

Nonparametric estimation

• Density estimation

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

Goal: estimate \hat{p}^* (under some metric)

e.g. $|\hat{p}_n(x_0) - p^*(x_0)|$ for given x_0

$$\|\hat{p}_n - p^*\|_{L^\infty}, d_{TV}(\hat{p}_n, p^*)$$

$$D_{KL}(\hat{p}_n \| p^*), D_{KL}(p^* \| \hat{p}_n).$$

$$W_2(\hat{p}_n, p^*) \quad \cdots$$

, Regression. $(X_i, Y_i)_{i=1}^n \stackrel{\text{iid.}}{\sim}$

$$f^*(x) = \mathbb{E}[Y_i \mid X_i = x] \quad (f^* \in \mathcal{F})$$

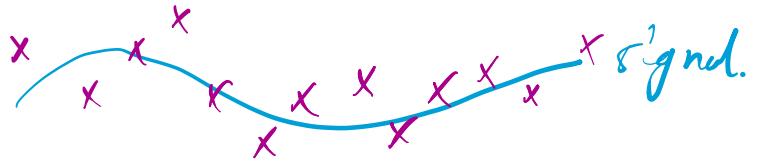
(random design).

Simplification: fixed design

x_1, x_2, \dots, x_n deterministic

Try to recover f^* on x_1, x_2, \dots, x_n
 (n-dim vector).

e.g. denoising.



Fixed - design regression (i.e. denoising problem).

$$y_i = f^*(x_i) + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, 1).$$

(x_i 's deterministic).

MLE \Leftrightarrow least square.

$$\hat{f}_n := \underset{f \in \mathcal{F}}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right\}$$

(Efficiently computable when \mathcal{F} is convex)

First-order optimality condition.



$$Y = (y_i)_{i=1}^n$$

(In literature,
 convexity can be
 relaxed to star-convex)

$$\sum_{i=1}^n (y_i - \hat{f}_n(x_i)) \cdot (f'(x_i) - \hat{f}'_n(x_i)) \leq 0$$

$(\forall f' \in \mathcal{F})$. "Basic ineq".

Substitute with f^* , re-arranging.

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

Denote $\hat{\Delta}_n := \hat{f}_n - f^*$, $\|f\|_n := \sqrt{\frac{1}{n} \sum_i^n f(x_i)^2}$

we have $\|\hat{\Delta}_n\|_n^2 \leq \langle \hat{\Delta}_n, \varepsilon \rangle_n$

(naive bound: $\|\hat{\Delta}_n\|_n^2 \leq \|\hat{\Delta}_n\|_n \cdot \|\varepsilon\|_n \Rightarrow \|\hat{\Delta}_n\|_n \leq \|\varepsilon\|_n$
 this leads to useless bound, $\|\varepsilon\|_n = O(1)$)

We know $\hat{\Delta}_n \in \mathcal{F}^* = \{f - f^* : f \in \mathcal{F}\}$.

$$\langle \hat{\Delta}_n, \varepsilon \rangle_n \leq \sup_{\|h\|_n \leq \|\hat{\Delta}_n\|_n} \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(x_i).$$

Define: "Gaussian complexity"

$$G_n(r) := \mathbb{E} \left[\sup_{\substack{\|h\|_n \leq r \\ h \in \mathcal{F}^*}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(x_i) \right| \right]$$

where $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, 1)$.

We have $\|\Delta_n\|_n^2 \leq \langle (\varepsilon, \hat{\Delta}_n) \rangle_n \leq G_n(\|\hat{\Delta}_n\|_n)$.
 Similar "chicken-egg problem"

Thm: Suppose $G_n(r) \leq \phi_n(r)$

s.t. $\phi_n(cr) \leq C^\alpha \phi_n(r)$ for some $\alpha < 2$. $(*)$.

If f_n solves $f_n = \phi_n(f_n)$

then we have $\|\hat{f}_n - f^*\|_n \leq C \cdot \delta_n$

w.h.p.

(Proof: see lecture 8).

Verifying growth condition on Gaussian complexity.



$\forall f \in F^* \cap B(cr)$, we have $\frac{f}{c} \in F^* \cap B(r)$.

