

NONPARAMETRIC REGRESSION

Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be indep. realizations of $(X, Y) \in \mathbb{R} \times \mathbb{R}$, where

$$Y = m(X) + \varepsilon,$$

with ε independent of X with $\mathbb{E} \varepsilon = 0$, $\mathbb{E} \varepsilon^2 = \sigma^2$.

So we have $m(x) = \mathbb{E}[Y|X=x]$.

Linear regression assumes

$$m \in \{m: \mathbb{R} \rightarrow \mathbb{R}: m(x) = \beta_0 + \beta_1 x, \beta_0, \beta_1 \in \mathbb{R}\}$$

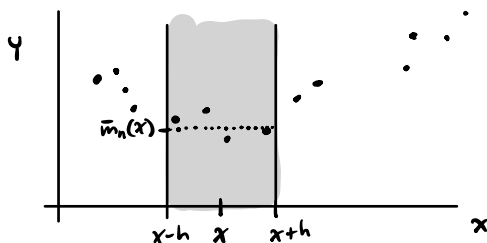
We will not assume any parametric form. Rather, we assume that m belongs to some class of functions of a certain smoothness.

There are a multitude of ways to estimate m nonparametrically!

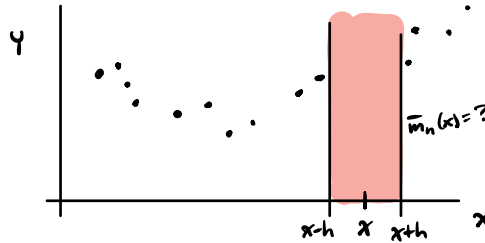
We first consider a local-averaging estimator

$$\bar{m}_n(x) = \frac{\sum_{i=1}^n Y_i \mathbb{1}(x-h \leq X_i \leq x+h)}{\sum_{j=1}^n \mathbb{1}(x-h \leq X_j \leq x+h)}.$$

At a point $x \in [0, 1]$, $\bar{m}_n(x)$ is the average of Y_i values for which the corresponding X_i values are near x . It is the average of points within a moving window:



Note that we need the denominator $\sum_{i=1}^n \mathbb{1}(x-h \leq X_i \leq x+h) = \#\{X_i \in (x \pm h)\}$ to be positive for all x . Otherwise the estimator will be undefined:



A more general version of the local averaging estimator is the Nadaraya-Watson estimator

$$\hat{m}_n^{NW}(x) = \frac{\sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h}\right)}{\sum_{j=1}^n K\left(\frac{X_j - x}{h}\right)},$$

where K is a kernel function like

$$K(u) = \mathbb{1}(|u| \leq 1)$$

$$K(u) = (1 - u^2) \mathbb{1}(|u| \leq 1)$$

$$K(u) = \frac{1}{\sqrt{2\pi}} \exp(-u^2/2)$$

Note that we do not need $\int_{\mathbb{R}} K(u) du = 1$ (as with KDEs).

Mean squared error of Nadaraya-Watson estimator

We consider $MSE \hat{m}_n^{NW}(x_0)$, $x_0 \in [0, 1]$, when $m \in \text{Lipschitz}(L)$ on $[0, 1]$.
(can shift/scale X_1, \dots, X_n to be on $[0, 1]$, so no generality is lost)

Let $MSE \hat{m}_n^{NW}(x_0) = \underbrace{b_n^2(x_0)}_{\text{squared bias}} + \underbrace{\sigma_n^2(x_0)}_{\text{variance}}$, with

$$b_n(x_0) = \mathbb{E} \hat{m}_n^{NW}(x_0) - m(x_0), \quad \sigma_n^2(x_0) = \text{Var} \hat{m}_n^{NW}(x_0). \quad \boxed{2}$$

We will make the following assumptions on the kernel and on X_1, \dots, X_n :

(K1) Let $K: \mathbb{R} \rightarrow \mathbb{R}$ have support on $[-1, 1]$ and satisfy $0 \leq K(u) \leq K_{\max} < \infty \forall u \in \mathbb{R}$.

(D1) Let $X_1, \dots, X_n \in [0, 1]$ be deterministic such that for some $n_0 > 0$

$$0 < f_{\min} \leq \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right) \leq f_{\max} < \infty \quad \text{for all } x \in [0, 1]$$

for all $n \geq n_0$.

