

Model selection

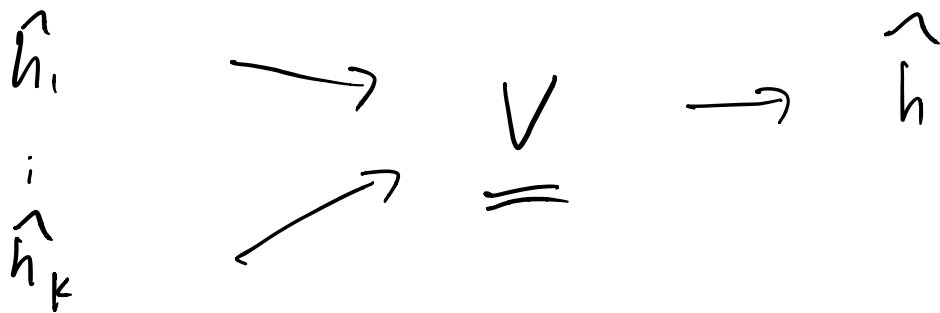
H_1, \dots, H_k

\hat{h}

whp.

$$\text{err}(\hat{h}, D) \leq \min_i \left(\text{err}(h_i^*, D) + O\left(\sqrt{\frac{\ln(|H_i|)}{m}}\right) \right)$$