By convexity, we have

$$\text{So } \mathbb{E} \left[\sup_{f \in F^* \cap B(cr)} \langle \varepsilon, f \rangle_n \right] \leq C \mathbb{E} \left[\sup_{f \in F^* \cap B(c)} \langle \varepsilon, f \rangle_n \right]$$

So Gaussian complexities grow sub-linearly.
 $(*)$ is satisfied with $\alpha \leq 1$.

Key observation:

MGF of Rade \leq MGF of $N(0, 1)$

and previous proofs for Rademacher complexity
 are basically bounding it using MGF of $N(0, 1)$.

Thm. $G_n(r) \leq \frac{C}{\sqrt{n}} \int_0^r \sqrt{\log N(\delta; \mathcal{F}^* \cap B_n(r), \| \cdot \|_n)} d\delta$
 (Proof: see lecture 6).

e.g. $\mathcal{F} = \{f: [0, 1] \rightarrow [-1, 1], |f(x) - f(y)| \leq |x - y|, \forall x, y\}$.

$$N(\delta; \mathcal{F}^* \cap B_n(r), \| \cdot \|_n)$$

$$\leq N(\delta; \mathcal{F}^*, \| \cdot \|_n) \leq N(\delta; \tilde{\mathcal{F}}, \| \cdot \|_n)$$

($\mathcal{F}^* \subseteq \tilde{\mathcal{F}} = \{f: [0, 1] \rightarrow [-1, 1], |f(x) - f(y)| \leq 2|x - y|, \forall x, y\}$)

$$\leq \exp(G \text{ fat}_{2\delta}(\tilde{\mathcal{F}})).$$

(Rudelson-Vershynin)

(From Lecture 8), $\text{fat}_{C\delta}(\tilde{F}) \leq C/\delta$.

Applying the results.

$$G_n(r) \leq \frac{C}{\sqrt{n}} \int_0^r \sqrt{\log N(\delta; \tilde{F}, \| \cdot \|_n)} d\delta$$

$$\leq \frac{C}{\sqrt{n}} \int_0^r \sqrt{\frac{C_1 C'}{\delta}} d\delta \leq C \sqrt{\frac{r}{n}}.$$

Solving for fixed pt, $r^2 = \sqrt{\frac{r}{n}}$

$$r_n \asymp n^{-1/3}.$$

So $\|\hat{f}_n - f^*\|_{l_n} \lesssim n^{-1/3}$ w.h.p.

Final exam: Dec 15, noon - 16, 23:59

e.g. In general, β -Hölder class.

$$\beta > 0, \quad \beta = k + \gamma \quad \begin{matrix} k \in \mathbb{N} \\ \gamma \in [0, 1]. \end{matrix}$$

$$\Sigma(\beta, L) := \left\{ f: [0, 1]^d \rightarrow [0, 1], \quad \begin{matrix} \forall \alpha \text{ multi-index} \\ \text{e.g. } |\alpha| = k \end{matrix} \right\}$$

$$|\partial^\alpha f(x) - \partial^\alpha f(y)| \leq L \cdot |x-y|^\beta$$

(A multi-index is

non-negative-integer-valued vector

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d),$$

$$|\alpha| = \sum_i \alpha_i$$

$$\partial^\alpha f = \frac{\partial^{|\alpha|} f}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \cdots \partial x_d^{\alpha_d}}.$$

Thm For any Q ,

$$N(\varepsilon; \sum(\beta_0), \| \cdot \|_{L^2(Q)}) \leq c \exp\left(\left(\frac{c_2}{\varepsilon}\right)^{\frac{d}{2\beta}}\right).$$

(c_1, c_2 may depend on (d, β))

Implication to Hölder regression:

$$G_n(r) \leq \frac{c}{\sqrt{n}} \int_0^r \sqrt{\log N(\delta; \mathcal{F}, \| \cdot \|_n)} d\delta$$

$$\left(\leq \frac{c'}{\sqrt{n}} \int_0^r \left(\frac{1}{\delta}\right)^{\frac{d}{2\beta}} d\delta \text{ may diverge} \right)$$

To avoid divergence, we note that

$$N(\delta; \mathcal{F}, \| \cdot \|_n) \leq N(\delta; [0, 1]^n, \| \cdot \|_n) \leq \left(\frac{c}{\delta}\right)^n$$

$$\int_0^r \sqrt{\log N(\delta; \mathcal{F}, \| \cdot \|_n)} d\delta$$

$$= \int_0^{r_0} \sqrt{n \log(\frac{1}{\delta})} d\delta + c \int_{r_0}^r \left(\frac{1}{\delta}\right)^{\frac{d}{2\beta}} d\delta$$