Note: (K1) excludes the beloved Gaussian kernel, but it makes our proof simpler!

Result: Under (K1) and (D1), if $m \in \text{Lipschitz}(L)$ on $[0, 1]$, we have

$$\text{MSE } \hat{m}_n^{\text{NW}}(x_0) \leq h^2 \cdot L^2 + \frac{\sigma^2}{nh} \frac{K_{\max}}{f_{\min}} \quad \text{for all } x_0 \in [0, 1],$$

provided $n \geq n_0$.

Proof: First write

$$\hat{m}_n^{\text{NW}}(x) = \frac{\sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h}\right)}{\sum_{j=1}^n K\left(\frac{X_j - x}{h}\right)} = \sum_{i=1}^n W_{ni}(x) Y_i,$$

with

$$W_{ni}(x) = \frac{K\left(\frac{X_i - x}{h}\right)}{\sum_{j=1}^n K\left(\frac{X_j - x}{h}\right)}.$$

Now we consider the variance:

$$\hat{\sigma}_n^2(x_0) = \text{Var} \left[\sum_{i=1}^n W_{ni}(x_0) Y_i \right]$$

$$= \sum_{i=1}^n W_{ni}^2(x_0) \sigma^2$$

$$\leq \sum_{i=1}^n W_{ni}(x_0) \cdot W_{ni}(x_0) \cdot \sigma^2 \quad (K(u) \geq 0 \forall u) \quad \boxed{3}$$

$$\begin{aligned}
&\leq \sup_{1 \leq j \leq n} W_{nj}(x_0) \sum_{i=1}^n W_{ni}(x_0) \sigma^2 \quad \left(\sum_{i=1}^n W_{ni}(x_0) = 1 \right) \\
&= \sup_{1 \leq j \leq n} W_{nj}(x_0) \sigma^2 \\
&\leq \sup_{1 \leq j \leq n} \left(\frac{K \left(\frac{x_j - x_0}{h} \right)}{\sum_{i=1}^n K \left(\frac{x_i - x_0}{h} \right)} \sigma^2 \right) \\
&\leq \frac{\frac{1}{nh} K_{\max}}{\frac{1}{nh} \sum_{i=1}^n K \left(\frac{x_i - x_0}{h} \right)} \cdot \sigma^2 \\
&\leq \frac{1}{nh} \cdot \frac{K_{\max}}{f_{\min}} \cdot \sigma^2.
\end{aligned}$$

Now the bias:

$$\begin{aligned}
b_n(x_0) &= \mathbb{E} \sum_{i=1}^n W_{ni}(x_0) Y_i - m(x_0) \\
&= \sum_{i=1}^n W_{ni}(x_0) m(x_i) - m(x_0) \\
&= \sum_{i=1}^n W_{ni}(x_0) [m(x_i) - m(x_0)]. \quad \left(\sum_{i=1}^n W_{ni}(x_0) = 1 \right)
\end{aligned}$$

From here we write

$$\begin{aligned}
|b_n(x_0)| &\leq \sum_{i=1}^n W_{ni}(x_0) L |x_i - x_0| \quad (m \text{ Lipschitz } (L) \text{ on } [0,1]) \\
&\stackrel{K \text{ supported on } [-1,1]}{=} \sum_{i=1}^n W_{ni}(x_0) L |x_i - x_0| \mathbb{1}(|x_i - x_0| \leq h) \\
&\leq \sum_{i=1}^n W_{ni}(x_0) L \cdot h \\
&= h \cdot L. \quad \square \quad \boxed{14}
\end{aligned}$$

Exercise: Find optimal bandwidth and then optimal MSE bound for $\hat{m}_n^{NW}(x_0)$.

$$\frac{\partial}{\partial h} (\text{bound on MSE } \hat{m}_n^{NW}(x_0)) = 2h \cdot L^2 - \frac{\sigma^2}{nh^2} \frac{K_{max}}{f_{min}} = 0$$

$$\Leftrightarrow h = n^{-\frac{1}{3}} \left(\frac{\sigma^2}{L^2} \frac{K_{max}}{2 \cdot f_{min}} \right)^{\frac{1}{3}} =: h_{opt}.$$

