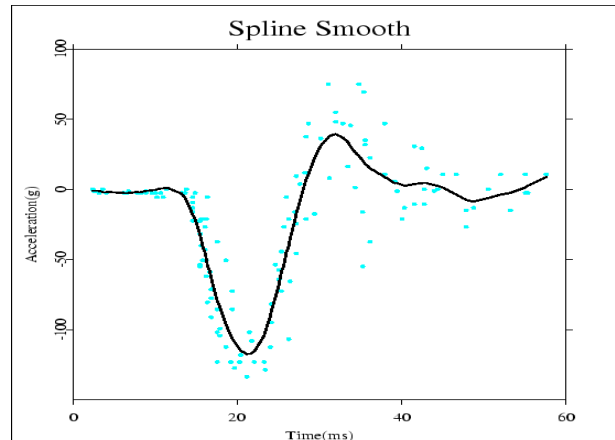
where  $m$  is twice-differentiable function on  $[a,b]$ , and  $\lambda$  represents the rate of exchange between residual error and roughness of the curve  $m$ .

## Spline Smoothing

- The second term,  $\int_a^b \{m''(x)\}^2 dx$ , represents a *roughness penalty*.
- The minimization problem over the class of all twice differentiable functions on  $[a,b]$  has a unique solution  $\hat{m}_\lambda(x)$ , which is defined as the *cubic spline*.
- $\hat{m}_\lambda(x)$  is a cubic polynomial between two successive  $X$ -values.
- At the  $x_i$ ,  $\hat{m}_\lambda(x)$  and its first two derivatives are continuous. At the boundary points  $x_{(1)}$  and  $x_{(N)}$ , the second derivative is zero.
- This properties follow from the choice of penalty for roughness. A different penalty produces different solutions.

## Spline Smoothing: Example

Figure 5. A spline smooth (Motorcycle data set). From Hardle (1990).



## Spline Smoothing

- Q: What is the spline doing to the data?

It can be shown that the spline is linear in the  $y_i$  observations, and there exists weights that

$$\hat{m}_\lambda(x) = N^{-1} \sum_{i=1}^N W_{\lambda i}(x) Y_i$$

- Silverman (1984) showed for large  $N$ , small  $\lambda$ , and  $x_i$ 's not too close to the boundary,

$$W_{\lambda i}(x) \approx f(X_i)^{-1} h(X_i)^{-1} K_s \left( \frac{x - X_i}{h(X_i)} \right)$$

where the local bandwidth  $h(X_i)$  satisfies

$$h(X_i) = \lambda^{1/4} N^{-1/4} f(X_i)^{-1/4}$$

- That is, the weight function looks like a kernel.

## Spline Smoothing: Weight Function - Example

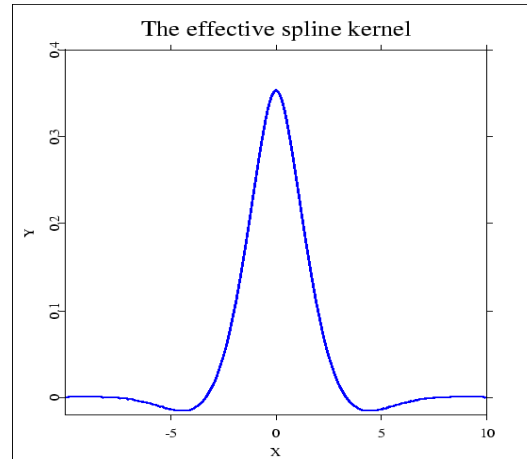


Figure 6. The asymptotic spline kernel function. From Hardle (1990).

$$K_s(u) = 1/2 \exp(-|u|/\sqrt{2}) \sin(|u|/\sqrt{2} + \pi/4).$$

## Spline Smoothing

- A variation to compute splines is to solve the equivalent problem

$$\min_m \int |m''(x)|^2 dx \quad \text{subject to} \quad \sum_{i=1}^n (Y_i - m(X_i))^2 \leq \Delta$$

- The parameters  $\lambda$  and  $\Delta$  have similar meanings, and are connected by the relationship

$$\lambda = -|G'(\Delta)|^{-1}$$

where

$$G(\Delta) = \int (\hat{m}''_{\Delta}(x))^2 dx$$

and  $\hat{m}_{\Delta}(x)$  solves the above problem.