$\beta > d/2$, pick $r_0 = 0$.

$$G_n(r) \leq \frac{c}{\sqrt{n}} \cdot \int_0^r \delta^{-\frac{d}{2\beta}} d\delta \leq \frac{c'}{\sqrt{n}} r^{1-\frac{d}{2\beta}}$$

Implication to rate: $r_n^2 = \frac{r_n^{1-\frac{d}{2\beta}}}{\sqrt{n}} \Rightarrow r_n \asymp n^{-\frac{\beta}{d+2\beta}}$

- $\beta < d/2$,

$$G_n(r) \leq r_0 + \frac{c}{\sqrt{n}} \int_{r_0}^r s^{-\frac{d}{2\beta}} ds$$

$$\leq r_0 + \frac{c}{\sqrt{n}} \cdot \left(r_0^{1-\frac{d}{2\beta}} - r^{1-\frac{d}{2\beta}} \right)$$

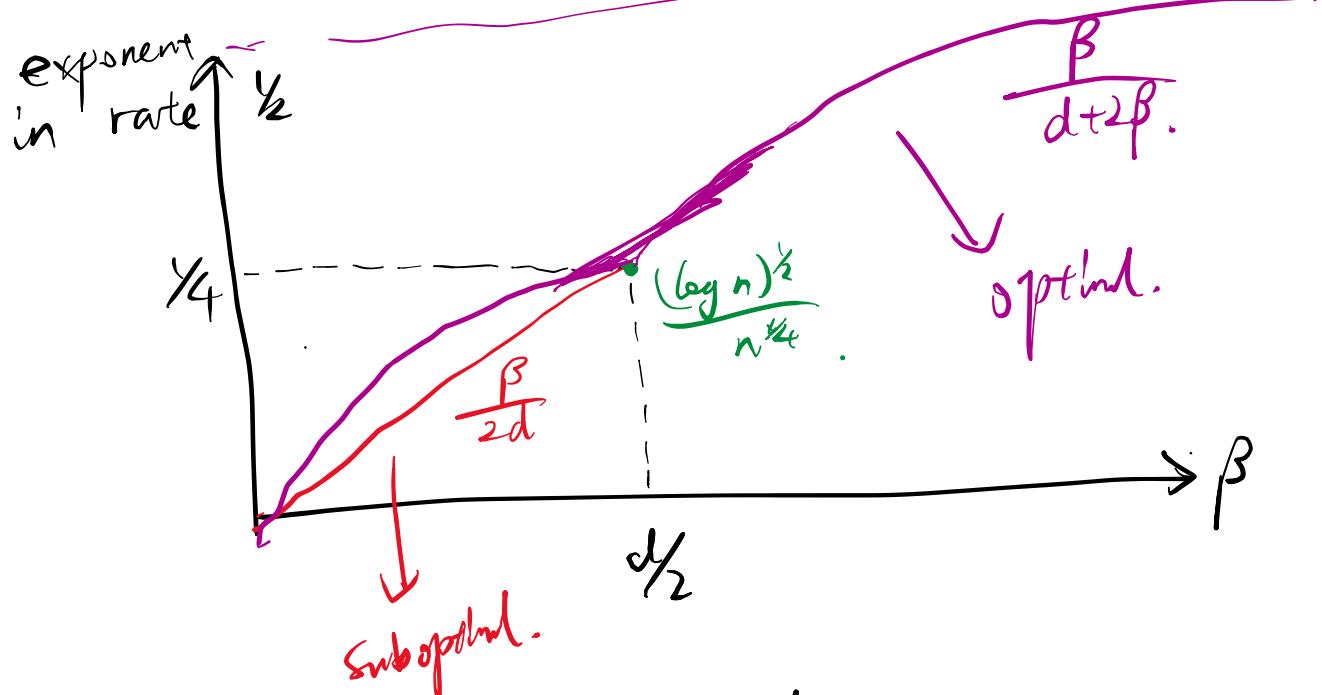
$$\leq r_0 + \frac{c}{\sqrt{n}} r_0^{1-\frac{d}{2\beta}}$$

Choose $r_0 = n^{-\beta/d}$ $G_n(r) \leq c \cdot n^{-\beta/d}$.

