

Fixed design: x_1, x_2, \dots, x_n deterministic

$$Y_i = f^*(x_i) + \varepsilon_i \quad (\text{e.g. denoising})$$

Random design. $(X_i, Y_i)_{i=1}^n \stackrel{iid}{\sim} P$

$$f^*(x) = E[Y_i | X_i = x]$$

Goal: generate an \hat{f}_n s.t. $\|\hat{f}_n - f^*\|_{L^2(P)}$ small.

(Operationally, given a new data point x
predict Y , s.t. err is small.)

$$(E[Y_i | X_i = x])$$

Constrained LS.

$$\hat{f}_n = \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} \frac{1}{n} \sum_{i=1}^n (Y_i - f(x_i))^2.$$

(Assuming $f^* \in \mathcal{F}$)

From first-order-condition

$$\frac{1}{n} \sum_{i=1}^n (\hat{f}_n - f^*)(x_i)^2 \leq \frac{1}{n} \sum_{i=1}^n (Y_i - f^*(x_i)) \cdot (\hat{f}_n - f^*)(x_i)$$

Question: relate $\frac{1}{n} \sum_i^n (\hat{f}_n - f^*)(x_i)^2$ to $\|\hat{f}_n - f^*\|_{L^2(P)}$?

Solution:

$$\begin{aligned} Z_n(r) &:= \mathbb{E} \left[\sup_{\substack{\|h\|_{L^2(P)} \leq r \\ h \in \mathcal{F}^*}} \left| (\hat{P}_n - P) h^2 \right| \right] \\ &= \mathbb{E} \left[\sup_{\dots} \left| \frac{1}{n} \sum_i^n h(x_i)^2 - \mathbb{E}[h(x)^2] \right| \right] \\ &\leq 2 \cdot \mathbb{E} \left[\sup_{\dots} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(x_i)^2 \right| \right]. \end{aligned}$$

Once we're able to bound this, we'll have

$$\begin{aligned} R_n^2 &= \|\hat{f}_n - f^*\|_{L^2(P)}^2 \leq \frac{1}{n} \sum_i^n \mathbb{E}[(\hat{f}_n - f^*)(x_i)^2] + Z_n(\|\hat{f}_n - f^*\|_{L^2}) \\ &\leq \mathbb{E} \left[\sup_{\substack{h \in \mathcal{F} \\ \|h\|_{L^2} \leq M}} \left| \frac{1}{n} \sum_i^n \varepsilon_i h(x_i)^2 \right| \right] + Z_n(M) \end{aligned}$$

(w.h.p.)

using the same argument in Lec 8.

From hw 2. Assuming $\forall f \in \mathcal{F}$, $\|f\|_\infty \leq M$.

$x \mapsto x^2$ is a $2M$ -Lipschitz mapping
when $x \in [-M, M]$.

So we have

$$\begin{aligned} & \mathbb{E} \left[\sup_{h \in \mathcal{F}^*(r)} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(x_i) \right| \right] \\ & \leq 2M \cdot \mathbb{E} \left[\sup_{h \in \mathcal{F}^*(r)} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(x_i) \right| \right]. \end{aligned}$$

Applying the arguments we have seen
yields the bounds.

- To improve failure prob, use concentration ineq for sup of empirical process.
- The contraction-based bound here is not tight.

e.g. If we have $y_i = f^*(x_i) + \varepsilon_i$
 $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$.

When σ is small (high signal, low noise).

Or even w/ $\sigma = 0$

Faster rates are expected, but not achieved
using above arguments.

(c.f. Mendelson (2012), "small ball argument").

Kernel density estimation.

$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} p^* \in \mathcal{P}$.

Possible methods:

• MLE $\hat{p} = \underset{p \in \mathcal{P}}{\operatorname{arg\max}} \frac{1}{n} \sum_{i=1}^n \log p(X_i)$

(then use empirical process to study the error).

• "skleton estimator" / "integral probability metric".

(see Hw 3)

(theoretically nice, but computationally infeasible).

• Local methods for smoothness classes.

Suppose we have $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} p^*$ on \mathbb{R}

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

where K is a kernel function

satisfying $\int_{\mathbb{R}} K(x) dx = 1$.

e.g. $K = \frac{1}{2} I_{[-1,1]}$. local averaging

Analys: at the point x_0 (deterministic)

$$\text{Var}(\hat{P}_n(x_0)) = \frac{1}{h^2 n} \text{Var}\left(K\left(\frac{x-x_0}{h}\right)\right).$$

$$\leq \frac{1}{h^2 n} \cdot \int K^2\left(\frac{y-x_0}{h}\right) \cdot p(y) dy$$

$$\leq \underbrace{\frac{p_{\max}}{h n} \int K^2(x) dx}_{\text{const.}}$$

Bias.

