

Nonparametric regression

- Least-square estimators (for any function class with covering/packing bounds)
- Fourier projection. (Sobolev class, requires basis functions)
- "Local" estimators. (Hölder class, pointwise bounds).
(does not need basis functions).

Warmup: Nadaraya-Watson estimator.

Problem: x_i 's deterministic,

(Fixed design analysis)

$$\begin{aligned} Y_i &= f^*(x_i) + \varepsilon_i \\ \varepsilon_i &\stackrel{\text{iid}}{\sim} P. \quad \mathbb{E}[\varepsilon_i] = 0 \\ &\quad \text{var}(\varepsilon_i) \leq \sigma^2. \end{aligned}$$

(For random design, first condition on $(X_i)_{i=1}^n$).

$$\hat{f}_n(x) = \frac{\sum_{i=1}^n Y_i \cdot K\left(\frac{x_i - x}{h}\right)}{\sum_{i=1}^n K\left(\frac{x_i - x}{h}\right)}.$$

(for simplicity, we define weight function).

$$W_{n,i}(x) = \frac{K\left(\frac{x-x_i}{h}\right)}{\sum_{j=1}^n K\left(\frac{x-x_j}{h}\right)}.$$

$$\hat{f}_n(x) = \sum_{i=1}^n W_{n,i}(x) Y_i.$$

Note: \hat{f}_n is a linear function of $(Y_i)_{i=1}^n$.

$W_{n,i}$ are deterministic.

(cf. Tsybakov Chapter 3. best linear estimation).

Bias: $b(x_0) = \sum_{i=1}^n W_{n,i}(x_0) \cdot (f^*(x_i) - f^*(x_0))$

Variance: $\sigma^2(x_0) = \sigma^2 \sum_{i=1}^n W_{n,i}^2(x_0)$.

Assuming $f^* \in \Sigma(\beta, L)$ for $0 < \beta \leq 1$

$$|b(x_0)| \leq \sum_{i=1}^n |W_{n,i}(x_0)| \cdot L \cdot |x_i - x_0|^\beta$$

(Assuming that K is supported on $[-1, 1]$)

$W_{n,i}(x_0)$ is non-zero only when

$$|x_i - x_0| \leq h$$

$$\leq \sum_{i=1}^n |W_{n,i}(x_0)| \cdot L h^\beta.$$

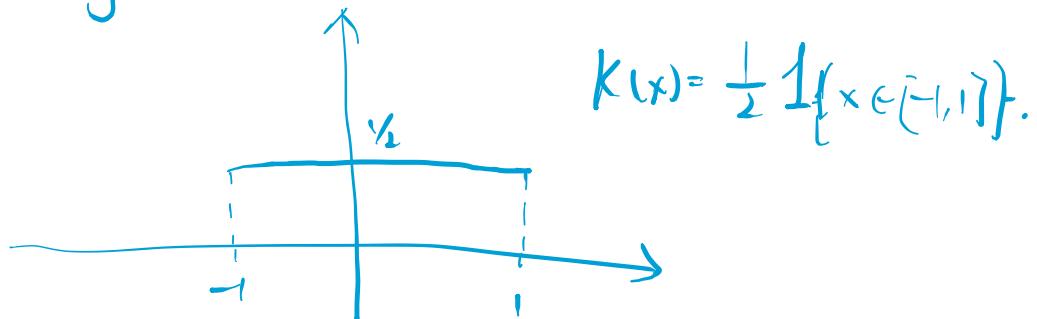
(Assuming $K \geq 0$).

$$= L \cdot h^\beta$$

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

$$\leq \sigma^2 \cdot \underbrace{\left(\sum_{i=1}^n W_{n,i}(x_0) \right)}_{\geq 1} \cdot \underbrace{\max_i |W_{n,i}(x_0)|}_{\leq \frac{1}{\#\{j : |x_j - x_0| \leq h\}}}.$$

Taking box-shaped kernel



$$W_{n,i}(x_0) = \frac{\mathbf{1}_{\{x_i \in [x_0-h, x_0+h]\}}}{\#\{j, \text{ s.t. } x_j \in [x_0-h, x_0+h]\}}$$

For equi-spaced design $x_i = i/n$.

