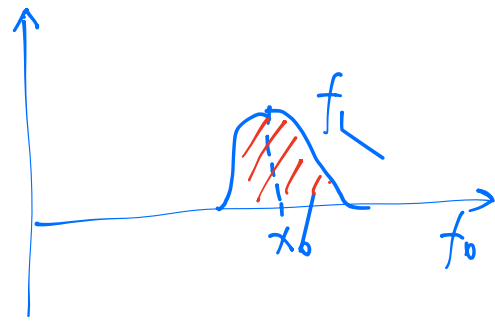


Question: minimax risk for estimating f

under integrated MSE

$$\inf_{\hat{f}} \sup_{f \in \mathcal{F}} \mathbb{E} \left[\int_0^1 |\hat{f}(x) - f^*(x)|^2 dx \right].$$

Recall construction from last time.



$$\int_0^1 |f_1(x) - f_0(x)|^2 dx$$

small.

Under Gaussian noise.

$$x_i = i/n,$$

$$D_{KL}(P_1^{(n)} \| P_0^{(n)}) = \frac{1}{2} \sum_{i=1}^n (f_1(i/n) - f_0(i/n))^2$$

$$\approx \frac{n}{2} \int_0^1 (f_1(x) - f_0(x))^2 dx.$$

If we want $D_{KL} \leq \frac{1}{2}$

$$\|f_1 - f_0\|_{L^2[0,1]} \lesssim \frac{1}{\sqrt{n}}.$$

Upper bound $\|\hat{f}_n - f^*\|_{L^2[0,1]} \lesssim n^{-\frac{\beta}{2\beta+1}}.$

Introduction to introduction to Info theory.

(C. Shannon).

• Compress information

• Send it over noiseless channel.

"Entropy"
 $H(X)$

$I(X;Y)$
"Mutual information"

X is discrete r.v.

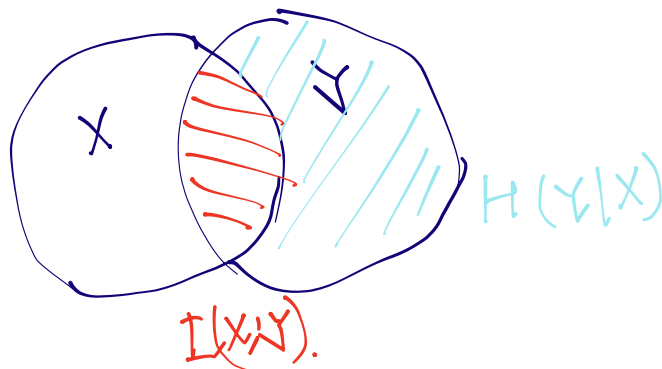
$$H(X) \left(\stackrel{\text{diff ent.}}{\approx} H(P_X) \right) \stackrel{\text{def}}{=} \sum_{x \in \mathcal{X}} \left(\log_2 \frac{1}{p(x)} \right) \cdot p(x).$$

On average, need $n \cdot H(X)$ bits to encode n iid copies of X .
 $\int p(x) \log_2 \frac{1}{p(x)} dx$

Joint entropy $H(X, Y)$. ($\leq H(X)$)

Conditional entropy $H(X|Y) = \sum_{y \in \mathcal{Y}} H(X|Y=y) \cdot P(Y=y)$

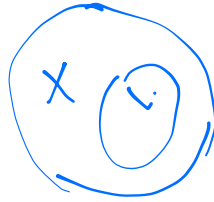
$$\begin{aligned} \text{Mutual Information} - I(X;Y) &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X). \end{aligned}$$



Independent: —



$Y = f(X)$: —



Shannon's theorem:

$$\text{Channel capacity} = \sup_X I(X; Y)$$

where $Y = \text{channel}(X)$.

In Stat's class

$$\theta \sim \pi, \quad X | \theta \sim P_\theta$$

$I(X; \theta)$ characterizes possibility of recovering θ from X .

Back to minimax lower bounds.

Discrete decision problem:

$$H_0: X \sim P_1 \quad \text{vs.} \quad H_2: X \sim P_2 \quad \dots \quad H_M: X \sim P_M$$

$$\text{Decision } T: X \rightarrow [M], \quad X \sim P_J, \quad J \sim \text{Unif}([M])$$

Quantity of interest.

$$\mathbb{P}(T(X) \neq \mathcal{J}).$$

Thm (Fano).

$$\mathbb{P}(T(X) \neq \mathcal{J}) \geq 1 - \frac{I(X; \mathcal{J}) + \log 2}{\log M}.$$

Proof of Fano.