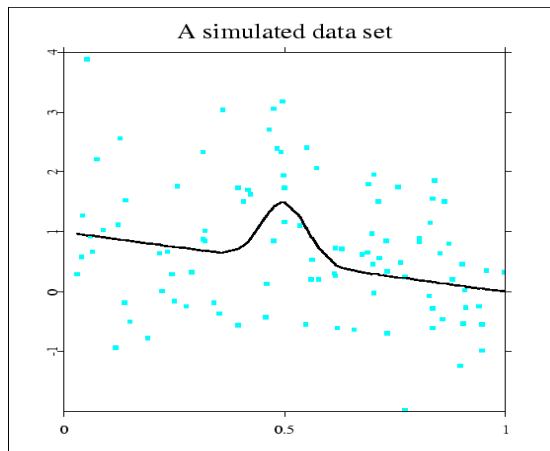


## Comparison: Kernel, $k$ -NN and Spline Smoothers

Table 1. Bias and variance of kernel and  $k$ -NN smoother

	kernel	$k$ -NN
bias	$h^2 \frac{(m''f + 2m'f')(x)}{2f(x)} d_K$	$(k/n)^2 \frac{(m''f + 2m'f')(x)}{8f^3(x)} d_K$
variance	$\frac{\sigma^2(x)}{nhf(x)} c_K$	$\frac{2\sigma^2(x)}{k} c_K$

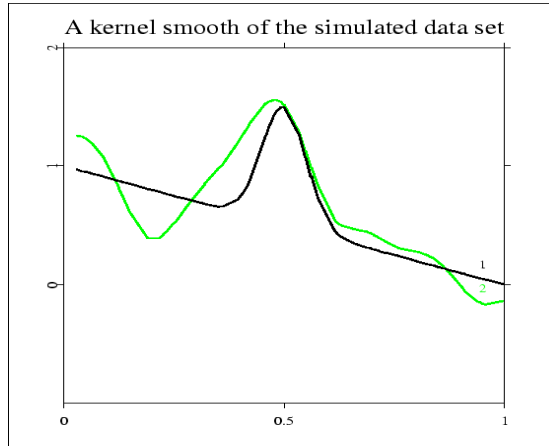
## Comparison: Kernel, $k$ -NN and Spline Smoothers



Note: Noisy Data.

Figure 7. Hardle (1990). A simulated data set. The raw data  $N=100$  were constructed from  $Y_i = m(X_i) + \varepsilon_i$ ,  $\varepsilon_i \sim N(0, 1)$ ,  $X_i \sim U(0, 1)$  and  $m(x) = 1 - x + e^{-200(x-1/2)^2}$

### Comparison: Kernel, $k$ -NN and Spline Smoothers

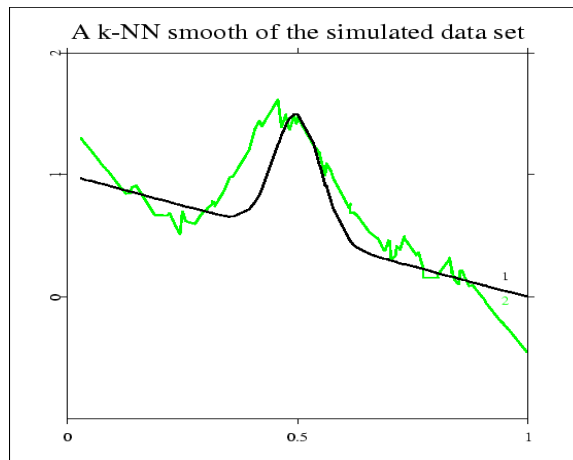


Note: As expected, kernel goes through the data.

Check smoother at boundaries (inaccurate at left).

Figure 8. A kernel smooth of the simulated data set. The black line (label 1) denotes the underlying regression curve  $m(x) = 1 - x + e^{-200(x-1/2)^2}$ . The green line (label 2) is the Gaussian kernel smooth  $\hat{m}_h(x)$ ,  $h = 0.05$ .

