

Final exam. Dec 19. 2-5 pm.

- Similar format to midterm
- No electronics
- Bring your ID.
- 4 pages (double-sided) cheat sheet.

Before the exam.

- Office hours (TA/ instructor).
- Practice question (released next week).

Nonparametric estimation.

$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} P$, learn something about P .

- Nonparametric model: P is not indexed by a finite-dimensional parameter.

e.g. cdf estimation.

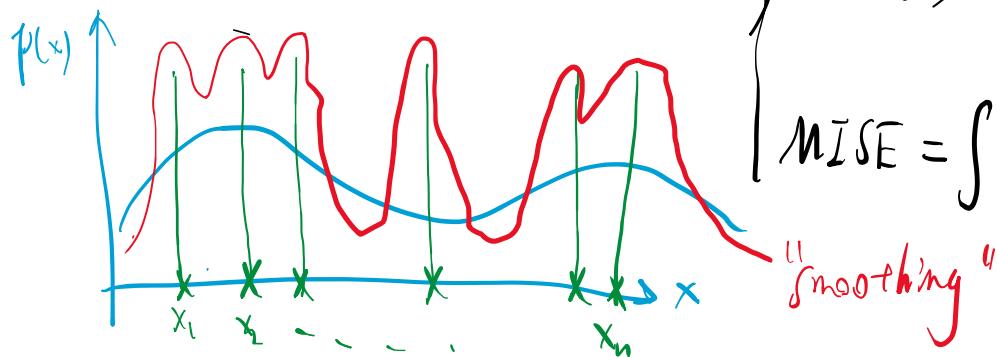
Today, we'll discuss

- (i) density estimation
- (ii) nonparametric regression.

(i) $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} p^*$ ($\in \mathcal{P}$, class of densities)
(p^* is the pdf)

Goal: recover p^* , more concretely, we want to

Empirical estimator won't work! minimize



$$\begin{cases} \text{MSE}(x_0) = \mathbb{E}[(\hat{p}(x_0) - p^*(x_0))^2] \\ \text{MISE} = \int \text{MSE}(x) dx. \end{cases}$$

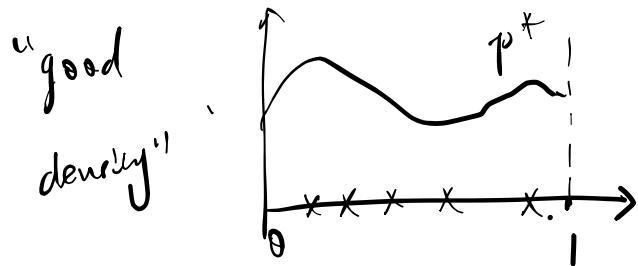
involving bias-variance trade-offs.

(ii) Nonparametric regression.

$$(X_i, Y_i)_{i=1}^n \quad \begin{cases} \text{Fixed design, } X_i's \text{ are deterministic} \\ Y_i = f^*(X_i) + \epsilon_i \\ \text{Random design, } (X_i, Y_i)_{i=1}^n \text{ follows a joint distribution.} \end{cases}$$

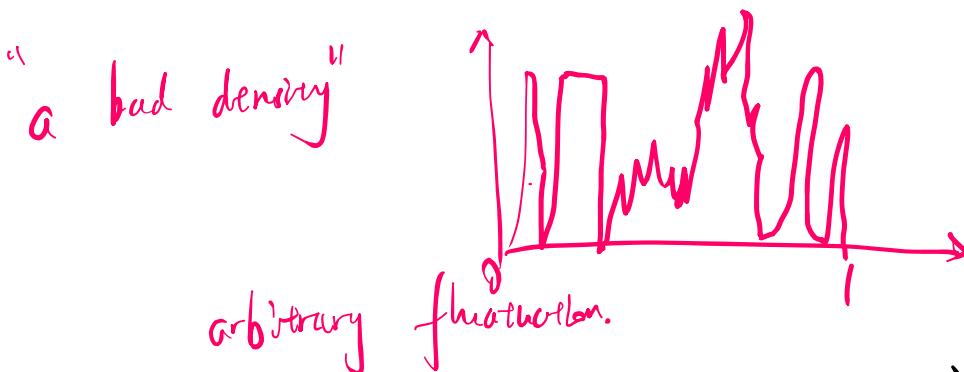
Density estimation.