Assuming p is Hölder w/ exponent β .
 $(0 < \beta \leq 1)$.

$$|p(x) - p(y)| \leq L \cdot |x-y|^\beta \quad (\forall x, y).$$

$$\mathbb{E}[\hat{P}_n(x_0)] - p(x_0)$$

$$= \frac{1}{h} \int_{\mathbb{R}} K\left(\frac{y-x_0}{h}\right) \cdot (p(y) - p(x_0)) dy.$$

$$\leq \frac{L}{h} \int_{\mathbb{R}} K\left(\frac{y-x_0}{h}\right) \cdot |y-x_0|^\beta dy$$

$$\leq h^\beta \cdot \underbrace{L \int K(u) \cdot |u|^\beta du}_{\text{const.}} \quad (\text{Assume } \int K(u) \cdot |u|^\beta du < \infty)$$

Bias-var trade off.

$$E[\hat{P}_n(x_0) - P^*(x_0)]^2 \leq C \cdot \frac{1}{nh} + C' \cdot h^{2\beta}$$

(Optimal $h_n = C_1 \cdot n^{-\frac{1}{2\beta+1}}$)

$$E[\hat{P}_n(x_0) - P^*(x_0)]^2 \leq C_2 n^{-\frac{2\beta}{2\beta+1}}.$$

What if we have more smoothness?

$$\begin{aligned} & p(x_0 + uh) - p(x_0) \\ &= p'(x_0) \cdot uh + \frac{p''(x_0)}{2!} (uh)^2 + \dots + \frac{p^{(t)}(x_0)}{(t-1)!} (uh)^{t-1} \\ & \quad + \underbrace{\frac{p^{(t)}(x_0 + \tau_u uh)}{t!} \cdot (uh)^t}. \end{aligned}$$

(for some $\tau_u \in [0, 1]$).

Here we assume $p \in \text{H\"older}(\beta, L)$

where $\beta = t + \gamma$, $t \in \mathbb{N}$, $\gamma \in (0, 1]$.

$$B_{\text{less}}(x_0) = \int_{IR} K(u) \cdot \left[p'(x_0) \cdot uh + p''(x_0) \cdot \frac{(uh)^2}{2} + \dots + p^{(\ell)}(x_0) \cdot \frac{(uh)^\ell}{\ell!} \right] du$$

If K is
symmetric,

$$+ \int_{IR} K(u) \frac{p^{(\ell)}(x_0 + tuh) - p^{(\ell)}(x_0)}{(uh)^\ell} du.$$

$$\int K(u) uh du = 0.$$

The entire sum can be
made 0 by careful
choice of K .

$$| \dots | \leq L \int_{IR} |K(u)| \cdot t u |uh|^\beta \cdot (uh)^\ell du$$

$$\leq L \cdot h^\beta \int_{IR} |K(u)| \cdot |u|^\beta du.$$

Note: to cancel (≥ 2) terms, we need to choose
 K that is negative somewhere.

Def. ℓ -th order kernel: K such that

$$j=1, 2, \dots, \ell,$$

$$\int_{IR} u^j K(u) du = 0.$$

(e.g. box and Gaussian are both 1st order).

Construction via Legendre polynomials.

$$(P_k)_{k=0}^{+\infty}$$

orthonormal basis on $L^2([-1, 1])$.

$$K(u) = \sum_{m=0}^l \varphi_m(0) \varphi_m(u) \cdot \{u \in [-1, 1]\}.$$

Easy to show

- $\int_{-1}^1 K(u) du = 1.$

$$K = \sum_{m=0}^l \varphi_m(0) \varphi_m$$

- $\int_{-1}^1 K(u) w^j du = \langle K, w^j \rangle_{L^2}.$

$$w^j = \sum_{m=0}^l b_m \varphi_m$$

$$= \sum_{m=0}^l b_m \cdot \varphi_m(0) = w^j \Big|_{u=0} = 0.$$

(for $j=1, 2, \dots, l$).

Easy to see, $\int_{-1}^1 K(x)^2 dx < +\infty$, $\int_{-1}^1 |K(x)| \cdot |x|^\beta < +\infty$ (H_β).

So we get

$$|\text{bias}(x_0)| \leq C \cdot h^\beta$$

For any $\beta > 0$, we always have

$$\text{MSE}(x_0) \leq C \cdot \left(\frac{1}{nh} + h^{2\beta} \right)$$

Choosing $h = n^{-\frac{1}{2\beta+1}}$, we get the rate $n^{-\frac{2\beta}{2\beta+1}}$.

Not hard to extend to multivariate case
 you'll get the rate $n^{-\frac{2\beta}{2\beta+d}}$.

$$MISE = \int_{\mathbb{R}} MSE(x) dx$$

(For Hölder class, just integrate the bound above.)

$\beta \in \mathbb{N}_+$, Sobolev class

$$\int |P^\beta(x)|^2 dx \leq L^2 < +\infty$$

• Variance.

$$\int_{\mathbb{R}} \sigma^2(x) dx = \frac{1}{nh} \int \text{Var}\left(K\left(\frac{x-x_i}{h}\right)\right) dx$$

$$\leq \frac{1}{nh} \int K^2(x) dx.$$

• Bias (assuming $(\beta-1)$ -th order Kernel).

$$b(x) = \int K(u) \frac{(uh)^\beta}{\beta!} \int_0^1 (1-\tau)^{\beta-1} P^\beta(x+\tau uh) d\tau du$$

$$\int b(x)^2 dx$$

$$\leq \int \left(\int |K(u)| \cdot \frac{|uh|^\beta}{\beta!} \int_0^1 (1-\tau)^{\beta-1} P^\beta(x+\tau uh) d\tau du \right)^2 dx$$

Detour: generalized Minkowski ineq.