Under this choice of h , we have

$$\begin{aligned} \text{MSE } \hat{m}_n^{NW}(x_0) &\leq \left[n^{-\frac{1}{3}} \left(\frac{\sigma^2}{L^2} \frac{K_{max}}{2 \cdot f_{min}} \right)^{\frac{1}{3}} \right]^2 \cdot L^2 + \frac{\sigma^2}{n} \frac{K_{max}}{f_{min}} n^{\frac{1}{3}} \left(\frac{L^2 \cdot 2 \cdot f_{min}}{\sigma^2 \cdot K_{max}} \right)^{\frac{1}{3}} \\ &= n^{-\frac{2}{3}} L^{\frac{2}{3}} \left\{ \left(\frac{\sigma^2 K_{max}}{2 f_{min}} \right)^{\frac{2}{3}} + \left(\frac{\sigma^2 K_{max}}{f_{min}} \right)^{\frac{2}{3}} \right\}. \end{aligned}$$

We summarize the above by giving a uniform result:

Result: Under (K1) and (D1), with $h = \alpha n^{-\frac{1}{3}}$ for some $\alpha > 0$, we have

$$\sup_{m \in \text{Lipshitz}(L) \text{ on } [0,1]} \sup_{x \in [0,1]} \mathbb{E} \left[\hat{m}_n^{NW}(x) - m(x) \right]^2 \leq n^{-\frac{2}{3}} \cdot C,$$

for $n \geq n_0$, where $C > 0$ depends on K_{max} , L , σ^2 , and f_{min} .

Comparison to MSE in simple linear regression:

What would the MSE bound be in parametric regression?

Let $\hat{m}_n^{Lin}(x)$ be the estimator given by

$$\hat{m}_n^{Lin}(x) = \hat{\beta}_0 + \hat{\beta}_1 x, \quad (\hat{\beta}_0, \hat{\beta}_1) = \underset{\beta_0, \beta_1 \in \mathbb{R}}{\text{argmin}} \sum_{i=1}^n [Y_i - (\beta_0 + \beta_1 X_i)]^2 \quad \boxed{5}$$

and define the class of linear functions on an interval $T \subset \mathbb{R}$ as

$$\text{Lin}(T) = \{m: T \rightarrow \mathbb{R} : m(x) = ax + b, a, b \in \mathbb{R}\}.$$

We will make the following assumption:

(D1 lin) Let $X_1, \dots, X_n \in [0, 1]$ be deterministic such that for some $n_0 > 0$

$$\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \geq S_{\min} > 0$$

for all $n \geq n_0$.

Result: Under (D1 lin), for $n \geq n_0$, we have

$$\sup_{m \in \text{Lin}([0,1])} \sup_{x \in [0,1]} \mathbb{E} \left[\hat{m}_n^{\text{lin}}(x_0) - m(x_0) \right]^2 \leq n^{-1} \left[1 + \frac{1}{S_{\min}} \right] \sigma^2.$$

Proof: For every $x_0 \in [0, 1]$, $\hat{m}_n^{\text{lin}}(x_0)$ is unbiased, since

$$\mathbb{E} \hat{m}_n^{\text{lin}}(x_0) = \beta_0 + \beta_1 x_0 = m(x_0).$$

Next, we have

$$\begin{aligned} \text{Var} \hat{m}_n^{\text{lin}}(x_0) &= \text{Var} \left(\hat{\beta}_0 + \hat{\beta}_1 x_0 \right) \\ &= \left[\frac{1}{n} + \frac{(x_0 - \bar{X}_n)^2}{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \right] \sigma^2 \quad (\text{See my STAT 513 notes}) \\ &= \frac{1}{n} \left[1 + \frac{(x_0 - \bar{X}_n)^2}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2} \right] \sigma^2 \end{aligned}$$

$$\leq \frac{1}{n} \left[1 + \frac{1}{S_{\min}} \right] \sigma^2. \quad \left((D1 \text{ lin}) \text{ and } (x_0 - \bar{x}_n) \leq 1 \right)$$

Since this holds for all $m \in \text{Lin}([0,1])$, the proof is complete.

Summary: The MSE of the parametric estimator is smaller than that of the non parametric estimator. But parametric models might be misspecified (this is what motivates non parametric statistics).