$$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} p^* \text{ on } [0, 1].$$



Assumption on p^* .

$$|p^*(x) - p^*(y)| \leq L|x-y| \quad \forall x, y \in [0, 1].$$



(L is a constant characterizing modulus of continuity).

$$\frac{1}{nh} \cdot \sum_{i=1}^n \mathbf{1}_{\{X_i \in B_j\}}.$$

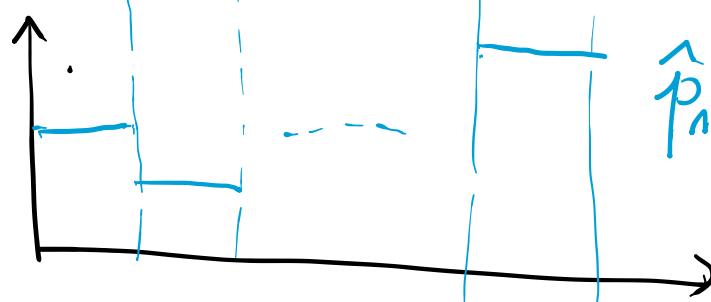
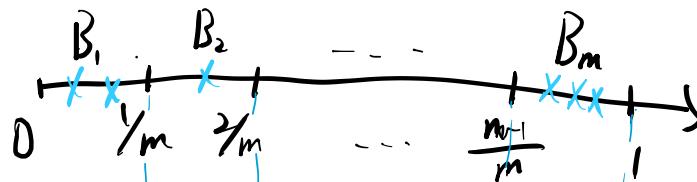
(assume L is fixed and known).

"Histogram estimator".

m bins

$$\hat{p}_n(x) = \frac{|B_j \cap \{x_1, \dots, x_n\}|}{n/m}$$

for $x \in B_j$.



$\hat{p}_n(x)$: piecewise constant function.

Idea: $\hat{P}_n(x) = \frac{\frac{1}{n} \cdot (\# \text{points in } B_j)}{\text{width of } B_j}$

$$p^*(x) = \lim_{h \rightarrow 0} \frac{P(X \in [x, x+h])}{h}$$

Remark: we can use this regardless of Lipschitz assumption.

But we expect it to work well only with

certain regularities on p^* (recall the bad example).

• Choice of m

— Larger m , width y_m smaller

This gives larger variance and smaller bias.

Theoretical analysis. (Under Lipschitz assumption).

Fact: $E[\hat{P}_n(x)] = \frac{P(X \in B_j)}{h}$. ($h := y_m$)

(for $x \in B_j$)

$$\text{Var}(\hat{P}_n(x)) = \frac{P(X \in B_j) \cdot (1 - P(X \in B_j))}{n h^2}$$

Upper bounds for bias & variance.

$$\cdot b(x) := E[\hat{P}_n(x)] - p^*(x).$$

$$= \frac{\mathbb{P}(x \in B_j)}{h} - p^*(x)$$

$$= \frac{1}{h} \int_{B_j} p^*(u) du - p^*(x).$$

$$= \frac{1}{h} \int_{B_j} (p^*(u) - p^*(x)) du.$$

$$|b(x)| \leq \frac{1}{h} \int_{B_j} |p^*(u) - p^*(x)| du$$

$$\leq \frac{L}{h} \int_{B_j} |x-u| du \leq \frac{L}{h} \int_{B_j} h du$$

$$\leq L \cdot h.$$

$$\cdot \sigma^2(x) = \text{var}(\hat{P}_n(x)) = \frac{\mathbb{P}(x \in B_j) \cdot (1 - \mathbb{P}(x \in B_j))}{nh^2}$$

$$\leq \frac{\mathbb{P}(x \in B_j)}{nh^2} \leq \frac{p_{\max}}{nh}$$

(Assuming $p^*(x) \leq p_{\max}$ & $x \in [0,1]$)

(Actually this is implied by Lipschitz assumption).

$$P(X \in B_j) = \int_{B_j} p^*(x) dx \leq \int_{B_j} P_{\max} dx = P_{\max} h.$$