Implication to rate: $r_n = n^{-\frac{\beta}{2d}}$

- $\beta = d/2$, we can similarly get

$$r_n = (\log n)^{1/2} \cdot n^{-1/4}$$



- Suboptimality in $\beta \leq \frac{d}{2}$ regime

(optimal rate for Hölder class
is always $n^{-\frac{\beta}{d+2\beta}}$).

introduce for constrained least-square
(regularized).

- For Hölder class, projection/local poly leads to optimal rates.
- In general, divergent integral bounds
↳ known as "non-Donsker",
usually strong indicator of suboptimality
of least squares method.
(construction of optimal estimator
can be ad-hoc).

e.g. multivariate convex functions.

$$\mathcal{F} = \left\{ f: \bar{[0,1]^d} \rightarrow \bar{[0,1]}, |f(x) - f(y)| \leq \|x-y\|, \quad \begin{matrix} f \text{ is convex} \end{matrix} \right\}.$$

Optimal rate similar to $\sum(z, l)$

$$n^{-\frac{2}{d+4}} \quad \begin{matrix} \text{achieved by LS} \\ \text{when } d \leq 3. \end{matrix}$$

LS becomes suboptimal for $d \geq 4$.

Projection estimators for Sobolev classes.

$$f^* \in \mathcal{F} = S(\beta, L).$$

$$\beta \in N^+, \quad S(\beta, L) = \left\{ f : (\bar{0}, 1)^d \rightarrow [0, 1], \quad \sum_{|\alpha| \leq \beta} \int |\partial^\alpha f(x)|^2 dx < L^2 \right\}$$

Under Fourier transform.

Let $\varphi_1, \varphi_2, \dots, \varphi_j, \dots$ be Fourier basis in $[0, 1]^d$.

$$\text{d=1} \quad \sum_{j=1}^{+\infty} \langle f, \varphi_j \rangle_{L^2}^2 \cdot j^{2\beta} \leq L^2.$$

$$\text{In general d-dim, } \varphi_z \text{ for } z \in \mathbb{Z}^d.$$

$$\sum_{z \in \mathbb{Z}^d} \langle f, \varphi_z \rangle_{L^2}^2 \cdot |z|^{2\beta} \leq L^2.$$

Fixed design setting

$$\text{in 1-D, } x_i = i/n \quad (\text{equispaced})$$

In general dimensions,

$$x_d = \frac{\alpha}{m} \quad \text{where } \alpha = (\alpha_1, \alpha_2, \dots, \alpha_m) \in \mathbb{N}^d$$

$$m^d = n$$

$\varphi_1, \varphi_2, \dots, \varphi_n$ "Discrete Fourier basis".

Orthonormal basis of $(\mathbb{R}^n, \| \cdot \|_n)$.

$$\mathcal{F} = \left\{ f : \sum_{j=1}^n \varphi_j \cdot c_j \mid \sum_{j=1}^n c_j^2 \cdot j^{\frac{2\beta}{d}} \leq L^2 \right\}.$$

We know that constrained LS
is suboptimal for \mathcal{F} (when $\beta < \frac{d}{2}$).

Alternative idea: truncate in Fourier domain.

Observe $\gamma \in \mathbb{R}^n$

$$\begin{aligned} \hat{c}_j &:= \langle \gamma, \varphi_j \rangle_n \quad (\text{for } j=1, 2, \dots, n) \\ \hat{f}_n &:= \sum_{j=1}^N \hat{c}_j \varphi_j \quad (N < n). \end{aligned}$$

Analysys:

$$\mathbb{E}[\|\hat{f}_n - f^*\|_n^2] = \frac{1}{n} \sum_{j=1}^N \mathbb{E}[\|\hat{c}_j - c_j^*\|^2] + \frac{1}{n} \sum_{j=N+1}^n (c_j^*)^2$$

$$= \frac{N}{n} + \sum_{j=N+1}^n (c_j^*)^2 \quad \leq L^2$$

$$\leq \frac{N}{n} + \frac{1}{(N+1)^{\frac{2\beta}{d}}} \cdot \boxed{\sum_{j=N+1}^n j^{\frac{2\beta}{d}} \cdot (c_j^*)^2}$$

$$\leq \frac{N}{n} + \frac{L^2}{N^{2\beta/d}}$$

Choose $N = n^{\frac{d}{2\beta+d}}$ $\Rightarrow r_n = n^{-\frac{\beta}{2\beta+d}}$