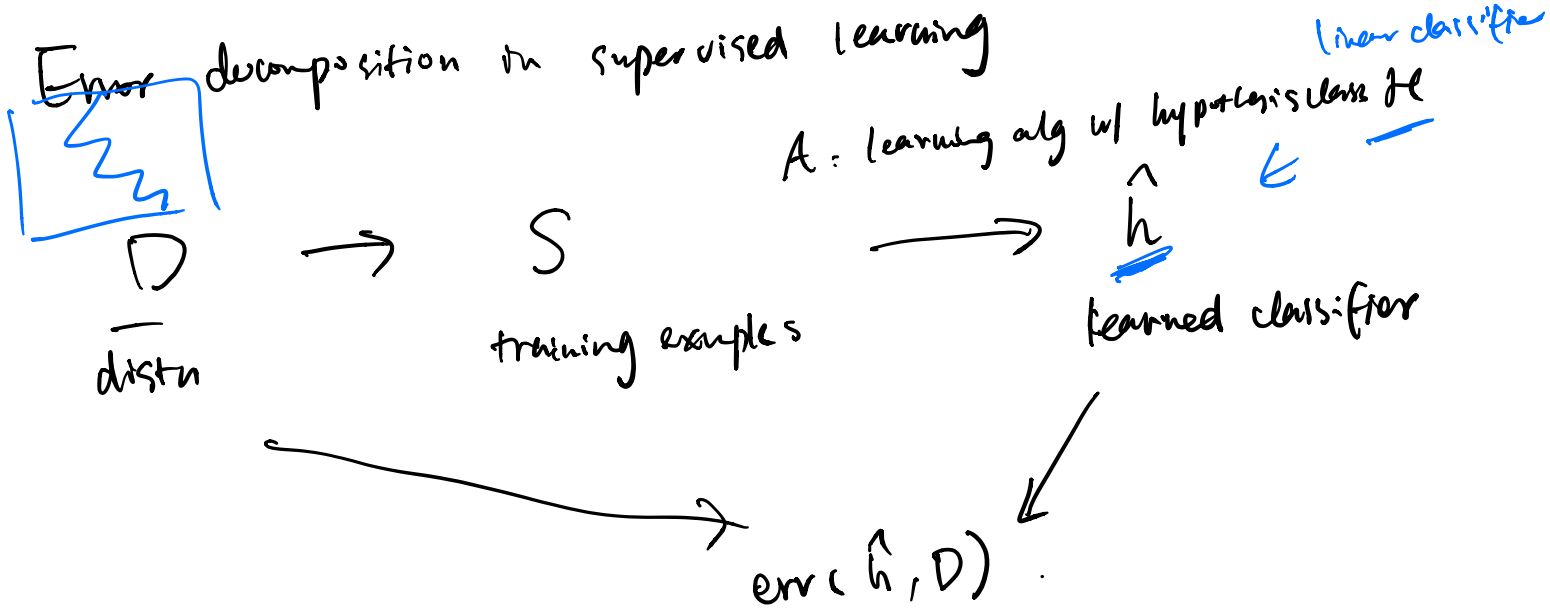


Model Selection

fixed hypothesis class \mathcal{H}

choose $h \in \mathcal{H}$ that are helpful for learning.



Q: what are some important factors that contribute to the generalization error of \hat{h} ?

- representativeness of training example
- complexity of \hat{h} , (\mathcal{H})
- optimization accuracy of A .
- expressiveness of \mathcal{H} relative to D .

Notation: $h' = \underset{h \in \mathcal{H}}{\text{argmin}} \text{err}(h, S)$

$h^* = \underset{h \in \mathcal{H}}{\text{argmin}} \text{err}(h^*, D)$

$\text{err}(\hat{h}, D)$

Thm: w.p. $1 - \delta$

$$\text{err}(\hat{h}, D) \leq \underbrace{\sqrt{\frac{\ln(1/\delta)}{2m}}}_{\text{bias}} + \underbrace{\Sigma_{\text{opt}}}_{\text{opt}} + \underbrace{\text{err}(h^*, D)}_{\text{bias}} + \sqrt{\frac{\ln(1/\delta)}{2m}}$$

$\hat{h} = h^*$ (ERM)