(this can be relaxed, e.g. distance between adjacent design points).

$$\left| \{j : |x_j - x_0| \leq h\} \right| \geq \lfloor 2n/h \rfloor \quad \text{for } h > n.$$

$$J^2(x_0) \leq \frac{\sigma^2}{nh}.$$

$$\cdot h_n = n^{\frac{-1}{2\beta+1}} \Rightarrow \text{MSE}(x_0) = C \cdot n^{\frac{-2\beta}{2\beta+1}}.$$

How about $\beta > 1$? Local polynomial estimation.

Recall NW estimator

$$\hat{f}_n(x) = \arg_{\theta \in \mathbb{R}} \min \sum_{i=1}^n (Y_i - \theta)^2 K\left(\frac{x_i - x}{h}\right).$$

(In general, we may use a polynomial instead of const)

$$f(x_i) \approx f(x) + f'(x)(x_i - x) + \frac{1}{2!} f''(x)(x_i - x)^2 + \dots + \frac{f^{(l)}(x)}{l!}(x_i - x)^l$$

$$= \begin{bmatrix} 1 \\ x_i - x \\ (x_i - x)^2 / 2! \\ \vdots \\ (x_i - x)^l / l! \end{bmatrix}^T \cdot \begin{bmatrix} f(x) \\ f'(x) \\ \vdots \\ f^{(l)}(x) \end{bmatrix}.$$

Notation: $U(t) = \begin{bmatrix} 1 \\ t \\ t^2 / 2! \\ \vdots \\ t^l / l! \end{bmatrix} \in \mathbb{R}^{l+1}$

$$\hat{\theta}_n(x) = \arg \min_{\theta \in \mathbb{R}^{d+1}} \sum_{i=1}^n \left(Y_i - \theta^T U\left(\frac{x_i-x}{h}\right) \right)^2 \cdot K\left(\frac{x_i-x}{h}\right)$$

(By taking K as box function, we're approximating f^* in $[x-h, x+h]$ using polynomials).

$$\hat{f}_n(x) = e_i^T \hat{\theta}_n(x) \quad (\text{first coordinate})$$

(We can also use $e_i^T \hat{\theta}_n(x)$ to estimate $f'(x)$, etc.)

Analysis. We define

$$B_{n,x} := \frac{1}{nh} \sum_{i=1}^n U\left(\frac{x_i-x}{h}\right) U\left(\frac{x_i-x}{h}\right)^T \cdot K\left(\frac{x_i-x}{h}\right)$$

(matrix in the quadratic form)

then we have

$$W_{n,i}(x) = \frac{1}{nh} e_i^T \cdot B_{n,x}^{-1} \cdot U\left(\frac{x_i-x}{h}\right) \cdot K\left(\frac{x_i-x}{h}\right)$$

$$\hat{f}_n(x) = \sum_{i=1}^n W_{n,i}(x) Y_i.$$

Variance:

$$\begin{aligned} \sigma^2(x_0) &= \sigma^2 \cdot \sum_{i=1}^n W_{n,i}(x_0)^2 \\ &\leq \sigma^2 \cdot \sum_{i=1}^n |W_{n,i}(x_0)| \cdot \max_{1 \leq i \leq n} |W_{n,i}(x_0)|. \end{aligned}$$

$$\max_{1 \leq i \leq n} |W_{n,i}(x_0)| \leq \frac{1}{nh} \cdot \max_{i: |x_i - x_0| \leq h} \|B_{n,x}^{-1} U\left(\frac{x_i - x_0}{h}\right)\|_2.$$

$$\leq \frac{1}{nh} \cdot \|B_{n,x}^{-1}\|_{op} \cdot \max_{|x_i - x_0| \leq h} \|U\left(\frac{x_i - x_0}{h}\right)\|_2.$$

$$\|U\left(\frac{x_i - x_0}{h}\right)\|_2^2 \leq \sum_{k=0}^l \frac{1}{(k!)^2} \leq 3.$$

(See Tsybakov). If using equispaced design,

$B_{n,x} \rightarrow$ something positive definite
(defined by integral)

$\exists n_0$, s.t. when $n > n_0$ $B_{n,x} \succcurlyeq \lambda_0 I_{d+1}$.