1. Validation

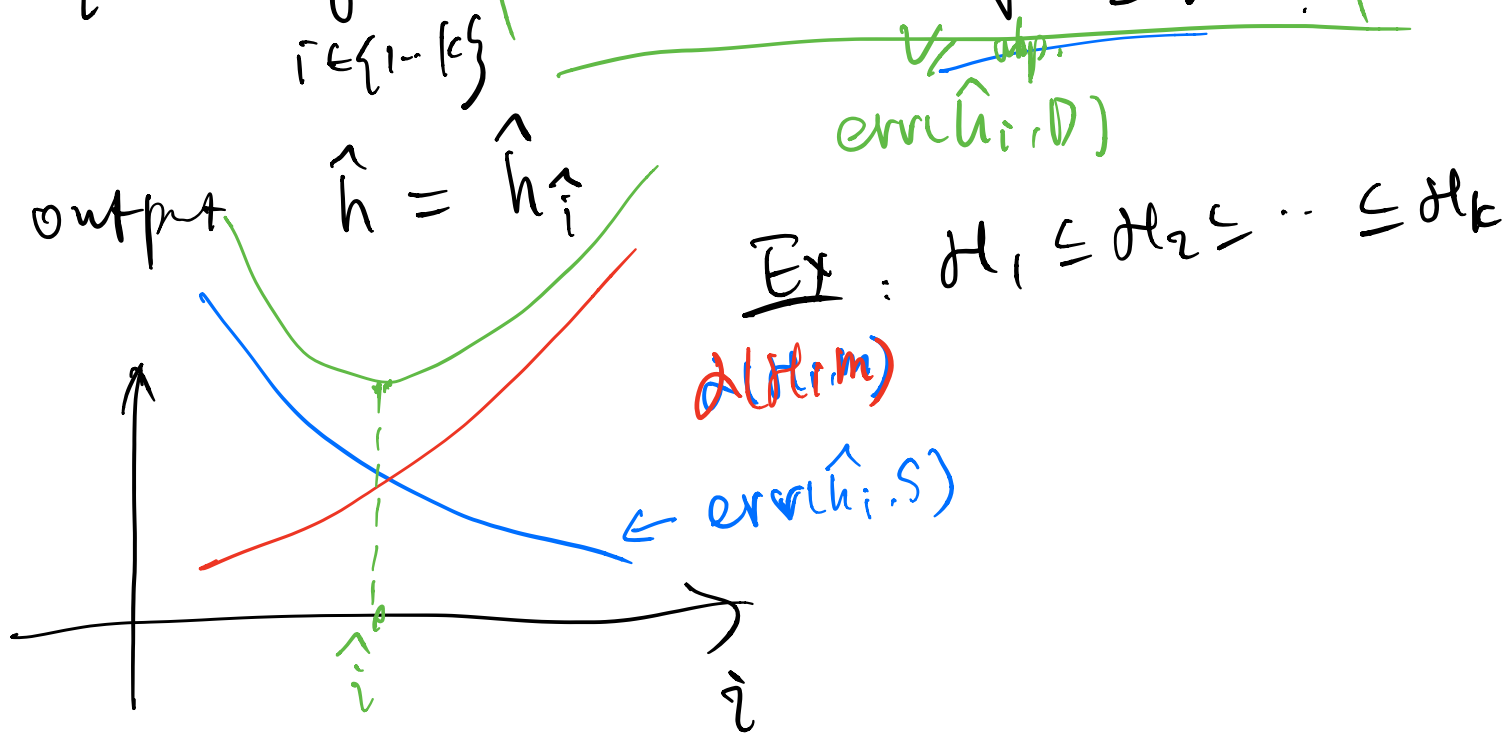


Idea 2: Structural risk minimization

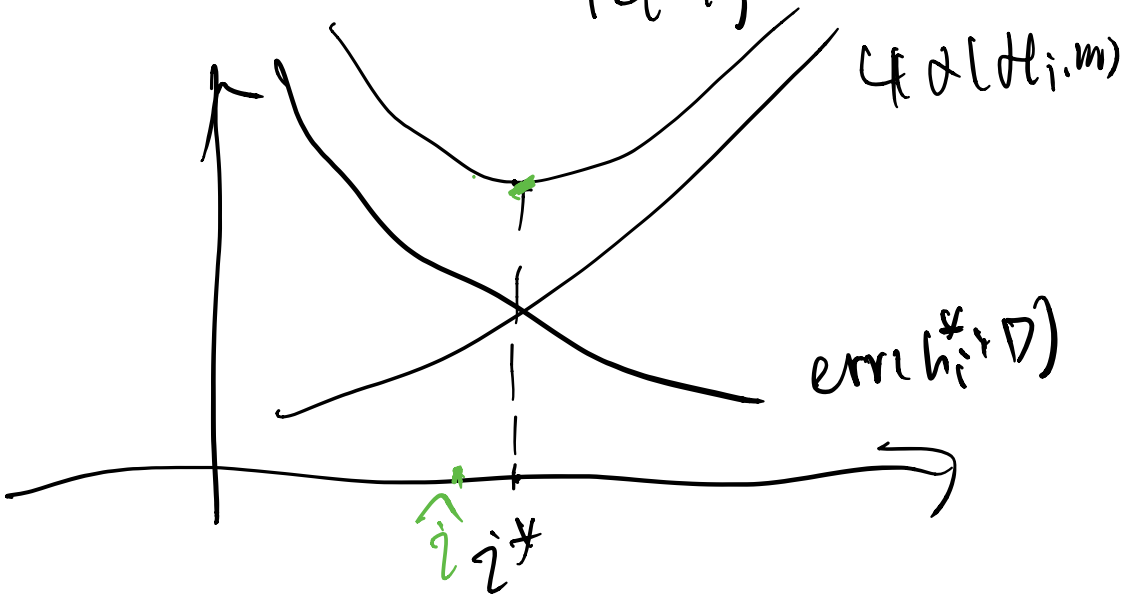
(penalized ERM)

$$\hat{h} = \underset{h}{\text{argmin}} \left(\text{err}(\hat{h}_i, S) + \sqrt{\frac{\ln \frac{2k|H_i|}{\delta}}{2m}} \right)$$

(Note: $\text{err}(\hat{h}_i, D)$ is crossed out in red, and $2(H_i, m)$ is written in blue above the fraction.)



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



oracle inequality

$$i^* = \underset{i}{\text{argmin}} (\text{err}(h_i^*, D) + 4 \alpha(H_i, m))$$

$\text{err}(\hat{h}, D)$ competitive w/ $\text{err}(h_{i^*}, D)$ w/o

knowing i^* .