$\text{err}(\hat{h}, D) - \text{err}(\hat{h}, S)$ $\text{err}(\hat{h}, S) - \text{err}(h^*, S)$
 generalization gap of \hat{h} optimization error of \hat{h}

Pf:

$$\text{err}(\hat{h}, D) = \text{err}(\hat{h}, S) + \Sigma_{\text{gen}}$$

$$= \text{err}(h^*, S) + \Sigma_{\text{opt}} + \Sigma_{\text{gen}}$$

$$= \text{err}(h^*, S) + \Sigma_{\text{opt}} + \Sigma_{\text{gen}} + \underbrace{(\text{err}(\hat{h}, S) - \text{err}(h^*, S))}_{\leq 0}$$

w.p. $1 - \delta$. by Hoeffding within $\sqrt{\frac{\ln(1/\delta)}{2m}}$

$$\leq \text{err}(h^*, D) + \sqrt{\frac{\ln(1/\delta)}{2m}} + \Sigma_{\text{opt}} + \Sigma_{\text{gen}} + \underbrace{(\text{err}(\hat{h}, S) - \text{err}(h^*, S))}_{\leq 0}$$

$$\leq \text{err}(h^*, D) + \sqrt{\frac{\ln(1/\delta)}{2m}} + \Sigma_{\text{opt}} + \Sigma_{\text{gen}}$$

can be quite loose

Remark: $\text{err}(h^*, D)$ is called the bias of \mathcal{H} on D .

$\min_{h \in \mathcal{H}} \text{err}(h, D)$

2. when m is large. $\sqrt{\frac{\ln(1/\delta)}{2m}}$ can be ignored.

3. tightness of the above bound

$\text{err}(h^*, S) - \text{err}(h^1, S)$ can be quite

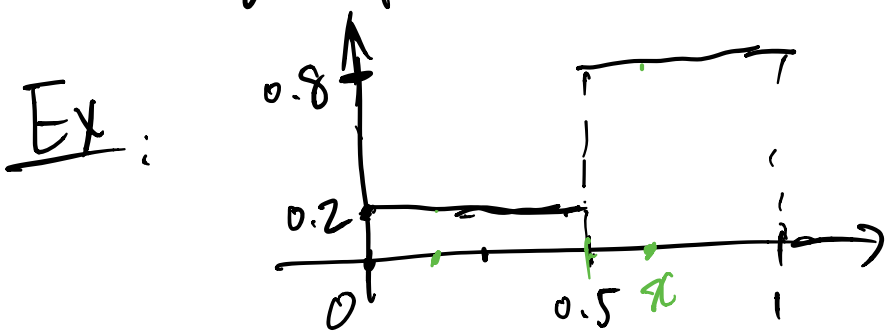
large. In this case,

$$\text{err}(h^1, D) \leq \underbrace{\Sigma_{\text{gen}} + \Sigma_{\text{opt}} + \sqrt{\frac{1}{m}}}_{\text{large}} + (\text{err}(h^1, S) - \text{err}(h^*, S)) + \underbrace{\text{err}(h^*, D)}$$

$$\underbrace{\text{err}(h^*, S) - \text{err}(h^1, S)}_{\text{large}} \leq \underbrace{\Sigma_{\text{gen}} + \Sigma_{\text{opt}}}_{\text{large}} + \underbrace{\sqrt{\frac{1}{m}}}_{\text{small}}$$

\Rightarrow at least one of Σ_{opt} and Σ_{gen} would be

large. $P(Y=1|X)$

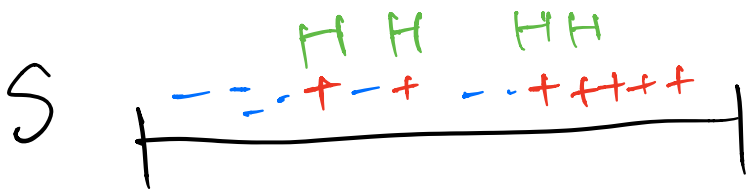


$$P_X(x) = \begin{cases} 1 & x \in [0, 1] \\ 0 & \text{o.w.} \end{cases}$$

optimal classification decision:
 $\begin{cases} x \geq 0.5 & +1 \\ x < 0.5 & -1 \end{cases}$

$$\mathcal{H} = \left\{ \sum_{i=1}^k I(x \in [a_i, b_i]) : k \in \mathbb{N}, a_i, b_i \in \mathbb{R} \right\}$$