Confidence Interval for $m(x_0)$ at $x_0 \in [0,1]$:

To construct a C.I. for $m(x_0)$ at some $x_0 \in [0,1]$, we would like a result like

$$\left(\hat{m}_n^{NW}(x_0) - m(x_0) \right) / \vartheta_n(x_0) \rightarrow N(0,1) \text{ in dist. as } n \rightarrow \infty$$

for some sequence $\vartheta_n(x_0)$. We will break this into two pieces.

$$\left(\hat{m}_n^{NW}(x_0) - m(x_0) \right) / \vartheta_n(x_0) = \underbrace{\left(\hat{m}_n^{NW}(x_0) - \mathbb{E} \hat{m}_n^{NW}(x_0) \right) / \vartheta_n(x_0)}_{\text{Asymptotically Normal?}} + \underbrace{\left(\mathbb{E} \hat{m}_n^{NW}(x_0) - m(x_0) \right) / \vartheta_n(x_0)}_{\text{Vanishing bias term?}}$$

Remark: Note that for consistency ($MSE \hat{m}_n^{NW}(x_0) \rightarrow 0$) we need $h \rightarrow 0$ and $nh \rightarrow \infty$.

For the Normality part, we will need the following result:

Corollary to Lindeberg CLT:

For seq. of iid rvs ξ_1, ξ_2, \dots with mean 0 and variance 1 and a seq. of real numbers a_1, a_2, \dots that satisfy

$$\max_{1 \leq i \leq n} |a_i| / \left[\sum_{j=1}^n a_j^2 \right]^{1/2} \rightarrow 0 \text{ as } n \rightarrow \infty,$$

we have

$$\frac{\sum_{i=1}^n a_i \xi_i}{\left[\sum_{j=1}^n a_j^2 \right]^{1/2}} \rightarrow N(0,1) \text{ in distribution as } n \rightarrow \infty.$$

□

We now present a theorem:

Theorem: Let K have support on $[-1, 1]$ and satisfy $0 \leq K(u) \leq K_{\max} < \infty$ for all $u \in \mathbb{R}$ and let X_1, \dots, X_n be deterministic such that for any sequence of bandwidths h_n satisfying $h_n \rightarrow 0$ and $nh_n \rightarrow \infty$ there exists an $n_0 > 0$ such that

$$\left(\text{flower} \right) \frac{1}{nh_n} \sum_{i=1}^n K^2\left(\frac{X_i - x}{h_n}\right) \geq M > 0 \quad \text{and} \quad \left(\text{needed?} \right) \frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right) \leq f_{\max} < \infty$$

for all $x \in [0, 1]$ for all $n \geq n_0$. Then for $m \in \text{Lipschitz}(L)$ on $[0, 1]$ and any sequence $h_n \rightarrow 0$ with $nh_n \rightarrow \infty$, we have

$$\textcircled{1} \quad \frac{\hat{m}_n^{\text{NW}}(x_0) - \mathbb{E} \hat{m}_n^{\text{NW}}(x_0)}{\sigma \left[\sum_{i=1}^n W_{ni}^2(x_0) \right]^{1/2}} \rightarrow N(0, 1) \text{ in dist. as } n \rightarrow \infty$$

$$\textcircled{2} \quad \frac{|\mathbb{E} \hat{m}_n^{\text{NW}}(x_0) - m(x_0)|}{\sigma \left[\sum_{i=1}^n W_{ni}^2(x_0) \right]^{1/2}} \leq \frac{n^{1/2} h_n^{3/2} f_{\max}}{\sigma M} L \quad \text{for all } n \geq n_0.$$

Proof of ①: Note that we can write

$$\begin{aligned} \hat{m}_n^{\text{NW}}(x_0) - \mathbb{E} \hat{m}_n^{\text{NW}}(x_0) &= \sum_{i=1}^n W_{ni}(x_0) \gamma_i - \sum_{i=1}^n W_{ni}(x_0) \mathbb{E} \gamma_i \\ &= \sum_{i=1}^n W_{ni}(x_0) \varepsilon_i, \end{aligned}$$

$$\text{where } W_{ni}(x_0) = K\left(\frac{X_i - x_0}{h_n}\right) / \sum_{j=1}^n K\left(\frac{X_j - x_0}{h_n}\right).$$