• Step I. data processing inequality

$$\mathcal{J} \rightarrow X \rightarrow T(X).$$

$$\text{Then } I(X; \mathcal{J}) \geq I(T(X); \mathcal{J}).$$



$$I(X; \mathcal{J}) = I(X, T(X); \mathcal{J}) = I(T(X); \mathcal{J}) + I(X; \mathcal{J} | T(X))$$

• Step II. $f(t, j) = \mathbb{1}\{t \neq j\}$.

Detour

$$I(X; Y) = D_{KL}(P_{X, Y} \parallel P_X \times P_Y)$$

$$I(T(X); J) = D_{KL} \left(P_{T(X), J} \parallel P_{T(X)} \times \text{Uniform}(\bar{M}) \right).$$

$$D_{KL}(P_X \parallel P_Y) \geq D_{KL}(P_{f(X)} \parallel P_{f(Y)})$$

$$\geq D_{KL} \left(P_{f(T(X), J)} \parallel P_{f(T(X), J')} \right) \quad (J' \perp X)$$

$$= D_{KL} \left(\text{Ber}(P_e) \parallel \text{Ber}\left(1 - \frac{1}{M}\right) \right).$$

$$(P_e = \mathbb{P}(T(X) \neq J))$$

$$D_{KL} \left(\text{Ber}(P_e) \parallel \text{Ber}\left(1 - \frac{1}{M}\right) \right)$$

$$= P_e \log \frac{P_e}{1 - \frac{1}{M}} + (1 - P_e) \log \frac{1 - P_e}{1/M}.$$

$$\geq -\log 2 + \log M - P_e \log M.$$

$$\Rightarrow P_e \geq 1 - \frac{I(X; J) + \log 2}{\log M}.$$

Corollary. If $f_1, f_2, \dots, f_M \in \mathcal{F}$.

$$J \sim \text{Unif}(\bar{M})$$

If $\|f_i - f_j\| \geq 2\delta$ for each i, j pair

Then Bayes risk lower bound

$$\inf_{\hat{f}} \mathbb{E} \left[\|\hat{f} - f_J\|^2 \right] \geq \delta^2 \cdot \left(1 - \frac{\mathbb{I}(X; J) + \log 2}{\log M} \right)$$

Proof: reduce to testing.

Suppose we have estimator \hat{f} ,

Construct test $T = \operatorname{argmin}_{j \in [M]} \|\hat{f} - f_j\|$

$$\mathbb{P}(T \neq J) \leq \frac{1}{\delta^2} \mathbb{E} \left[\|\hat{f} - f_J\|^2 \right].$$

Construct $f_1, f_2, \dots, f_M \in \mathcal{F}$. and make sure.

- M large
 - $\mathbb{I}(X; J)$ bounded;
 - f_i, f_j separated from each other.
- make $\frac{1}{M} \sum_{j=1}^M D_{KL}(P_j \| Q)$ bounded.

$$\text{Bound } \mathbb{I}(X; J) = D_{KL} \left(P_{X, J} \parallel P_X \times P_J \right)$$

$J \sim \text{Unif}([M])$

$$= \frac{1}{M} \cdot \sum_{j=1}^M D_{KL} (P_j \parallel P)$$

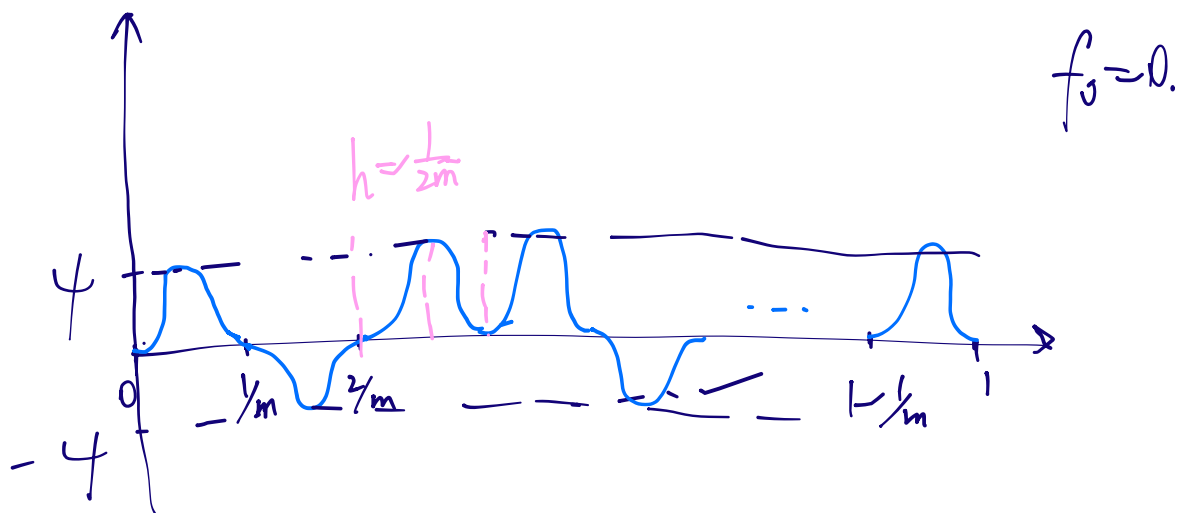
where $P_j = P(X|J=j)$.