Pf of (*)

$$\hat{h} = \underset{i \in \{1, \dots, k\}}{\operatorname{argmin}} \min_{h \in \mathcal{H}_i} \left(\operatorname{err}(h, S) + \sqrt{\frac{\ln \frac{|\mathcal{H}_i| k}{\delta}}{2m}} \right)$$

\forall for $h \in \mathcal{H}_i$ w.p.
 $\operatorname{err}(h, D)$

Step 1: w.p. $1 - \delta$: for all i , for all $h \in \mathcal{H}_i$:

$$\left| \operatorname{err}(h, S) - \operatorname{err}(h, D) \right| \leq \sqrt{\frac{\ln \frac{|\mathcal{H}_i| k}{\delta}}{2m}} \quad (\Delta_i)$$

fix \mathcal{H}_i . w.p. $1 - \delta/k$. (Δ_i) true

$$P\left(\bigcap_i (\Delta_i)\right) \geq 1 - \delta.$$

non-uniform
concentration

given this happens:

$$\operatorname{err}(\hat{h}, D) = \operatorname{err}(\hat{h}_i^*, D)$$

$$\leq \operatorname{err}(\hat{h}_i^*, S) + d(\mathcal{H}_i^*, m)$$

for all i

$$\leq \operatorname{err}(\hat{h}_i, S) + d(\mathcal{H}_i, m)$$

concentration within \mathcal{H}_i

$$\leq \operatorname{err}(\hat{h}_i, D) + 2d(\mathcal{H}_i, m)$$

ERM analysis

$$\leq \operatorname{err}(\hat{h}_i^*, D) + 4d(\mathcal{H}_i, m)$$

✱

Remark: 1. $\mathcal{L}(H, m)$ can be pessimistic in practice.

2. It may be possible to use more refined generalization error bounds.

Boosting

Motivation: combine weak classification rules to obtain strong ones.

Ex spam filtering

$x, y \rightarrow$ spam / no spam
text description

"free offer" \rightarrow spam

"a million \$"

Weak PAC learning: H, A is a δ -weak PAC learner for H , if \exists a fn $f: (0,1) \rightarrow \mathbb{N}$, for any $\delta > 0$, for any D realizable by \mathcal{R} , w/ $m \geq f(\delta)$, A produces a classifier h , s.t. w.p. $1-\delta$:

$$\text{err}(h, D) \leq \frac{1}{2} - \delta$$

\mathcal{H} is said to be δ -weak-DAC learnable if there is a δ -weak-DAC learner for \mathcal{H} .

\mathcal{H} PAC learnable \Rightarrow weak PAC learnable
 \nLeftarrow

history:

1988: Kearns asked this question

1990: Schapire: boosting:

builds a PAC learner w/ black-box access to a weak PAC learner. based on recursion.

1990: Freund, boost by majority

combining the outputs of weak learners by a weighted majority vote.

1997: Freund and Schapire: AdaBoost (no need to know δ)

Since: empirical success: X & boost

AdaBoost:

given training examples: $(x_1, y_1) \dots (x_n, y_n)$

Maintain a weighting on them, adjust weights. call weak learner.

Alg: Initialize $(D_1(i) = \frac{1}{m})_{i=1}^m$.

For $t=1, 2, \dots, T$:

$h_t \leftarrow B$ trained on weighted examples $(x_i, y_i, D_t(i))_{i=1}^m$.

weighted error: $\epsilon_t = P_{(x,y) \sim D_t} (h_t(x) \neq y)$

$$= \sum_{i=1}^m I(h_t(x_i) \neq y_i) D_t(i).$$

$$\leq \frac{1}{2} - \gamma.$$

classifier wt:

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$$

$$D_1(i) \propto 1$$
$$D_2(i) \propto e^{-\alpha_1 y_i h_1(x_i)}$$
$$D_3(i) \propto e^{-\alpha_1 y_i h_1(x_i) - \alpha_2 y_i h_2(x_i)}$$

update wt on training examples.

$$D_{t+1}(i) = D_t(i) \cdot e^{-\alpha_t y_i h_t(x_i)} / Z_t$$

up weighting incorrect examples

normalization factor.

Final classifier:

$$H_T(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right).$$

Analysis:

Thm: suppose for every t : $\sum_t \leq \frac{1}{2} - \gamma$,

then, $\text{err}(H_T, S) \leq \exp(-2T \cdot \gamma^2)$

pf: ① relate to exponential loss

② AdaBoost optimizes exponential loss.