By the Corollary to the Lindeberg CLT, it is suff. to show that

$$\max_{1 \leq i \leq n} |W_{ni}(x_0)| / \left[\sum_{j=1}^n W_{nj}^2(x_0) \right]^{1/2} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

We have

$$\frac{\max_{1 \leq i \leq n} |W_{ni}(x_0)|}{\left[\sum_{j=1}^n W_{nj}^2(x_0) \right]^{1/2}} = \frac{\max_{1 \leq i \leq n} K\left(\frac{x_i - x_0}{h_n}\right)}{\sum_{j=1}^n K\left(\frac{x_j - x_0}{h_n}\right)} \bigg/ \left(\sum_{i=1}^n \frac{K^2\left(\frac{x_i - x_0}{h_n}\right)}{\left[\sum_{j=1}^n K\left(\frac{x_j - x_0}{h_n}\right) \right]^2} \right)^{1/2}$$

$$= \frac{\max_{1 \leq i \leq n} K\left(\frac{x_i - x_0}{h_n}\right)}{\left[\sum_{j=1}^n K^2\left(\frac{x_j - x_0}{h_n}\right) \right]^{1/2}}$$

$$\leq \frac{K_{\max}}{\sqrt{nh_n} \underbrace{\left[\frac{1}{nh_n} \sum_{j=1}^n K^2\left(\frac{x_j - x_0}{h_n}\right) \right]^{1/2}}_{\geq M \quad \forall n \geq n_0}} \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

by $(*)$ and because $nh_n \rightarrow \infty$.

Proof of (2): By $K(x) = 0$ for $x \notin [-1, 1]$ and $m \in \text{Lip}(\alpha; L)$, we have

$$\left| \mathbb{E} m_n^{ANW}(x_0) - m(x_0) \right| \leq h_n \cdot L. \quad (\text{previous work})$$

Also, for all $n \geq n_0$ we have

$$\frac{1}{\left[\sum_{i=1}^n W_{ni}^2(x_0) \right]^{1/2}} = \frac{\sum_{i=1}^n K\left(\frac{x_i - x_0}{h_n}\right)}{\left[\sum_{j=1}^n K^2\left(\frac{x_j - x_0}{h_n}\right) \right]^{1/2}} = \frac{\frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{x_i - x_0}{h_n}\right)}{\left[\frac{1}{nh_n} \sum_{j=1}^n K^2\left(\frac{x_j - x_0}{h_n}\right) \right]^{1/2}} \leq \frac{\sqrt{nh_n} f_{\max}}{M},$$

by $(*)$. Combining these bounds, we obtain (2).

Remark: The condition $(*)$ just means that the x_1, \dots, x_n are "nicely" spread across the interval $[0, 1]$, not too concentrated around any one point and not too sparse over any intervals.

Corollary: Under the conditions of the Theorem, we have

$$\frac{[\hat{m}_n^{NW}(x_0) - m(x_0)]}{\sigma \left[\sum_{i=1}^n W_{ni}^2(x_0) \right]^{1/2}} \rightarrow N(0,1) \text{ in dist. as } n \rightarrow \infty$$

provided $n^{1/2} h_n^{3/2} \rightarrow 0$.

Exercise: Determine whether we can build a valid C.I. for $m(x_0)$ under the MSE-optimal choice of h_n .

Solution: The MSE-optimal choice of h_n is $h_n = \alpha n^{-1/3}$ for some $\alpha > 0$.
Under this choice, $n^{1/2} h_n^{3/2} = n^{1/2} (\alpha^{-1/3})^{3/2} = \alpha^{1/2} = 1$, so the bias does not vanish.

Remark: We must choose h_n to be smaller than the MSE-optimal bandwidth in order to build valid C.I.s. This is called "undersmoothing".

Estimating σ^2 :

Let $(X_{(1)}, Y_{(1)}), \dots, (X_{(n)}, Y_{(n)})$ represent the data reordered such that

$$X_{(1)} < X_{(2)} < \dots < X_{(n)}.$$

Consider the variance estimator:

$$\hat{\sigma}_n^2 = \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (Y_{(i+1)} - Y_{(i)})^2$$

Result: If $X_1, \dots, X_n \in [0,1]$, $\mathbb{E} \varepsilon_i^4 = \mu_4 < \infty$, $m \in \text{Lipschitz}(L)$ on $[0,1]$, then

$$\hat{\sigma}_n^2 \rightarrow \sigma^2 \text{ in probability}$$

as $n \rightarrow \infty$.