### Comparison: Kernel, $k$ -NN and Spline Smoothers

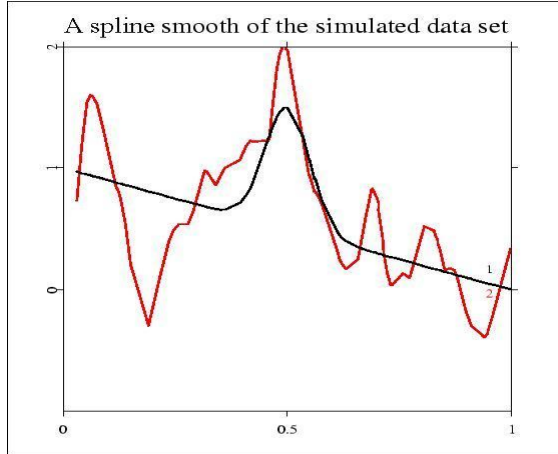


Note: Rougher curve..

Check smoother at boundaries (more points averaged).

Figure 9. Hardle (1990). A  $k$ -NN kernel smooth of the simulated data set. The black line (label 1) denotes the underlying regression curve. The green line (label 2) is the  $k$ -NN smoother.  $\hat{m}_k(x)$ ,  $k = 11$ .

### Comparison: Kernel, $k$ -NN and Spline Smoothers

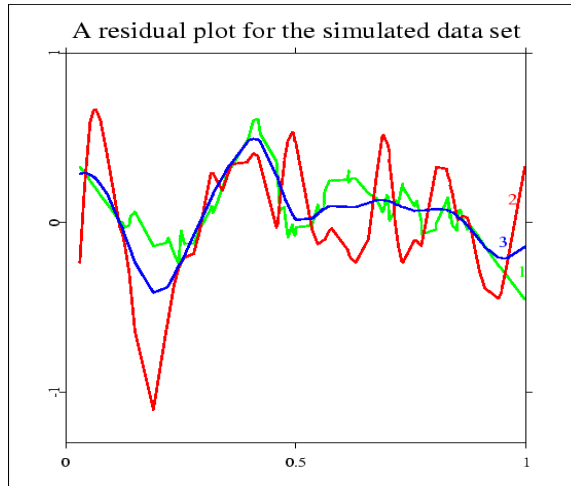


Note: As expected, very good track of observations.

Negative smooth (possible, even when all observations positive, check weights).

Figure 10. Hardle (1990). A spline smooth of the simulated data set. The black line (label 1) denotes the underlying regression curve. The green line (label 2) is the spline smoother  $\hat{m}_\Delta(x)$ ,  $\Delta = 75$ .

### Comparison: kernel, $k$ -NN and Spline smoothers



Note: Similar overall pattern. Artificial bump at  $x \approx 0.2$ .

Figure 11. Hardle (1990). Residual plot of  $k$ -NN, kernel and spline smoother for the simulated data set.

## Semiparametric Methods (SPM)

- A model is called *semiparametric* if it is described by  $\theta$  and  $\tau$ , where  $\theta$  is finite-dimensional (parametric) and  $\tau$  is infinite-dimensional (nonparametric).
- All moment condition models are semiparametric in the sense that the distribution of the data ( $\tau$ ) is unspecified and infinite dimensional. But the settings more typically called semiparametric are those where there is explicit estimation of  $\tau$ .
- In many contexts the nonparametric part  $\tau$  is a conditional mean, variance, density or distribution function.
- Often  $\theta$  is the parameter of interest, and  $\tau$  is a nuisance parameter, but this is not necessarily the case.

## Semiparametric Methods – Example 1

Example: Feasible Nonparametric GLS

$$\text{DGP: } \mathbf{y} = \mathbf{X} \theta + \boldsymbol{\varepsilon} \quad (\dim(\mathbf{X}) = N \times q)$$

$$E[\boldsymbol{\varepsilon} | \mathbf{X}] = 0$$

$$E[\boldsymbol{\varepsilon}_i^2 | \mathbf{X}] = \sigma^2(X_i) \quad (\tau(X_i) = \sigma^2(X_i))$$

where the variance function  $\sigma^2(X_i)$  is unknown but smooth in  $X$ .

We want to estimate  $\theta$ . GLS is the efficient method, but it is not feasible.

Feasible GLS is possible. Replace  $\sigma^2(X_i)$  using a nonparametric estimator (a kernel or a  $k$ -NN estimator).

Q: What is the asymptotic distribution of the GLS estimator?

