

Structure Learning



Lecture 2.

Preamble

Setting

Data: $(x_i, y_i)_{i=1 \dots m}$

Model: $y_i = f^*(x_i) + \varepsilon_i$ with
 $\rightarrow \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, 1) \perp (x_i)_{i=1 \dots m}$
 $\rightarrow f^* \in \mathcal{F} \leftarrow$ functional class

Empirical risk minimizer:

$$\hat{f} \in \underset{f \in \mathcal{F}}{\text{argmin}} \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2$$

Notation: $d_m(f, g)^2 = \frac{1}{m} \sum_{i=1}^m (f(x_i) - g(x_i))^2$

A simple bound

By definition, for any $f \in \mathcal{F}$:

$$\frac{1}{m} \sum_i (\overbrace{f^*(x_i)}^{y_i} + \varepsilon_i - \hat{f}(x_i))^2 \leq \frac{1}{m} \sum_i (f^*(x_i) + \varepsilon_i - f(x_i))^2$$

so

$$d_m(\hat{f}, f^*)^2 \leq d_m(\hat{f}, f)^2 + \frac{2}{m} \sum_{i=1}^m (\hat{f}(x_i) - f(x_i)) \varepsilon_i$$

So for $f = f^*$:

$$d_m(\hat{f}, f^*)^2 \leq \frac{2}{m} \sum_{i=1}^m (\hat{f}(x_i) - f^*(x_i)) \varepsilon_i$$

Assume for simplicity that

$$\{f - f^* : f \in \mathcal{F}\} \subset \{dg : \rightarrow d \geq 0 \rightarrow g \in \mathcal{G}\}$$

with $\frac{1}{m} \sum_{i=1}^m g(x_i)^2 = 1$ for all $g \in \mathcal{G}$.

Then, we have $\hat{f} - f^* = d_m(\hat{f}, f^*) \hat{g}$, with $\hat{g} \in \mathcal{G}$
and

$$d_m(\hat{f}, f^*)^2 \leq 2 d_m(\hat{f}, f) \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m g(x_i) \varepsilon_i$$

Hence

$$d_m(\hat{f}, f^*) \leq 2 \sup_{g \in \mathcal{G}} \underbrace{\frac{1}{m} \sum_{i=1}^m g(x_i) \varepsilon_i}_{\sim \mathcal{N}(0, 1/m)}$$

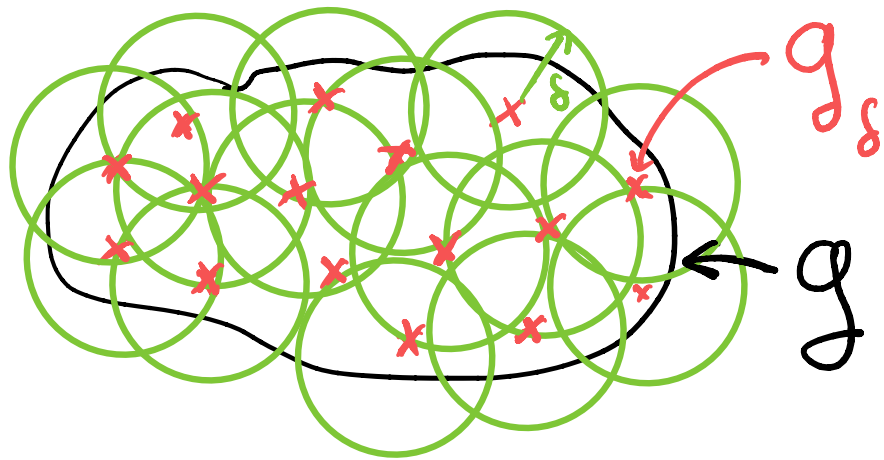
$=: \mathcal{L}(\mathcal{G})$

How large is $\mathcal{L}(\mathcal{G})$? problem: if \mathcal{G} is infinite, union bound is useless...

We can provide a simple upper-bound based on covering numbers.