Proof:

We have

$$\begin{aligned}\hat{\sigma}_n^2 &= \frac{1}{2(n-1)} \sum_{i=1}^{n-1} \left[(m(X_{(i+1)}) + \varepsilon_{(i+1)}) - (m(X_{(i)}) + \varepsilon_{(i)}) \right]^2 \\ &= \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (\varepsilon_{(i+1)} - \varepsilon_{(i)})^2 + \frac{1}{2(n-1)} \sum_{i=1}^{n-1} [m(X_{(i+1)}) - m(X_{(i)})]^2 \\ &\quad + 2 \cdot \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (\varepsilon_{(i+1)} - \varepsilon_{(i)}) [m(X_{(i+1)}) - m(X_{(i)})].\end{aligned}$$

The first term has expectation σ^2 and the third term has expectation 0.

The second term satisfies

$$\begin{aligned}\frac{1}{2(n-1)} \sum_{i=1}^{n-1} [m(X_{(i+1)}) - m(X_{(i)})]^2 &\leq \frac{1}{2(n-1)} \sum_{i=1}^{n-1} L^2 \cdot (X_{(i+1)} - X_{(i)})^2 \quad (m \in \text{Lipschitz}(L)) \\ &\leq \frac{L^2}{2(n-1)}. \quad (X_1, \dots, X_n \in [0, 1])\end{aligned}$$

Therefore $\mathbb{E} \hat{\sigma}_n^2 - \sigma^2 \rightarrow 0$ as $n \rightarrow \infty$.

Moreover

$$\begin{aligned}\text{Var} \hat{\sigma}_n^2 &= \frac{1}{4(n-1)^2} \left[\sum_{i=1}^{n-1} \text{Var} \left((Y_{(i+1)} - Y_{(i)})^2 \right) + \sum_{|i-j|=1} \text{Cov} \left((Y_{(i+1)} - Y_{(i)})^2, (Y_{(j+1)} - Y_{(j)})^2 \right) \right] \\ &\leq \frac{1}{4(n-1)^2} \left[\sum_{i=1}^{n-1} \text{Var} \left((Y_{(i+1)} - Y_{(i)})^2 \right) + \sum_{|i-j|=1} \sqrt{\text{Var} \left((Y_{(i+1)} - Y_{(i)})^2 \right) \text{Var} \left((Y_{(j+1)} - Y_{(j)})^2 \right)} \right] \\ &\stackrel{\text{(see Appendix)}}{\leq} \frac{1}{4(n-1)^2} \left[(n-1) + 2(n-1) \right] \left[2 \underbrace{\left(2\mu_4 + 6\sigma^4 \right)}_{\text{Upper bound for } \text{Var} \left((Y_{(i+1)} - Y_{(i)})^2 \right) \forall i} + 16\sigma^2 L^2 \right] \\ &\rightarrow 0 \quad \text{as } n \rightarrow \infty.\end{aligned}$$

Since the bias and variance go to zero, we have consistency. \square

\square

Note: Other estimators exist, but this one is nice because it does not depend on a bandwidth choice.

Corollary: Under the conditions of the theorem and $\mathbb{E} \Sigma_i^4 < \infty$,

$$P \left(m(x_0) \in \left(\hat{m}_n^{NW}(x_0) \pm z_{\alpha/2} \hat{\sigma}_n^2 \sqrt{\sum_{i=1}^n W_{hi}^2(x_0)} \right) \right) \rightarrow 1 - \alpha$$

as $n \rightarrow \infty$, provided $\frac{1}{n} \frac{3}{2} \frac{1}{h_n} \rightarrow 0$.

Local-Polynomial Estimators

We now introduce a generalization of the N-W estimator and consider its performance over $m \in \mathcal{H}(p, L)$ on $[0, 1]$.

This estimator fits a polynomial function to the data locally.

We have in mind, for x close to x_0 , the approximation

$$m(x) \approx a_0 + a_1(x-x_0) + \dots + a_p \frac{(x-x_0)^p}{p!}$$

for some coefficients a_0, a_1, \dots, a_p .