## Semiparametric Methods – Example 2

Example: Generated Regressors

DGP:  $y_i = \theta \tau(\mathbf{x}_i) + \varepsilon_i$

$E[\varepsilon | \mathbf{X}] = 0$

$\theta$  is finite dimensional and  $\tau$  is an unknown function.

- Suppose  $\tau$  is identified by another equation. We have consistent estimate,  $\tau^\wedge(\mathbf{x})$ . (Imagine a non-parametric Heckman estimator).
- Then, OLS is possible to estimate  $\theta$ . This problem is called *generated regressors*, as the regressor is a (consistent) estimate of an infeasible regressor.
- Q: In general, the OLS estimator is consistent. But what is its distribution?

## SPM – Asymptotic Distribution

- Based on Andrew's (1994) MINPIN paper.

Setting:  $\theta^\wedge$  MINimizes a criterion function,  $Q_N(\theta, \tau^\wedge)$ , which depends on a Preliminary Infinite dimensional Nuisance parameter estimator.

$\Rightarrow \theta^\wedge$  is a two-step estimator

- The usual derivation of asymptotic distributions expands the f.o.c.  $m(\theta, \tau) = 0$ , We can do this for  $\theta$ , but not for  $\tau$  (it is infinite dimensional).
- To proceed, Andrews uses a stochastic equicontinuity assumption. Now, we work with the population version of  $m(\theta, \tau) = E[m_i(\theta, \tau)]$  and study the convergence of

$$\begin{aligned} \nu_n(\tau) &= \sqrt{n}(\bar{m}_n(\theta_0, \tau) - m(\theta_0, \tau)) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (m_i(\theta_0, \tau) - E m_i(\theta_0, \tau)) \end{aligned}$$

## SPM – Asymptotic Distribution

- Under a lot of assumptions:  $\hat{\theta}$  and  $\hat{\tau}(\mathbf{x}) \rightarrow^p$  to  $\theta_0$  and  $\tau_0$ ; f.o.c. equal to 0 at  $(\theta_0, \tau_0)$  –i.e., identification condition-, convergence of f.o.c.; smoothness of underlying functions; and existence of moments),

$$\sqrt{N}(\hat{\theta} - \theta_0) \rightarrow^d N(0, V),$$

where

$$\begin{aligned} V &= M^{-1} \Omega M^{-1'} \\ M &= E \frac{\partial}{\partial \theta'} m_i(\theta_0, \tau_0) \\ \Omega &= E m_i(\theta_0, \tau_0) m_i(\theta_0, \tau_0)' \end{aligned}$$

- The theorem says that  $\hat{\theta}$  has the same asymptotic distribution as the idealized estimator obtained by replacing the nonparametric estimate  $\hat{\tau}$  with the true function  $\tau_0$ .

=> the estimator is adaptive.

## SPM – Asymptotic Distribution

- But the assumptions are not trivial. The convergence in probability assumptions need to be verified. The key assumption is

$$m((\theta_0, \tau_0) = \delta Q_N(\theta, \tau) / \delta \theta | (\theta = \theta_0, \tau = \tau_0) = 0.$$

- This assumption does not always hold. It turns out, it requires a sort of orthogonality condition between the estimation of  $\theta$  and  $\tau$ .
- It holds for example 1 (FGLS with nonparametric variance), but not for Example 2 (generated regressors).

## SPM – Partially Linear Regression Model

- It is easy to define a “partially linear” regression model:

$$y_i = m_{\tilde{x}}(\mathbf{z}_i) + \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad (\dim(\mathbf{Z})=N \times q)$$

$$E[\varepsilon_i | \mathbf{X}_i, \mathbf{Z}_i] = 0$$

$$E[\varepsilon_i^2 | \mathbf{X}_i = \mathbf{x}, \mathbf{Z}_i = \mathbf{z}] = \sigma^2(x, \tilde{z})$$

- The regressors are  $(\mathbf{X}; \mathbf{Z})$ .
- The conditional mean is linear in  $X_i$ , but possibly non-linear in  $Z_i$ .
- Dummy variables are usually put in the  $X$  vector
- To keep things simple, we assume just one nonlinear variable:  $q=1$ .
- Goal: Estimate  $\boldsymbol{\beta}$  and  $m_{\tilde{x}}(\cdot)$ ; and to obtain C.I.
- Issues: Identification, Distribution of estimates.