Let $\delta > 0$, and let $\mathcal{G}_\delta \subset \mathcal{G}$ be a δ -covering of \mathcal{G} :

$$\forall g \in \mathcal{G}, \exists g' \in \mathcal{G}_\delta \text{ s.t. } d_m(g, g') \leq \delta$$



Then, we have:

$$\frac{1}{m} \sum_{i=1}^m g(x_i) \varepsilon_i = \frac{1}{m} \sum_{i=1}^m \underbrace{g'(x_i)}_{\in \mathcal{G}_\delta} \varepsilon_i + \frac{1}{m} \sum_{i=1}^m (g(x_i) - g'(x_i)) \varepsilon_i \leq d_m(g, g') \cdot \|\varepsilon\|_m$$

So choosing $g' \in \mathcal{G}_\delta$ s.t. $d_m(g, g') \leq \delta$ we get

$$\mathcal{Y}(g) \leq \delta \|\varepsilon\|_m + \sup_{g \in \mathcal{G}_\delta} \underbrace{\frac{1}{m} \sum_{i=1}^m g(x_i) \varepsilon_i}_{\sim \mathcal{N}(0, 1/m) \text{ conditionally to } x_i}$$

so

$$\mathbb{E}_\varepsilon [\mathcal{Y}(g)] \leq \delta + \sqrt{\frac{2}{m} \log \text{Card}(\mathcal{G}_\delta)}$$

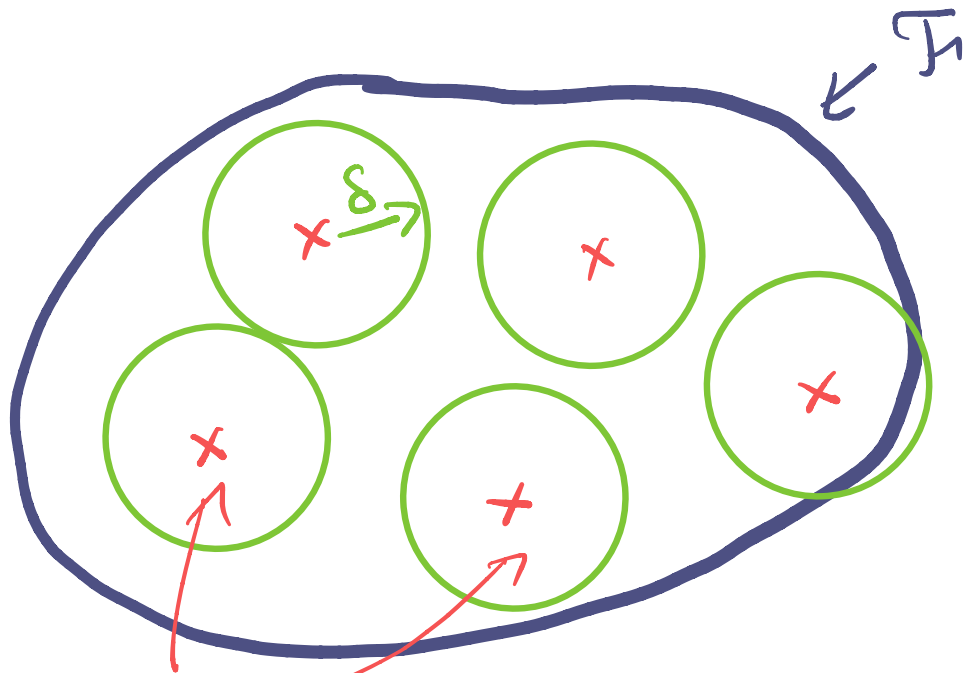
For example, if $\log \text{Card}(\mathcal{G}_\delta) \asymp 1/\delta^\alpha$:

$$\mathbb{E}_\varepsilon [d_m(\hat{g}, g^*)] \leq c \inf_{\delta > 0} \left(\delta + \frac{1}{\sqrt{m} \delta^\alpha} \right)$$

$$\delta = \frac{1}{m^{\frac{1}{2+\alpha}}} \Rightarrow \frac{c'}{m^{\frac{1}{2+\alpha}}}$$

\leadsto the massiveness of \mathcal{G} (value α) governs the convergence rate.