• Minkowski ineq.

$$\|g(\cdot, u_1) + g(\cdot, u_2) + \dots + g(\cdot, u_m)\|_2 \leq \sum_{j=1}^m \|g(\cdot, u_j)\|_2.$$

Generalized Minkowski:

$$\left(\int \left| \int g(x,u) du \right|^2 dx \right)^{1/2} \leq \int \left(\int |g(x,u)|^2 dx \right)^{1/2} du.$$

Using this ineq, we have (sup of x):

$$\int b(x)^2 dx \leq h^{2\beta} \left[\int \left(\int K(u) \cdot |u|^\beta \left(\int_0^1 p^{(\beta)}(x+\tau uh) d\tau \right)^2 dx \right)^{1/2} du \right]^2$$

by Cauchy-Schwarz

$$\left(\int_0^1 p^{(\beta)}(x+\tau uh) d\tau \right)^2$$

$$\leq \int_0^1 p^{(\beta)}(x+\tau uh)^2 d\tau$$

$$\leq h^{2\beta} \left[\int_{\mathbb{R}} |K(u)| \cdot |u|^\beta \left(\int_{\mathbb{R}} \int_0^1 p^{(\beta)}(x+\tau uh)^2 d\tau dx \right)^{1/2} du \right]^2 \leq L$$

$$\leq L^2 h^{2\beta} \cdot \left(\int |K(u)| \cdot |u|^\beta du \right)^2$$

Then apply the same trade off.

Question: asymptotic optimal bandwidth?

e.g. in parametric estimation.

$$\|\hat{\theta}_n - \theta^*\|_2 \leq O_p(\sqrt{n}).$$

We can actually show that

$$\text{Asymptotic MSE} \quad \lim_{n \rightarrow \infty} \mathbb{E}[\| \hat{\theta}_n - \theta^* \|_2^2] \geq I(\theta^*)^{-1}$$

achieved asymptotically by MLE.

- Operationally, we can choose estimators by comparing asymptotic MSE.

Can we do the same for KDE

(e.g. choosing h, K).

Do we have sth. like

$$\lim_{n \rightarrow \infty} \mathbb{E}\left[n^{\frac{2\beta}{2\beta+1}} \|\hat{P}_n - P^*\|_2^2\right] \geq \text{sth. ?}$$

Lack of asymptotic optimality.

Let P be a density on \mathbb{R} , $\int (P''(x))^2 dx \leq C$.

Optimal rate: $n^{-4/5}$ (MISE)

Achieved using KDE, 1st order Kernel.

Thm: Suppose we use a second-order kernel,

$$\forall \varepsilon > 0. \text{ take } h = \frac{n^{-4/5}}{\varepsilon} \int K(u)^2 du$$

then we have

$$\limsup_{n \rightarrow \infty} n^{4/5} E \left[\int |\hat{p}_n(x) - p(x)|^2 dx \right] \leq \epsilon.$$

Rmk:

- This is for fixed p , $n \rightarrow \infty$.
If we take finite n minimax risk,
we'll still get $C_0 n^{-4/5}$, where $C_0 > 0$
 $\quad \quad \quad$ is a constant.
- Fixed n , the "hardest" problem in $S(2,L)$
at sample size n depends on n .
- For $S(2,L)$, every density is asymptotically regular.
than the average
(every density in this class is
a little bit more smooth).

Proof idea:

variance:

$$\int \sigma^2(x) dx = \frac{1}{nh} \int K(u)^2 du + o\left(\frac{1}{nh}\right)$$

• bias. $\int b^2(x) dx = \frac{h^4}{4} \left(\int u^2 K(u) du \right)^2 \cdot \left(\int p''(x)^2 dx \right) + o(h^4)$
 by using second-order kernel.

So we trade off $\frac{1}{nh} + o(h^4)$.

Proof idea for the bias term:

$$b(x) = h^2 \int u^2 K(u) \left[\int_0^1 ((1-\tau)p''(x+\tau uh)) d\tau \right] du.$$

$$\tilde{b}(x) = h^2 \int u^2 K(u) \left[\int_0^1 ((1-\tau)p''(x)) d\tau \right] du$$

We have

$$\int_{\mathbb{R}} \tilde{b}(x)^2 dx = \frac{h^4}{4} \left(\int u^2 K(u) du \right)^2 \cdot \left(\int p''(x)^2 dx \right).$$

and $\int |b(x) - \tilde{b}(x)|^2 dx \rightarrow 0$.

Key step: $\|p''(\cdot + h) - p''\|_{L^2} \xrightarrow{h \rightarrow 0} 0$

(by density argument in L^2 ,
 approximation using "nice" functions).