- Bias-variance tradeoff

$$E[|\hat{p}_n(x_0) - p^*(x_0)|^2]$$

$$= b(x_0)^2 + \sigma^2(x_0).$$

$$\leq L^2 h^2 + \frac{P_{\max}}{nh}$$

(larger $n \Rightarrow$ smaller h)

\downarrow
smaller

\downarrow

larger
variance.

$$\text{Minimize the sum. optimal } h_n = \left(\frac{P_{\max}}{2L^2 n} \right)^{1/3}.$$

This gives the error bound

$$MSE(x_0) \leq 2 \cdot \left(\frac{P_{\max} L^2}{n} \right)^{2/3}.$$

- Unlike parametric estimators, the convergence rate

is of order $n^{-2/3}$ in MSE.

- Here we mainly focus on the rate w/o caring too much about the constant..

So we can choose h w/o knowing (L, P_{\max}) .

By choosing $h_n = Cn^{-1/3}$, we obtain

$$MSE(x_0) \leq C' \cdot n^{-2/3}$$

where C' is a constant depending on (c_0, p_{\max}, L) .

Integrating the bound w.r.t. x .

$$MISE = \int_0^1 MSE(x) dx \leq C' \cdot n^{-2/3}.$$

Both MSE and MISE bounds are information-theoretically (rate) optimal for Lipschitz densities.

p^* may be more regular than just Lipschitz.

$$\text{eg } \left| \frac{d^2 p^*(x)}{dx^2} \right| \leq L \quad (\forall x \in [0,1]).$$

(this is more than Lipschitz).

Better choice: kernel density estimator.

Given a kernel function $K: \mathbb{R} \rightarrow \mathbb{R}$.

satisfying

$$(i) \int K(x) dx = 1.$$

$$(ii) \int x K(x) dx = 0 \quad (\text{satisfied by even function } K).$$

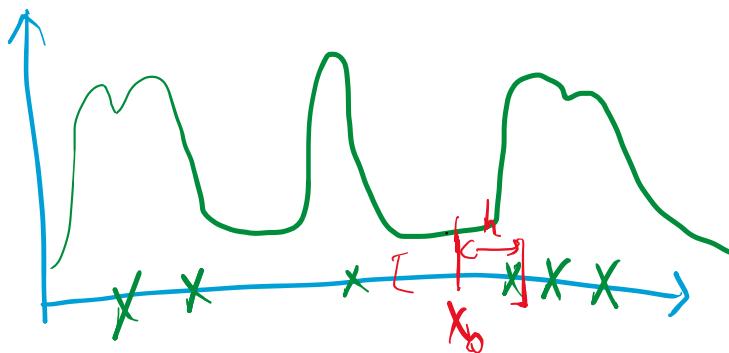
We estimate $p^*(x)$ by

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

h : bandwidth parameter.

$\frac{1}{nh}$ normalization

ensures $\int \hat{P}_n(x) dx = 1$.



Choice of K :

- For Lipschitz / second-order smooth p^* , choice of K does not matter too much.

You could simply use box-shaped K

$$K(x) = \frac{1}{2} \cdot \mathbb{1}_{x \in [-1, 1]}$$

bandwidth choice is more important.

(when taking box-shaped K , KDE becomes local averaging).

- Choice of K matters for (≥ 3) -rd order smoothness p^* , beyond the scope.

Analysis of KDE.

Bias $b(x) = \frac{1}{h} \int K\left(\frac{u-x}{h}\right) \cdot p^*(u) du - p^*(x)$

$$= \frac{1}{h} \int K\left(\frac{u-x}{h}\right) \cdot (p^*(u) - p^*(x)) du.$$

Variance $\sigma^2(x) = \frac{1}{nh^2} \cdot \text{var}\left(K\left(\frac{x-u}{h}\right)\right)$

(similar to histogram estimator).

$$\leq \frac{1}{nh^2} \cdot \int K^2\left(\frac{u-x}{h}\right) \cdot p^*(u) du.$$

(change of var: $y = \frac{u-x}{h}$).

$$= \frac{1}{nh} \int K^2(y) \cdot p^*(x+yh) dy.$$

$$\leq \frac{P_{\max}}{nh} \cdot \underbrace{\int K^2(y) dy}_{\text{Constant.}}$$