(Assuming $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} P$, with decoupling bounds)

We can verify using empirical process tools

$$\|B_{n,x} - \frac{\mathbb{E}[U\left(\frac{X-x}{h}\right) U\left(\frac{X-x}{h}\right)^T K\left(\frac{X-x}{h}\right)]}{h}\|_{op} \leq O\left(\frac{1}{\sqrt{nh}}\right)$$

w.h.p. then we can verify PD of $\mathbb{E}[\dots]$.

For our analysis we assume design points are good enough to satisfy $B_{n,x} \succcurlyeq \lambda_0 I_{d+1}$.

$$\max_{1 \leq i \leq n} |W_{n,i}(x_0)| \leq \frac{C}{\lambda_0 nh}.$$

$$\sum_{i=1}^n |W_{n,i}(x_0)| \leq \frac{a_0 \cdot C}{\lambda_0}.$$

(Assuming that for any interval $A \subseteq [0,1]$

$$\frac{1}{n} |\{i : x_i \in A\}| \leq a_0 \cdot \max\{1, \frac{1}{n}\}$$

This can be verified using prob assumptions).

$$\text{So } \sigma^2(x_0) \leq \frac{a_0 \cdot C^2}{\lambda_0^2 nh}.$$

$$\text{Bias. } b(x_0) = \sum_{i=1}^n W_{n,i}(x_0) \cdot \left(f^*(x_i) - f^*(x_0) \right)$$

$$= \sum_{i=1}^n W_{n,i}(x_0) \cdot \sum_{k=1}^l \boxed{(f^*)^{(k)}(x_0) \cdot \frac{(x_i - x_0)^k}{k!}}$$

$$+ \sum_{i=1}^n (x_i - x_0)^l W_{n,i}(x_0) \cdot \boxed{\frac{(f^*)^{(l)}(x_0 + \tau_i(x_i - x_0)) - (f^*)^{(l)}(x_0)}{l!}}$$

$$|...| \leq |x_i - x_0|^{\beta-l} \cdot L.$$

$$|...| \leq L \cdot h^\beta.$$

Hope: make the first term cancel out.

Lemma. For a degree- t polynomial Q .

We have $\sum_{i=1}^n W_{n,i}(x) \cdot Q(x_i) = Q(x)$.

(Proof. see Tsybakov).

Using this lemma, we have

$$\sum_{i=1}^n W_{n,i}(x) \cdot (x_i - x)^k = 0 \quad \text{for } k=1, 2, \dots, t$$

$$\sum_{k=1}^t \frac{(f^*)^{(k)}(x_0)}{k!} \sum_{i=1}^n W_{n,i}(x_0) \cdot (x_i - x_0)^k = 0$$

Therefore, we can bound the bias as

$$|b(x_0)| \leq L \cdot h^\beta.$$

By choosing $h_n = n^{-\frac{1}{2\beta+1}}$, $MSE(x_0) \leq n^{-\frac{2\beta}{2\beta+1}}$.

Question: minimax optimality?

Recap: Le Cam's two-point method:

Construct $f_1, f_2 \in \mathcal{F}$.

$$\inf_{\hat{T}} \sup_{f \in \mathcal{F}} \mathbb{E}[(\hat{T} - T(f))^2] \geq \frac{1}{8} (T(f_1) - T(f_2))^2 \cdot [1 - d_{TV}(P_1, P_2)]$$

Application to nonparametric regression.

(results for density estimation are similar).