We can actually prove the following: Let $N(\mathcal{F}, d_m, \delta)$ be the maximum number of δ -separated points in (\mathcal{F}, d_m)



δ -separated points

Assume that

$$\log N(\mathcal{F}, d_m, \delta) \stackrel{[\delta \rightarrow 0]}{\asymp} \frac{1}{\delta^\alpha}$$

Then

$$\inf_{\hat{f}} \max_{f^* \in \mathcal{F}} \mathbb{E}_P [d_m(\hat{f}, f^*)^2] \geq \frac{c}{m^{2/2+\alpha}}$$

[exercise 3.6.4.]

Examples: for $\beta \in (0, 1]$, $\sigma^2 = 1$

① β -Hölder functions in \mathbb{R}^p

Set $\mathcal{F} = \{f: [0, 1]^p \rightarrow \mathbb{R} : |f(x) - f(y)| \leq |x - y|^\beta\}$.

Then

$$\inf_{\hat{f}} \max_{f^* \in \mathcal{F}} \mathbb{E}_P [d(\hat{f}, f^*)^2] \stackrel{m \rightarrow \infty}{\asymp} \frac{1}{m^{\frac{2}{2+p/\beta}}}$$

② Single index model

Set $\mathcal{F} = \left\{ f(x) = h(\langle x, w \rangle) : \begin{array}{l} -\|w\| = 1 \\ h: \mathbb{R} \rightarrow \mathbb{R} \\ \beta\text{-Hölder} \end{array} \right\}$. Then,

$$\inf_{\hat{f}} \sup_{f^* \in \mathcal{F}} \mathbb{E}_{f^*} [d_n(\hat{f}, f^*)^2] \stackrel{n \rightarrow \infty}{\asymp} \frac{1}{n^{2/2+1/\beta}}$$

Benefit of learning
the direction w

Learning structures
in
the linear model

The linear model is ubiquitous in data analysis. Predictions often take the form

$$\hat{f}(x) = \langle \hat{\beta}, \phi(x) \rangle$$

↑
features (learnt or
chosen)
or L^2 -basis, etc...

→ we focus on the linear model

$$f^*(x) = \langle \beta^*, x \rangle$$

$\beta^* \in \mathbb{R}^p$, unknown

• Model: $y_i = \langle \beta^*, x_i \rangle + \varepsilon_i$, $i=1, \dots, n$
with $\varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$

• Notation:
 $Y := \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$; $f^* := \begin{bmatrix} f^*(x_1) \\ \vdots \\ f^*(x_n) \end{bmatrix}$; $\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$

and $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times p}$

$$\Rightarrow Y = X\beta^* + \varepsilon = f^* + \varepsilon.$$

• Hidden structure: we assume that
 $|\beta^*|_0 := \text{card} \{j: \beta_j^* \neq 0\}$ is small
Coordinate sparse assumption.

How can we benefit from this assumption?

We set $S^* := \text{supp}(\beta^*) \leftarrow$ unknown

linear span $\rightarrow \bar{S}^* = \{X\beta: \text{supp}(\beta) \subset S^*\} \subset \mathbb{R}^n$

• If S^* was known: Then, we could

Solve $\hat{\beta}^{(S)} := \underset{\text{supp}(\beta) \subset S}{\text{argmin}} \|Y - X\beta\|^2$

for $S = S^*$. \leadsto back to low dimensional regression.