Closer look at bias.

$$b(x) = \frac{1}{h} \int K\left(\frac{u-x}{h}\right) (p^*(u) - p^*(x)) du$$

- If we assume p^* is L -Lipschitz,

$$|b(x)| \leq \frac{1}{h} \int K\left(\frac{u-x}{h}\right) |u-x| du.$$

(Assuming that K is supported on $[-1, 1]$
 i.e. $K(x)=0$ when $|x|>1$).

$$= \frac{L}{h} \int_{x-h}^{x+h} K\left(\frac{u-x}{h}\right) |u-x| du.$$

(Assuming $|K(x)| \leq K_{\max}$)

$$\leq \frac{LK_{\max}}{h} \int_{x-h}^{x+h} |u-x| du = LK_{\max} \cdot h.$$

Same as histogram method.

Assuming $|(\hat{p}^*)''(x)| \leq L$ $\forall x \in [0, 1]$.

KDE can do better.

$$b(x) = \int_{-1}^1 K(y) \cdot (p(x+yh) - p(x)) dy.$$

$$(u-x = yh).$$

$$p(x+yh) - p(x) = yh \cdot p'(x) + \frac{y^2 h^2}{2} p''(x + \tau_y yh)$$

where $\tau_y \in [0, 1]$ depends on y .

Cancels out after integration.

$$|b(x)| \leq \int_{-1}^1 \frac{|K(y)|y^2 h^2}{2} \cdot |p''(x + tyh)| dy \leq L \cdot h^2 K_{\max}$$

This is because

$$\begin{aligned} & \int_{-1}^1 K(y) \cdot y \cdot h p'(x) dy \\ &= \left(\int_{-1}^1 K(y) \cdot y dy \right) \cdot h p'(x) = 0. \end{aligned}$$

Bias-var trade off.

$$\begin{aligned} \mathbb{E}[|\hat{P}_n(x_0) - p^*(x_0)|^2] &= b(x_0)^2 + \sigma^2(x_0) \\ &\leq (L h^2 K_{\max})^2 + \frac{P_{\max}}{nh} \cdot \int K(x)^2 dx \\ &= C_1 \cdot h^4 + C_2 \cdot \frac{1}{nh}. \end{aligned}$$

Choose $h = C_0 \cdot n^{-1/5}$, we get

$$\mathbb{E}[|\hat{P}_n(x_0) - p^*(x_0)|^2] \leq C' \cdot n^{-4/5}.$$

Remark: faster rate available assuming more derivatives
But this requires clever choice of
(i) K (ii) h .

Brief intro to nonparametric regression.

Nadaraya-Watson estimator

Idea: "local averaging".

$$Y_i = f^*(x_i) + \varepsilon_i \quad (i=1, 2, \dots, n)$$

f^* is 1st/2nd order smooth.

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

(e.g. when K is box-shaped,
 $\hat{f}_n(x)$ is local averaging).

Achieves similar convergence rate.

- $n^{-2/3}$ for 1st order differentiable f^*

- $n^{-4/5}$ for 2nd order differentiable f^* .

Final review.

- cdf / empirical estimation / bootstrap.
 - Concepts. How these are implemented.
 - Theory. (no ε - δ proofs,
though ε - δ ideas are useful
for understanding).
 - eg. When bootstrap works.
(What it means for bootstrap to work).
- M-estimation, MLE, finite-dimensional asymptotics.
 - Consistency. $\hat{\theta}_n \xrightarrow{P} \theta^*$ } "big theorems"
 - Asymptotic normality.
(— Error bounds).
- Testing.
 - Concepts of type I/II.
 - Wald / χ^2 tests / permutation tests.
(— Testing radius).

- Decision theory / Bayes methods.
 - Concepts of minimax / Bayes estimators.
 - For calculations about Bayes estimator.
Calculation won't be the barrier.
(we'll provide formula for conjugate prior).
- Regression / nonparametrics.
 - (You'll get partial credit for solving $-D$).
 - Calculate concretely. (e.g. estimation error, inference etc.)
 - Nonparametric density estimation
(need to be comfortable w/
convergence rate analysis).