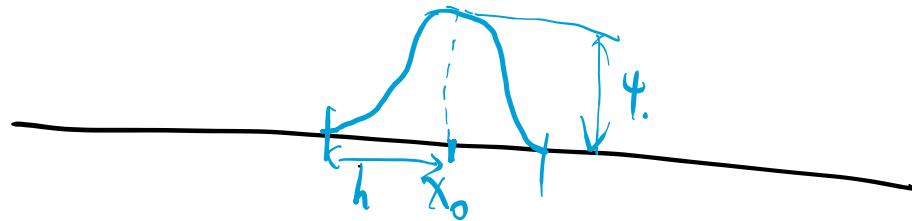
x_1, x_2, \dots, x_n deterministic, satisfying spacing assumptions
(e.g. equispaced).

$$Y_i = f^*(x_i) + \varepsilon_i, \quad \varepsilon_i \sim N(0, 1).$$

$$f^* \in \mathcal{F} = \Sigma(\rho, L). \quad T(f) = f(x_0).$$

WLOG $f_2(x_0) = 0$ ($\forall x$)

$$f_1(x) = \psi \cdot K\left(\frac{x-x_0}{h}\right)$$



Idea: we want large ψ and small h .

Choose K : C^∞ smooth, and bounded support.

$$K(u) = \exp\left(\frac{-1}{1-u^2}\right) \cdot \mathbb{1}_{\{|u| < 1\}}.$$

Under this construction.

$$\cdot |f_1(x_0) - f_2(x_0)| = \psi.$$

$$\begin{aligned}
 d_{TV}(P_1, P_2) &\leq \sqrt{\frac{1}{2} D_{KL}(P_1 \parallel P_2)} \\
 &= \sqrt{\frac{1}{2} \sum_{j=1}^n D_{KL}(P_{1,j} \parallel P_{2,j})} \\
 (\text{where } P_{j,i} &\text{ is the dist of } i^{\text{th}} \text{ obs under } P_j) \\
 &= \sqrt{\frac{1}{2} \sum_{j=1}^n (f_1(x_i) - f_2(x_i))^2} \\
 &\leq \sqrt{\frac{4^2}{2} \cdot | \{i : |x_i - x_0| \leq h\} |} \leq \sqrt{\frac{a_0 4^2 n h}{2}}
 \end{aligned}$$

Still need to verify $f_i \in \Sigma(\theta, L)$.

$$f_i(x) = 4 \cdot K\left(\frac{x - x_0}{h}\right).$$

$$f^{(1)}(x) = \frac{4}{h^2} K'\left(\frac{x - x_0}{h}\right)$$

$$|f^{(1)}(x) - f^{(1)}(y)| \leq \frac{4}{h^2} \cdot C_1 \cdot \frac{|x - y|}{h}$$

$$\leq \frac{4}{h^2} \cdot C_1 \cdot \frac{|x - y|^{\beta-1} \cdot (2h)^{1-(\beta-1)}}{h}$$

$$\leq 2C_1 \cdot |x - y|^{\beta-1} \cdot \frac{4}{h^\beta} \quad (\leq L|x - y|^{\beta-1})$$

need

$$\text{Choose } \psi = \frac{L h^\beta}{2 C_1}$$

Choose c_0 small enough with $h_n = c_0 n^{-\frac{1}{2\beta+1}}$.

$$d_{TV}(P_1, P_2) \leq \sqrt{\frac{a_0}{2} \cdot n \cdot h_n \cdot \left(\frac{L h_n^\beta}{2 C_1}\right)^2}$$

$$\sqrt{\text{Const. } n \cdot (h_n)^{2\beta+1}} < \frac{1}{2}.$$

$$\inf_{\hat{T}} \sup_{f \in \Sigma(P, U)} \mathbb{E}[|f(x_0) - \hat{T}|^2] \geq \frac{G^2}{16} \cdot C' \cdot n^{-\frac{2\beta}{2\beta+1}}$$

- Remark: for MISE, we need something more "Fano's method"
- Le Cam finds a pair of hard instances.
- Fano deals with difficulties from a large collection of problem instances using mutual info arguments.