$$\bar{P} = \frac{1}{M} \sum_{j=1}^M \bar{P}_j.$$

$$\frac{1}{M} \sum_{j=1}^M D_{KL}(P_j \| \bar{P}) \leq \frac{1}{M} \sum_{j=1}^M D_{KL}(P_j \| Q)$$

So we have

$$I(X;J) \leq \frac{1}{M} \sum_{j=1}^M D_{KL}(P_j \| Q).$$



For sign vector $z \in \{-1, 1\}^m$

$$\text{Define } f_z(x) = \psi K\left(\frac{x - \frac{z_i - 1}{2m}}{h}\right) z_i$$

$$\text{for } x \in \left[\frac{i-1}{m}, \frac{i}{m}\right].$$

$$D_{KL}(P_z \| P_0) = \frac{n}{2} \|f_z\|_n^2 \leq n \cdot \psi^2$$

For each z, z' pair

$$\|f_z - f_{z'}\|_{L^2(\mathcal{Q}, \mu)}^2 = c \cdot \psi^2 \cdot h \cdot \sum_{i=1}^m \mathbb{1}_{z_i = z'_i}$$

Sample $z^{(1)}, z^{(2)}, \dots, z^{(m)}$ iid $\text{Unif}(\mathcal{Z})^m$.

From Hw1.

$$\mathbb{P}\left(\exists i, j, \|z^{(i)} - z^{(j)}\|_1 \leq \frac{m}{8}\right) \leq \frac{1}{2}$$

Only for
conclusion.

$$\text{For } M = \exp(c'm) \quad (c' > 0).$$

seen as
deterministic

for the rest
of proof.

Use them as our prior distribution.

$$\forall i, j \quad \|f_{z^{(i)}} - f_{z^{(j)}}\|_{L^2}^2 \geq c \cdot \psi^2 \cdot h \cdot \frac{m}{8} \\ = \frac{c}{16} \psi^2.$$

$$\log M = c'm.$$

$$\frac{1}{M} \sum_{j=1}^M D_{\text{KL}}(P_j \| Q) \leq n \cdot \psi^2.$$

From last lecture.

for $f_z \in \Sigma(\beta, L)$.

$$\psi \leq C_1 h^\beta$$

$$\inf_{\hat{f}} \sup_{f^* \in \mathcal{F}} \mathbb{E} \left[\|\hat{f} - f^*\|_{L^2}^2 \right] \geq C \cdot \psi^2 \left\{ 1 - \underbrace{\frac{1}{m} \sum_{j=1}^m D_{\text{KL}}(P_j \| Q) + \log 2}_m \right\}$$

$$\geq C \cdot \psi^2 \left\{ 1 - \frac{n\psi^2 + \log 2}{m} \right\}$$

$$\psi \leq h^\beta = m^{-\beta}$$

Choose - $\psi = n^{-\frac{\beta}{2\beta+1}}$ (sup to const factors)

$$m = n^{\frac{1}{2\beta+1}}$$

So minimax rate $\geq n^{-\frac{2\beta}{2\beta+1}}$

Adaptivity.

Recall in nonpara regression

$$f^* \in \Sigma(\beta, L)$$

then

$$\text{MSE}(\hat{f}_n) \lesssim n^{-\frac{2\beta}{2\beta+1}}$$

Require knowledge about β .

Unknown β , hope to get optimal rate for underlying β .

$$\forall \beta \in [\beta_{\min}, \beta_{\max}] \quad (\beta_{\min} > 0)$$

Use a single estimator \hat{f}

$$\text{se. } \forall f^* \in \Sigma(\beta, 1)$$

$$\mathbb{E} \left[\left| \hat{f}(x_0) - f^*(x_0) \right|^2 \right] \leq C n^{-\frac{2\beta}{2\beta+1}}?$$

Surprisingly, this is not possible

(but we're not too far off).