We thus define $\hat{m}_n^{LP}(x)$ as the first element of the vector

$$(\hat{a}_0, \hat{a}_1, \dots, \hat{a}_p)(x) = \underset{a_0, a_1, \dots, a_p \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n \left[y_i - \left(a_0 + a_1(x_i - x) + \dots + a_p \frac{(x_i - x)^p}{p!} \right) \right]^2 K \left(\frac{x_i - x}{h} \right).$$

That is, we fit a local order- p polynomial approximation

$$\hat{a}_0(x) + \hat{a}_1(x)(\cdot - x) + \dots + \hat{a}_p(x) \frac{(\cdot - x)^p}{p!}$$

to the function m at x and evaluate it at x , obtaining $\hat{m}_n^{LP}(x) = \hat{a}_0(x)$. 12

Remark: The N-W estimator is the LP estimator of order $l=0$.

Exercise: (i) Find matrices U_x and K_x and the vector \underline{Y} such that

$$\hat{g}_2(x) := (\hat{a}_0, \hat{a}_1, \dots, \hat{a}_2)(x) = (U_x^T K_x U_x)^{-1} U_x^T K_x \underline{Y}$$

(ii) Show that we may write

$$\hat{m}_n^{LP}(x) = \sum_{i=1}^n W_{ni}(x) Y_i$$

with

$$W_{ni}(x) = \frac{1}{nh} e_i^T \left(\frac{1}{nh} U_x^T K_x U_x \right)^{-1} \underbrace{U_{x,i}}_{\text{row } i \text{ of } U_x} \cdot K\left(\frac{X_i - x}{h}\right).$$

Solution: Let

$$U_x = \begin{bmatrix} 1 & \dots & (X_1 - x) & \dots & \frac{1}{2} (X_1 - x)^2 & \dots & \frac{1}{l!} (X_1 - x)^l \\ \vdots & & \vdots & & \vdots & & \vdots \\ 1 & \dots & (X_n - x) & \dots & \frac{1}{2} (X_n - x)^2 & \dots & \frac{1}{l!} (X_n - x)^l \end{bmatrix}$$

$$K_x = \text{diag} \left(K\left(\frac{X_1 - x}{h}\right), \dots, K\left(\frac{X_n - x}{h}\right) \right).$$

We present a result from Tsybakov's book, which depends on these assumptions:

(LP1) There exists a $\lambda_0 > 0$ and $n_0 > 0$ such that

$$\lambda_{\min} \left(\frac{1}{nh} \mathcal{U}_x^T K_x \mathcal{U}_x \right) \geq \lambda_0 > 0 \quad \text{for all } n \geq n_0.$$

[Points X_1, \dots, X_n distributed such that no interval in $[0, 1]$ stays empty.]

(LP2) There exists a $q_0 > 0$ such that for every $A \subset [0, 1]$,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \in A) \leq q_0 \left(\text{Leb}(A) \vee \frac{1}{n} \right).$$

[Points X_1, \dots, X_n not too concentrated in any one small interval]

(LP3) K with support on $[-1, 1]$ and $0 \leq K(u) \leq K_{\max} < \infty \quad \forall u \in \mathbb{R}$.

Theorem 2: let X_1, \dots, X_n be deterministic such that (LP1), (LP2), and (LP3) hold under $h = \alpha n^{-1/(2\beta+1)}$. Let $m \in \mathcal{H}(\beta, L)$ on $[0, 1]$. Then

$$\sup_{m \in \mathcal{H}(\beta, L)} \sup_{x \in [0, 1]} \mathbb{E} \left[\hat{m}_n^{\text{LP}}(x) - m(x) \right]^2 \leq n^{-\frac{2\beta}{2\beta+1}} \cdot C$$

for all $n \geq n_0$, where $\hat{m}_n^{\text{LP}}(x)$ is of order $\ell = \beta - 1$. The constant C depends only on $\beta, L, \lambda_0, q_0, K_{\max}, \sigma^2$, and α .

In practice, the N-W and local linear ($\ell = 1$) polynomial estimators are used very often. Making ℓ larger can lead to numerical issues...

MSE of N-W estimator under bounded 2nd derivative

Much of the nonparametric literature is written under the settings of the following theorem (filched from Larry Wasserman's lecture notes).