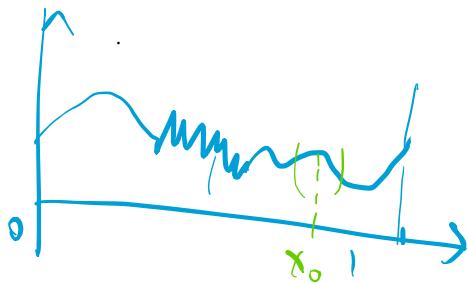
Adaptivity:

Question: Can we do well w/o knowing β ?

Want to construct a single \hat{f}_n .

$$\text{st. } \forall \beta, \forall f^* \in \Sigma(\beta, L), \quad \text{MSI}(x_0) \leq n^{-\frac{2\beta}{2\beta+1}}.$$

In general, f^* may have different smoothness at different places.



This also allows us to adapt to local smoothness around x_0 .

Surprisingly, this is not possible. (w/o additional log factor).

Thm (Lepski, Lower bound).

Suppose $0 < \beta_1 < \alpha < \beta_2 \leq 1$.

$$r_1^2(n) = \left(\frac{\log n}{n} \right)^{\frac{2\beta_1}{2\beta_1+1}}, \quad r_2^2(n) = \left(\frac{1}{n} \right)^{\frac{2\alpha}{2\alpha+1}}$$

$$\inf_{\hat{f}} \sup_{\{i=1,2\}} \sup_{f^* \in \Sigma(\beta_i, L)} \frac{\mathbb{E}[\hat{f}_n(x_0) - f^*(x_0)]^2}{r_i^2(n)} \geq c.$$

- Asymetric lower bound. Penalize more

for mls-estimating the easier class.

- If we want to achieve anything faster than $n^{-\frac{2\beta_1}{2\beta_1+1}}$ under $\Sigma(\beta_1, L)$.

then this estimator must pay log factor under $\Sigma(\beta_1, L)$.

Proof: from Hw!, asymmetric two-point lemma.

$$f_2^* \in \Sigma(\beta_1, L), \quad f_1^* \in \Sigma(\beta_1, L).$$

$$\inf_{\hat{f}} \sup_{i \in \{1, 2\}} \frac{\mathbb{E}[\hat{f}_n(x_0) - f_i^*(x_0)]^2}{r_i^2(n)} = \varphi^2, \text{ under construction.}$$
$$\geq \frac{(f_1^*(x_0) - f_2^*(x_0))^2}{8 \cdot r_1^2(n)} \left[1 - \frac{r_2^2(n)}{r_1^2(n)} \left(+ \chi^2(P_f, \| P_{f_1} \right) \right].$$

Conclusion: $f_2^* = 0$.

$$f_1^* = \varphi \cdot K\left(\frac{x - x_0}{h}\right).$$

Need to ensure.

$$\left(\frac{r_2^2(n)}{r_1^2(n)} \right) \cdot \left\{ 1 + \chi^2(P_{f_1} \parallel P_{f_2}) \right\} < \frac{1}{2}.$$

$$\geq \frac{1}{\text{poly}(n)}.$$

We only need

$$\chi^2(P_{f_1} \parallel P_{f_2}) \leq \text{poly}(n)$$

(e.g. $n^{\frac{a-\beta_1}{2}}$)

(Cf. for symmetric version, we need $D_{KL}(P_{f_1} \parallel P_{f_2}) \leq \frac{1}{2}$).

By choosing $\begin{cases} h = C_0 \cdot \left(\frac{\log n}{n}\right)^{\frac{1}{2\beta_1+1}} \\ \psi = C_1 \cdot h^{\beta_1} \end{cases}$

We can ensure $\chi^2(\dots) \leq \text{poly}(n)$.

Substituting back, we prove the result.