Setup $Y_i = f^*(i/n) + \varepsilon_i$ $\varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$

$$f^* \in \Sigma(\beta, 1) \quad \text{with } \beta \in (0, 1].$$

Thm. $0 < \beta_1 < a < \beta_2 \leq 1$, $r_1^2(n) = \left(\frac{\log n}{n} \right)^{\frac{2\beta_1}{2\beta_1+1}}$

$$r_2^2(n) = n^{-\frac{2\alpha}{2\alpha+1}}$$

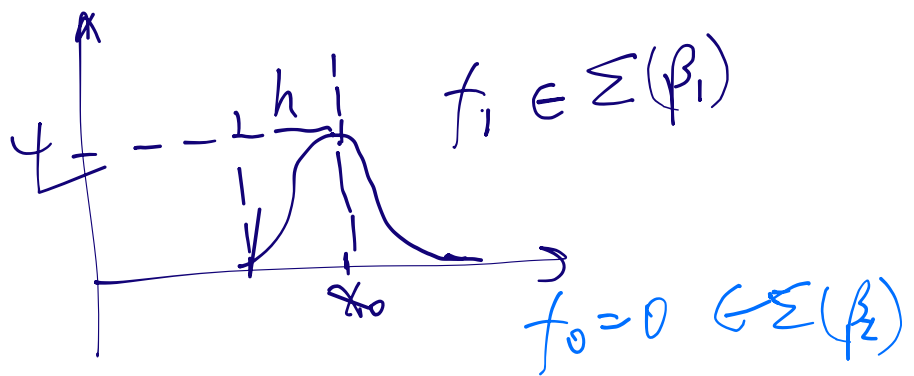
$$\inf_{\hat{f}} \sup_{i \in \{1,2\}} \sup_{f \in \Sigma(\beta_i)} \mathbb{E}_f \left[\frac{1}{r_i^2(n)} \left| \hat{f}_n(x_0) - f(x_0) \right|^2 \right] \geq C.$$

If we want 'idealized adaptivity,

need $r_i^2(n) = n^{-\frac{2\beta_i}{2\beta_i+1}}$ for $i \in \{1,2\}$

Proof sketch.

$$\psi \leq h^{\beta_1}.$$



"Asymmetric version of two-pt method!"

$$q_n = \frac{r_1^2(n)}{r_2^2(n)} \quad (= \text{poly}(n)).$$

$$R_{\text{adaptive-minimax}} \geq \frac{\psi_n^2}{c \cdot r_1(n)^2} \left(\mathbb{P}_1(\hat{T} = 0) + q_n \cdot \mathbb{P}_0(\hat{T} = 1) \right).$$

$$P_1(\hat{T} = 0) + q_n P_0(\hat{T} = 1)$$

$$\geq \int \min(P_1(x), q_n P_0(x)) dx$$

$$\geq 1 - \frac{\chi^2(P_1 \| P_0) + 1}{q_n}$$

$$1 + \chi^2(P_1 \| P_0) \leq \exp(2\psi_n^2 h_n \cdot n)$$

$$\text{Allow } \psi_n^2 \cdot h_n \cdot n \leq \log n.$$

Thm. \exists an estimator \hat{f}_{Lepski} s.t.

$$\mathbb{E} \left| \hat{f}_{\text{Lepski}}(x_0) - f(x_0) \right|^2$$

$$\sup_{\beta_{\min} \leq \beta \leq \beta_{\max}}$$

$$\sup_{f \in \Sigma(\beta, \nu)}$$

$$\frac{(\log n)^{\frac{2\beta}{2\beta+1}}}{n}$$

$$\leq C.$$

If β^* is the ground truth,

for $\beta < \beta^*$ (over-smooth)

$$\left| \hat{f}_{h_\beta}(x_0) - \hat{f}_{h_{\beta^*}}(x_0) \right| \leq C \cdot h_\beta^\beta \quad \text{for any } \beta < \beta^*.$$

$$\left(h_\beta = \left(\frac{\log n}{n} \right)^{\frac{1}{2\beta+1}}, \quad f_h \text{ local avg estimator w/ bandwidth } h \right).$$

Lepski's method.

$$\text{Step 1. } B = \{ \beta_{\min} = \beta_1 < \beta_2 < \dots < \beta_N = \beta_{\max} \}$$

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

$$N = O(\log n) \quad \text{many grid pts}$$

$$\text{Rate} \quad \beta \in [\beta_j, \beta_{j+1}]$$

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

Step 2.

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

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

Proof Sketch.

$$\Sigma_j = \{ \hat{\beta} = \beta_j \}$$

Ground truth β_{j^*} $f \in \Sigma(\beta_{j^*}, 1)$.

$$\mathbb{E} \left[|\hat{f}(x_0) - f(x_0)|^2 \right] = \sum_{j=1}^N \mathbb{E} \left[|\hat{f}(x_0) - f(x_0)|^2 \mathbb{1}_{\Sigma_j} \right].$$

Case I: $j \geq j^*$. Event Σ_j makes it safe to use β_j

Case II: $j < j^*$. β_{j^*} does not pass the test happens with small prob (log n appears here).