Then $\hat{f}^{(S^*)} = X \hat{\beta}^{(S^*)} = \text{Proj}_{\bar{S}^*} Y$

• When S^* is unknown:

• We can try to compare $\hat{\beta}^{(S)}$ for different guessed support $S \in \mathcal{S}$
 \uparrow
Collection of guessed supports

Reminder: $d_m(f, g)^2 = \frac{1}{m} \sum_{i=1}^m (f(x_i) - g(x_i))^2$

Risk:

$$\frac{1}{m} \|\text{Proj}_{\bar{S}}(f^* + \varepsilon) - f^*\|^2$$

$$R(\hat{f}^{(S)}) := \mathbb{E} [d_m(\hat{f}^{(S)}, f^*)^2]$$

Pythagore

$$= \frac{1}{m} \mathbb{E} [\|\text{Proj}_{\bar{S}} \varepsilon\|^2] + \frac{1}{m} \mathbb{E} [\|f^* - \text{Proj}_{\bar{S}} f^*\|^2]$$

$$= \underbrace{\frac{\sigma^2}{m} \text{Tr}(\text{Proj}_{\bar{S}})}_{= \text{dim}(\bar{S})} + \frac{1}{m} \|f^* - \text{Proj}_{\bar{S}} f^*\|^2$$


Variance term

bias term

Best S ? oracle choice

$$S_0 = \underset{S \in \mathcal{S}}{\text{argmin}} \left\{ \frac{1}{m} \|f^* - \text{Proj}_{\bar{S}} f^*\|^2 + \frac{\sigma^2}{m} \text{dim}(\bar{S}) \right\}$$

\uparrow
unknown!!

 1) estimate $R(\hat{f}^{(s)})$ by some $\hat{R}(\hat{f}^{(s)})$

2) Choose $\hat{f}(\hat{s})$ with $\hat{s} \in \operatorname{argmin}_{s \in \mathcal{S}} \hat{R}(\hat{f}^{(s)})$

Questions:

- which $\hat{R}(\hat{f}^{(s)})$?
- which performance?
(can we bypass the curse of dimensionality?)



We can try to estimate the bias $\|f^* - \operatorname{Proj}_{\mathcal{S}} f^*\|^2$ with $\|Y - \underbrace{\operatorname{Proj}_{\mathcal{S}} Y}_{= \hat{f}^{(s)}}\|^2$.

Let us analyse it!

$$\mathbb{E} \left[\|Y - \text{Proj}_{\bar{S}} Y\|^2 \right] \stackrel{y = f^* + \varepsilon}{=} \underbrace{\|f^* - \text{Proj}_{\bar{S}} f^*\|^2}_{\hat{f}^{(S)}} + \underbrace{\mathbb{E} \left[\|\varepsilon - \text{Proj}_{\bar{S}} \varepsilon\|^2 \right]}_{\downarrow} + 2 \underbrace{\mathbb{E} \left[\langle f^* - \text{Proj}_{\bar{S}} f^*, \varepsilon - \text{Proj}_{\bar{S}} \varepsilon \rangle \right]}_{=0}$$

$$= (n - \dim(\bar{S})) \sigma^2$$

So we have the unbiased estimation of the risk

$$\hat{R}_{\text{AIC}}(\hat{f}^{(S)}) := \frac{1}{n} \|Y - \hat{f}^{(S)}\|^2 + \underbrace{\frac{2}{n} \sigma^2 \dim(\bar{S})}_{\text{correction term}} \underbrace{(-\sigma^2)}_{\text{can be dropped for computing } \hat{S}}$$



Problem

$$\min_{S \in \mathcal{S}} \|Y - \text{Proj}_{\bar{S}} Y\|^2 + 2 \sigma^2 \dim(\bar{S}) \quad \text{can vary much}$$

$$\text{deviate from } \min_{S \in \mathcal{S}} \mathbb{E} \left[\|Y - \text{Proj}_{\bar{S}} Y\|^2 \right] + 2 \sigma^2 \dim(\bar{S})$$

when \mathcal{S} is large e.g. $\mathcal{S} = \{S : S \subset \{1, \dots, p\}\}$ -