How to achieve it? Lepski's method.

Key observation:

for each β , let h_β be optimal bandwidth.

Let $\hat{f}_{h_\beta}(x_0)$ be the estimator using h_β

Let β^* be the true smoothness.

If we under-smooth, $\beta < \beta^*$.

$$|\hat{f}_{h\beta}(x_0) - \hat{f}_{h\beta^*}(x_0)| \leq C h_\beta^\beta$$

(w.h.p.) for any $0 < \beta < \beta^*$.

and if we use this test for $\beta > \beta^*$,
this may not be satisfied.

Step I. discretize $\beta^* \in [\beta_{\min}, \beta_{\max}]$

$$\mathcal{B} = \{\beta_{\min} = \beta_1 < \beta_2 < \dots < \beta_N = \beta_{\max}\}$$

where $N = \log n$. $\beta_j - \beta_{j-1} = \frac{1}{\log n}$.

For $\beta \in [\beta_{j-1}, \beta_j]$.

$$C \cdot n^{-\frac{2\beta_{j-1}}{2\beta_{j-1}+1}} \leq n^{-\frac{2\beta}{2\beta+1}} \leq n^{-\frac{2\beta_j}{2\beta_j+1}}$$

The rates within each interval are the same.

Step II.

$$\hat{\beta} = \max \left\{ \beta \in B : \begin{array}{l} \forall \beta' \in B, \beta' < \beta, \\ |\hat{f}_{h_{\beta'}}(x_0) - \hat{f}_{h_\beta}(x_0)| \leq C \cdot h_{\beta'}^{\beta'} \end{array} \right\}$$

then estimate by

$$\hat{f}_{\text{Lepick!}}(x_0) = \hat{f}_{h_{\hat{\beta}}}(x_0).$$

Thm:

$$E[|\hat{f}_{\text{Lepick!}}(x_0) - f^*(x_0)|^2] \leq C \left(\frac{\log n}{n} \right)^{\frac{2\beta_*}{2\beta_* + 1}}.$$

Proof idea: Define event

$$\mathcal{E}_j = \{\hat{\beta} = \beta_j\}.$$

$$E[|\hat{f}(x_0) - f^*(x_0)|^2] = \sum_{j=1}^N E[|\hat{f}(x_0) - f^*(x_0)|^2 \cdot 1_{\mathcal{E}_j}].$$

(i) For $j \geq j^*$. On the event \mathcal{E}_j ,

$$|\hat{f}_{h_{\beta^*}}(x_0) - \hat{f}_{h_{\beta_j}}(x_0)| \leq C \cdot h_{\beta^*}^{\beta^*}.$$

Even if we over estimate β^* , it doesn't matter
(by our screening criterion).

$$\begin{aligned}
& \sum_{j=j^*}^N \mathbb{E} \left[|\hat{f}(x_0) - f^*(x_0)|^2 \mathbf{1}_{\varepsilon_j} \right] \\
& \leq 2 \cdot \underbrace{\sum_{j=j^*}^N \mathbb{E} \left[|\hat{f}_{hp^*}(x_0) - f^*(x_0)|^2 \cdot \mathbf{1}_{\varepsilon_j} \right]}_{\text{by NW analysis}} \\
& \quad + 2 \cdot \underbrace{\sum_{j=j^*}^N \mathbb{E} \left[|\hat{f}_{hp_j}(x_0) - \hat{f}_{hp^*}(x_0)|^2 \cdot \mathbf{1}_{\varepsilon_j} \right]}_{\leq h_{p^*}^{2\beta^*} \text{ by def } \varepsilon_j}.
\end{aligned}$$

(ii) $j < j^*$. We can show using Gaussian tail bound.

$$\mathbb{P}(\varepsilon_j) \leq n^{-c'} \quad (\text{for } c' \text{ large enough}).$$

this requires $h_{p_j} = \left(\frac{\log n}{n} \right)^{\frac{1}{2\beta_j+1}}$.