## SPM – Estimation

- Robinson (*Econometrica*, 1988) shows we can concentrate out  $m_{\tilde{x}}(\mathbf{z}_i)$  by using a generalization of residual regression. Start with:

$$y_i = m_{\tilde{x}}(\mathbf{z}_i) + \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad (\dim(\mathbf{Z})=N \times q)$$

Taking conditional expectations in  $\mathbf{Z}$ :

$$E[y_i | \mathbf{z}_i] = E[m_{\tilde{x}}(\mathbf{z}_i) | \mathbf{z}_i] + E[\mathbf{x}_i' \boldsymbol{\beta} | \mathbf{z}_i] = m_{\tilde{x}}(\mathbf{z}_i) + E[\mathbf{x}_i' | \mathbf{z}_i] \boldsymbol{\beta}$$

- Two conditional means:

$$- m_y(\mathbf{z}_i) = E[y_i | \mathbf{z}_i]$$

$$- m_x(\mathbf{z}_i) = E[\mathbf{x}_i' | \mathbf{z}_i]$$

- Then,

$$m_y(\mathbf{z}_i) = m_{\tilde{x}}(\mathbf{z}_i) + m_x(\mathbf{z}_i)' \boldsymbol{\beta}$$

Subtract from the original equation ( $m_{\tilde{x}}(\mathbf{z}_i)$  disappears):

$$y_i - m_y(\mathbf{z}_i) = [\mathbf{x}_i' - m_x(\mathbf{z}_i)'] \boldsymbol{\beta} + \varepsilon_i$$

## SPM – Estimation

- Rewrite in terms of residuals:  $y_i - m_y(\mathbf{z}_i) = [\mathbf{x}_i' - m_x(\mathbf{z}_i)'] \boldsymbol{\beta} + \varepsilon_i$ 
  - $\varepsilon_{yi} = y_i - m_y(\mathbf{z}_i)$
  - $\varepsilon_{xi} = [\mathbf{x}_i' - m_x(\mathbf{z}_i)']$
  - $\varepsilon_{yi} = \varepsilon_{xi}' \boldsymbol{\beta} + \varepsilon_i$
- That is,  $\boldsymbol{\beta}$  is the coefficient of the regression of  $\varepsilon_{yi}$  on  $\varepsilon_{xi}$ . But, we do not observe the errors. It is an unfeasible LS estimator!
- Robinson suggests the following steps:
  - 1) Estimating  $m_y(\mathbf{z}_i)$  and  $m_x(\mathbf{z}_i)$  by NW regression (different  $b$ 's, OK).
  - 2) Get the residuals,  $\varepsilon_{xi}$  &  $\varepsilon_{yi}$ .
  - 3) Using the residuals, do OLS to estimate  $\boldsymbol{\beta}$ .

Note: We can use in 1) LL or weighted NW.

## SPM – Estimation: Trimming

- The nonparametric regression estimates depend inversely on  $f_z^{\wedge}(\tilde{z})$ .
- Problem: For values of  $\tilde{z}$  where  $f_z(\tilde{z})$  is close to 0,  $f_z^{\wedge}(\tilde{z})$  is not bounded away from 0. The NW estimates at this point can be poor.
- Solution: Trimming.  
Let  $b > 0$  be a trimming constant. The trimmed estimator of  $\boldsymbol{\beta}$  is:
 
$$\boldsymbol{\beta}^{\wedge} = (\sum_i \varepsilon_{xi} \varepsilon_{xi}' I[f_z^{\wedge}(\tilde{z}_i) \geq 0])^{-1} \sum_i \varepsilon_{xi} \varepsilon_{yi} I[f_z^{\wedge}(\tilde{z}_i) \geq 0]$$
 => This is a trimmed LS residual regression.
- The asymptotic theory requires that  $b = b_N \rightarrow 0$ , but it is not clear how to select  $b$  in practice. Often trimming is ignored in applications.  
Suggestion: Estimate model with and without trimming.



## SPM – Asymptotic Distribution

- The needed regularity conditions: the data are *i.i.d.*,  $Z_i$  has a density, and the regression functions, density, and conditional variance function are sufficiently smooth with respect to their arguments.