①: $\text{err}(H_T, S) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(H_T(x_i) \neq y_i)$

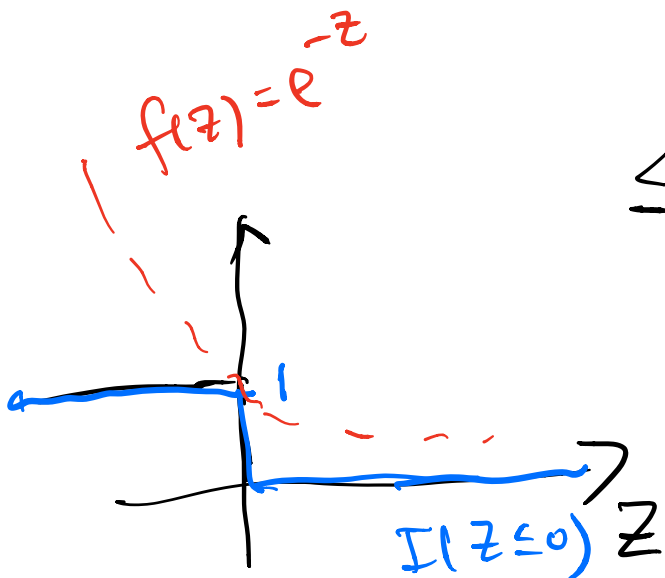
$$f_t = \sum_{s=1}^t d_s h_s(x)$$

\downarrow
 z

$$= \frac{1}{m} \sum_{i=1}^m \mathbb{I}(y_i \cdot f_T(x_i) \leq 0)$$

$$\leq \frac{1}{m} \sum_{i=1}^m e^{-y_i \cdot f_T(x_i)}$$

$= L_T$ (cumulative exponential loss)



Show L_t is decreasing exponentially in T .

$$\frac{L_t}{L_{t-1}} = \frac{\frac{1}{m} \sum_{i=1}^m e^{-y_i f_t(x_i) + d_t h_t(x_i)}}{\frac{1}{m} \sum_{i=1}^m e^{-y_i f_{t-1}(x_i)}} \quad (*)$$

Observation 1 :

$$D_t(i) \propto e^{-\sum_{s=1}^{t-1} d_s y_i h_s(x_i) - y_i f_{t-1}(x_i)} = e$$

$\Rightarrow \exists N_t$.

$$D_t(i) = e^{-y_i f_{t-1}(x_i)} / N_t$$

plugging into (*)

$$= \frac{\sum_{i=1}^m N_t D_t(i) e^{y_i d_t h_t(x_i)}}{\sum_{i=1}^m N_t \cdot D_t(i)} = N_t$$

$$= Z_t$$

$$\frac{L_t}{L_{t-1}} = z_t$$

$$L_T = L_{T-1} \cdot z_T$$

$$= L_{T-2} \cdot z_{T-1} \cdot z_T$$

$$= \dots$$

$$= \underline{L_0} \cdot \prod_{t=1}^T z_t$$

upper bounding z_t :

$$z_t = \sum_{i=1}^m D_t(i) \cdot e^{-\alpha_t y_i h_t(x_i)}$$

$$= \sum_{i: y_i = h_t(x_i)} D_t(i) \cdot e^{-\alpha_t} + \sum_{i: y_i \neq h_t(x_i)} D_t(i) \cdot e^{\alpha_t}$$

$$= e^{-\alpha_t} (1 - \zeta_t) + e^{\alpha_t} \cdot \zeta_t$$

$$d_z = \frac{1}{2} \ln \frac{1-\xi_t}{\xi_t}$$

$$= \sqrt{\frac{\xi_t}{(1-\xi_t)}} \cdot (1-\xi_t) + \sqrt{\frac{(1-\xi_t)}{\xi_t}} \cdot \xi_t$$

$$= 2 \cdot \sqrt{\xi_t (1-\xi_t)}$$

$$\xi_t \leq \frac{1}{2} - \delta.$$

$$\leq \sqrt{1 - 4\delta^2} \leq e^{-4\delta^2}$$

$$\leq e^{-2\delta^2}$$

$$\Rightarrow L_T \leq \prod_{t=1}^T z_t \leq e^{-2T\delta^2} \quad \star$$