Indeed,

$$\|Y - \text{Proj}_{\bar{S}} Y\|^2 = \underbrace{\|f^* - \text{Proj}_{\bar{S}} f^*\|^2}_{\text{what we want to estimate}} + \underbrace{\|\varepsilon\|^2}_{\text{does not depend on } S} - \underbrace{\|\text{Proj}_{\bar{S}} \varepsilon\|^2}_{\text{depends on } S} + \underbrace{\text{cross-product}}_{\text{let's forget it here}}$$

We have $\mathbb{E}[\|\text{Proj}_{\bar{S}} \varepsilon\|^2] = \dim(\bar{S}) \sigma^2$. Setting $d_S := \dim(\bar{S})$, we have

$$\mathbb{P}\left[\|\text{Proj}_{\bar{S}} \varepsilon\|^2 \leq \sigma^2 (\sqrt{d_S} + \sqrt{2L})^2\right] \geq 1 - e^{-L}$$

In particular, we have

$$\mathbb{P}\left[\min_{|S|=d} -\|\text{Proj}_{\bar{S}} \varepsilon\|^2 \geq -\sigma^2 (\sqrt{d} + \sqrt{2 \log C_P^d + 2L})\right] \geq 1 - e^{-L}$$

union bound

since $\log C_P^d \leq d \log\left(\frac{ep}{d}\right)$, we must correct $-\|\text{Proj}_{\bar{S}} \varepsilon\|^2$

not by $d_S \sigma^2$ but by $(\sqrt{d_S} + \sqrt{2d_S \log \frac{ep}{d_S}})^2 \sigma^2 \approx 2d_S \log\left(\frac{ep}{d_S}\right) \sigma^2$

$d_S \text{ small} \rightarrow \approx 2d_S \log(p) \sigma^2$

Conclusion: above analysis suggests to select

$$\hat{S} \in \operatorname{argmin}_S \|Y - \hat{f}^{(S)}\|^2 + \operatorname{pen}(S) \sigma^2$$

where $\operatorname{pen}(S) = K d_S \log\left(\frac{eP}{d_S}\right)$

for some $K \geq 2$.

Theorem: $\exists C_K, C'_K > 0$, such that

$$R(\hat{f}^{(\hat{S})}) \leq C_K \min_S \left\{ R(\hat{f}^{(S)}) + \frac{\sigma^2}{m} \operatorname{pen}(S) \right\}$$

$$\begin{aligned} \text{for } S=S^* \rightarrow & \leq C'_K \underbrace{\log\left(\frac{P}{d_{S^*}}\right)}_{\downarrow} \underbrace{R(\hat{f}^{(S^*)})}_{= d_{S^*} \frac{\sigma^2}{m}} \\ & = d_{S^*} \frac{\sigma^2}{m} \end{aligned}$$

Proof: See Theorem 2.2 \square

Can we avoid the $\log p$ factor?

Minimax estimation:

We can prove that

$$\min_{\hat{f}} \sup_{f \in \mathbb{R}^{n \times p}} \sup_{|S^*|=d} R(\hat{f}, f^*) \asymp \frac{\sigma^2}{m} d \log \frac{P}{d}$$

so the \log factor is unavoidable.

Yet, the risk of $\hat{f}^{(\hat{S})}$ is much smaller than the risk of vanilla least-square, if

$$|S^*| \log \frac{P}{|S^*|} \ll P \ (\leq m):$$

$$R(\hat{f}^{(\hat{S})}) \leq \frac{\sigma^2}{m} |S^*| \log \frac{P}{|S^*|} \ll \frac{\sigma^2}{m} P = R(\hat{f}^{LS})$$

\uparrow
if $\operatorname{rank}(X) = p$

\leadsto We have been able to take benefit of structure to get a better estimator



The complexity for computing \hat{S}
is $\geq 2^P$ in general

\rightsquigarrow prohibitive in practice!

Except when:

\rightarrow The columns of X are orthogonal
(simple thresholding: exercise 2.8.1)

$\rightarrow f(x)$ piecewise constant ($m=p$),
then complexity = $O(p^3)$ with dynamic
programming (exercise 2.8.4)

To be continued...