- Assume  $b$  is the same for all  $q$ . The important condition on the  $b$  sequence is

$$\sqrt{n} \left( h^4 + \frac{1}{nh^q} \right) \rightarrow 0$$

- Equivalently, what is essential is that the estimators themselves converge faster than  $N^{-1/4}$ . From the theory for nonparametric regression, these rates hold when  $b$ 's are picked optimally and  $q \leq 3$ .

## SPM – Asymptotic Distribution

- **Theorem** (Robinson). Under regularity conditions, including  $q \leq 3$ ; the trimmed estimator satisfies

$$\sqrt{n} (\hat{\beta} - \beta) \rightarrow_d N(0, V)$$

$$V = (E(e_{xi}e'_{xi}))^{-1} (E(e_{xi}e'_{xi}\sigma^2(X_i, Z_i)))^{-1} (E(e_{xi}e'_{xi}))^{-1}$$

- That is,  $\hat{\beta}$  is asymptotically equivalent to the infeasible LS estimator.
- Estimate the variance matrix  $V$  as usual, using residuals.

## SPM – Estimation of Nonparametric Part

- The model:

$$y_i = m_x(z_i) + \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad (\dim(\mathbf{Z}) = N \times q)$$

- We estimated  $\boldsymbol{\beta}$ . Now, we want to estimate  $m_x(z_i)$ . It looks like an iterative algorithm is needed, but since  $\boldsymbol{\beta}$  converges faster than the nonparametric rate, we can pretend it is fixed. Then,

$$\hat{m}_z(z_0) = \frac{\sum_{i=1}^N K_h\left(\frac{z_0 - z_i}{h}\right) (y_i - \mathbf{X}_i' \hat{\boldsymbol{\beta}})}{\sum_{i=1}^N K_h\left(\frac{z_0 - z_i}{h}\right)}$$

- The bandwidth  $h = (h_1, \dots, h_q)$  is distinct from those for the first-stage regressions. Standard errors for  $\hat{m}_x(z_i)$  as usual for standard nonparametric regression.

## SPM – Bandwidth Choice

- In a semiparametric context, it is important to study the effect a bandwidth has on the performance of the estimator of interest before determining the bandwidth.
- In many cases, this requires a nonconventional bandwidth rate.
- However, this problem does not occur in partially linear models. The first-step bandwidths  $h$  used for  $\hat{m}_y(z_i)$  and  $\hat{m}_x(z_i)$  are inputs for calculation of  $\hat{\boldsymbol{\beta}}$ .
- $h$  impacts the theory for  $\hat{\boldsymbol{\beta}}$ , through the uniform convergence rates for  $\hat{m}_y(z_i)$  and  $\hat{m}_x(z_i)$ , suggesting that we use conventional bandwidth rules, for example CV.

## Further Comments

- There are some specification tests that compare non-parametric regressions (“unconstrained” model) with parametric regressions (“constrained” model). See Blundell and Duncan (1998), Pagan and Ullah (1999) and Yatchew (Chapter 6).
- Recent research has focused on correcting for *endogeneity* (see Yatchew) and *heteroscedasticity* (see Yatchew). In general, the most promising approaches are two-step methods.
  - (1) Non-parametrically regress endogenous  $x$  variables on the IV  $z$  and calculate “errors” as the difference between those  $x$  variables and their (non-parametrically) predicted values.
  - (2) Add these errors into the equation of interest.

## Readings

- Blundell and Duncan (1998), “Kernel Regression in Empirical Microeconomics,” JEL.
- Blundell and Powell (2003) “Endogeneity in Nonparametric and Semiparametric Regression Models” in **Advances in Economics and Econometrics**, edited by Dewatripont, Hansen and Turnovsky.
- Cameron, A. and P. Trivedi (2003), **Microeconometrics: Methods and Applications**, Cambridge University Press.
- Hansen, B. (2013), **Econometrics**.
- Ichimura and Todd (2007) “Implementing Nonparametric and Semi-Parametric Estimators”, in *Handbook of Econometrics, Volume 6B*
- Pagan, A and A. Ullah (1999), **Nonparametric Econometrics**, Cambridge University Press.
- Yatchew, A (2003), **Semiparametric Regression for the Applied Econometrician**, Cambridge University Press.