Let $\hat{m}_n^{LP-2}(x)$ denote the LP estimator with $q=1$ ("local linear").

Theorem 3: Let X_1, \dots, X_n have a continuous, differentiable density f , $f(x) > 0$, and suppose m'' is continuous and bounded. Then if $h_n \rightarrow 0$ and $nh_n \rightarrow \infty$, we have

$$\text{Bias } \hat{m}_n^{NW}(x) = \frac{h^2}{2} \left(m''(x) + \frac{2 m'(x) f'(x)}{f(x)} \right) \int_{\mathbb{R}} u^2 K(u) du + o_p(h^2)$$

$$\text{Bias } \hat{m}_n^{LP-1}(x) = \frac{h^2}{2} m''(x) \int_{\mathbb{R}} u^2 K(u) du + o_p(h^2),$$

while $\text{Var } \hat{m}_n^{NW}(x)$ and $\text{Var } \hat{m}_n^{LP-1}(x)$ are given by

$$\frac{1}{nh} \frac{\sigma^2}{f(x)} \int_{\mathbb{R}} K^2(u) du + o_p((nh)^{-1})$$

The proof of this result is more cumbersome than that of our result under Lipschitz smoothness, but it offers some insights.

The optimal rate is $O(n^{-4/5})$ under these settings (better than under Lipschitz smoothness, due to existence of 2 derivatives).

We see how the curvature $m''(x)$ contributes to the bias.

Also how a scarcity of X values (small $f(x)$) can inflate both variance and bias.

Appendix:

For $Y_i = m(X_i) + \varepsilon_i$, $i=1, \dots, n$, $\varepsilon_1, \dots, \varepsilon_n$ iid with $\mathbb{E}\varepsilon_i = 0$, $\mathbb{E}\varepsilon_i^2 = \sigma^2$, $\mathbb{E}\varepsilon_i^4 = \mu_4 < \infty$
and $m \in \text{Lipschitz}(L)$ on $[0,1]$, $X_1, \dots, X_n \in [0,1]$ we have

$$\begin{aligned} \text{Var}\left(\left(Y_{(i+1)} - Y_{(i)}\right)^2\right) &= \text{Var}\left[\left(m(X_{(i+1)}) - m(X_{(i)}) + \varepsilon_{(i+1)} - \varepsilon_{(i)}\right)^2\right] \\ &= \text{Var}\left[\left(m(X_{(i+1)}) - m(X_{(i)})\right)^2 + \left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)^2\right. \\ &\quad \left.+ 2\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)\left(m(X_{(i+1)}) - m(X_{(i)})\right)\right] \\ &= \text{Var}\left[\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)^2 + 2\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)\left(m(X_{(i+1)}) - m(X_{(i)})\right)\right] \\ &\leq 2 \text{Var}\left[\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)^2\right] + 2 \text{Var}\left[2\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)\left(m(X_{(i+1)}) - m(X_{(i)})\right)\right] \\ &\leq 2 \mathbb{E}\left[\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)^4\right] + 8 \underbrace{\mathbb{E}\left[\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)^2\right]}_{2\sigma^2} \underbrace{\left|m(X_{(i+1)}) - m(X_{(i)})\right|^2}_{\leq L^2 |X_{(i+1)} - X_{(i)}|^2} \\ &\leq 2\left(2\mu_4 + 6\sigma^4\right) + 16\sigma^2 L^2, \quad \underbrace{\left(X_1, \dots, X_n \in [0,1]\right)}_{\leq L^2} \end{aligned}$$

for all $i=1, \dots, n$, since

$$\begin{aligned} \mathbb{E}\left(\varepsilon_{(i+1)} - \varepsilon_{(i)}\right)^4 &= \mathbb{E}\left[\varepsilon_{(i+1)}^4 - 4\varepsilon_{(i+1)}^3\varepsilon_{(i)} + 6\varepsilon_{(i+1)}^2\varepsilon_{(i)}^2 - 4\varepsilon_{(i+1)}\varepsilon_{(i)}^3 + \varepsilon_{(i)}^4\right] \\ &= \mu_4 + 6\sigma^2 \cdot \sigma^2 + \mu_4 \\ &= 2\mu_4 + 6\sigma^4. \end{aligned}$$