• h^*



$$\text{err}(h', S) = \min_{h \in \mathcal{H}} \text{err}(h, S) = 0.$$

$$\text{err}(h^*, D) = \min_{h \in \mathcal{H}} \text{err}(h, D) = 0.2$$

$$h^* = 2 \cdot \mathbb{I}(x \in [0.5, 1]) - 1$$

Hoeffding

$$\text{err}(h^*, S) \geq 0.2 - \sqrt{\frac{1}{m}}$$

$\Rightarrow \text{err}(h^*, S) - \text{err}(h', S)$ is large.

Bringing down $\hat{\Sigma}_{\text{opt}}, \hat{\Sigma}_{\text{gen}}, \text{err}(h^*, D)$?

$\hat{\Sigma}_{\text{opt}} \downarrow$: change the ML optimization alg
make it simple to optimize

$\hat{\Sigma}_{\text{gen}} \downarrow$: choose a less expressive \mathcal{H} .

collect more sample

$\text{err}(h^*, D) \downarrow$: choose a more expressive \mathcal{H} .

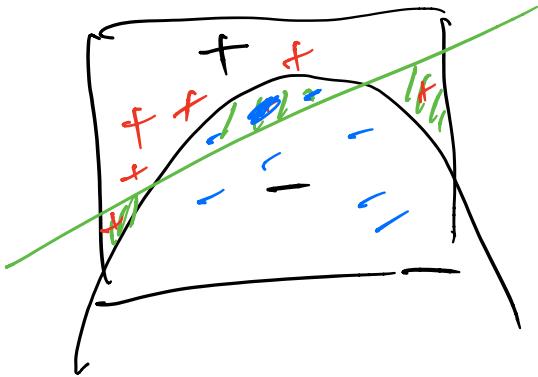
Important special case, $A \equiv \text{ERM}(\mathcal{H})$.

$\ln(\mathcal{H}) \rightarrow VC(\mathcal{H})$

$$\text{err}(\hat{h}, D) \leq \underbrace{\text{err}(h^*, D)}_{\text{bias}} + 2 \sqrt{\frac{\ln \frac{2|H|}{\delta}}{2m}} \underbrace{\quad}_{\text{complexity of } H}$$

bias-complexity tradeoff.

① underfitting: bias too large



$H = \{ \text{linear classifiers} \}$

some lines can be caught by seeing

$\text{err}(\hat{h}, S)$ is too large \Rightarrow $\text{err}(h^*, S)$ is also large
 \Downarrow $\text{err}(h^*, D)$ is large \Leftarrow $\text{err}(h^*, D)$

② overfitting: H is too large so that complexity term is too large

$\text{err}(\hat{h}, D) - \text{err}(\hat{h}, S)$ is large.



$\text{err}(\hat{h}, S) \Rightarrow 0$

Use fresh validation set V

$$\text{err}(\hat{h}, D) \approx \text{err}(\hat{h}, V) \quad (\text{Hoeffding})$$

$$\Sigma_{\text{gen}} \approx \text{err}(\hat{h}, V) - \text{err}(\hat{h}, S)$$

Model Selection

- How can we choose a good learning alg in practice?
- only considers EDM over hypothesis classes.

Setup: $\mathcal{H}_1, \dots, \mathcal{H}_k$
 $= \{ \text{depth} \leq k \text{ decision trees} \} = \{ \text{depth} \leq k \text{ decision tree} \}$

$$h_i^* = \underset{h \in \mathcal{H}_i}{\text{argmin}} \text{err}(h, D)$$

\hat{h}_k may not be the best among $\hat{h}_1, \dots, \hat{h}_k$

$$\hat{h}_i = \underset{h \in \mathcal{H}_i}{\text{argmin}} \text{err}(h, S)$$

How to use $\mathcal{H}_1, \dots, \mathcal{H}_k$ to find a good h w/ low error?

$$\hat{h} = \underset{h \in \bigcup_i \mathcal{H}_i}{\text{argmin}} \text{err}(h, S) \quad ?$$

Ideal: validation:

$$\hat{\mathcal{H}} = \{ \hat{h}_1, \dots, \hat{h}_k \}$$

$\hat{h} = \underset{h \in \hat{H}}{\operatorname{argmin}} \operatorname{err}(h, V)$ where V is
 a fresh validation sample.

Analysis

① w.p. $1 - \delta/2$: $\forall i$.

$$\operatorname{err}(\hat{h}_i, D) \leq \operatorname{err}(h_i^*, D) + 2 \sqrt{\frac{|\mathcal{H}_i| \ln 4/\delta}{2m}}$$

(standard ERM analysis + union bound over all i)

② w.p. $1 - \delta/2$.

$$\operatorname{err}(\hat{h}, D) \leq \min_i \operatorname{err}(\hat{h}_i, D) + 2 \sqrt{\frac{\ln 4/\delta}{|V|}}$$

① ② \Rightarrow w.p. $1 - \delta$:

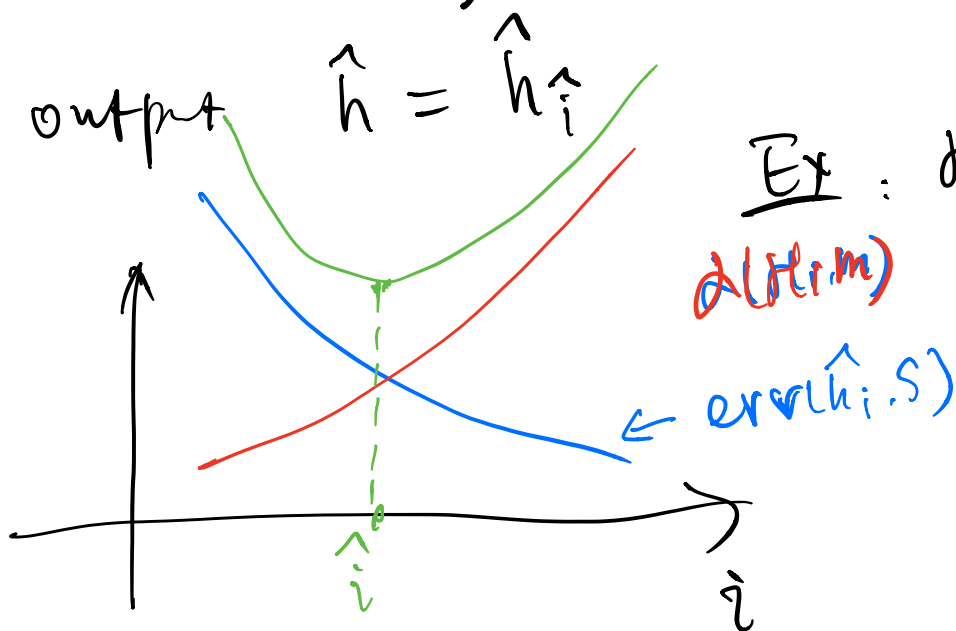
$$\operatorname{err}(\hat{h}, D) \leq \min_i \left(\operatorname{err}(h_i^*, D) + 2 \sqrt{\frac{|\mathcal{H}_i| \ln 4/\delta}{2m}} + 2 \sqrt{\frac{\ln 4/\delta}{|V|}} \right)$$

$|V| = \Theta(m)$ $\operatorname{opt}_{\mathcal{H}}(f(x)) \approx \operatorname{opt}_{\mathcal{H}}[c_1 f(x), c_2 f(x)]$ const c_1, c_2 .

$\Rightarrow \hat{h}$ has the best bias complexity tradeoff.

Idea 2: structural risk minimization
(penalized ERM)

$$\hat{i} = \operatorname{argmin}_{i \in \{1, \dots, k\}} \operatorname{err}(\hat{h}_i, S) + \sqrt{\frac{\ln \frac{2k}{\delta} |H_i|}{2m}}$$



$$\operatorname{err}(\hat{h}, D) \leq \min_{i \in \{1, \dots, k\}} \operatorname{err}(h_i^*, D) + 4 \cdot \alpha(H_i, m)$$

