

Recap

Fix margin $\rho > 0$.

Want to show: for every $h \in \mathcal{H}$, for every \mathcal{D} , for every $\delta > 0$,

$$R(h) \leq \hat{R}_{S,\rho}(h) + \frac{2}{\rho} \text{Rad}_S(\mathcal{H}) + 3 \sqrt{\frac{\log(4/\delta)}{2m}},$$

with probability $> 1 - \delta$.

$$\text{Take } \Phi(S) = \sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h))$$

Proof:

Recall:

$$R(h) = \mathbb{E}_{\substack{z \sim \mathcal{D} \\ (x,y)}} \mathbb{1}_{\{y \neq h(x)\}}$$

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{y_i \cdot h(x_i) < 0\}}$$

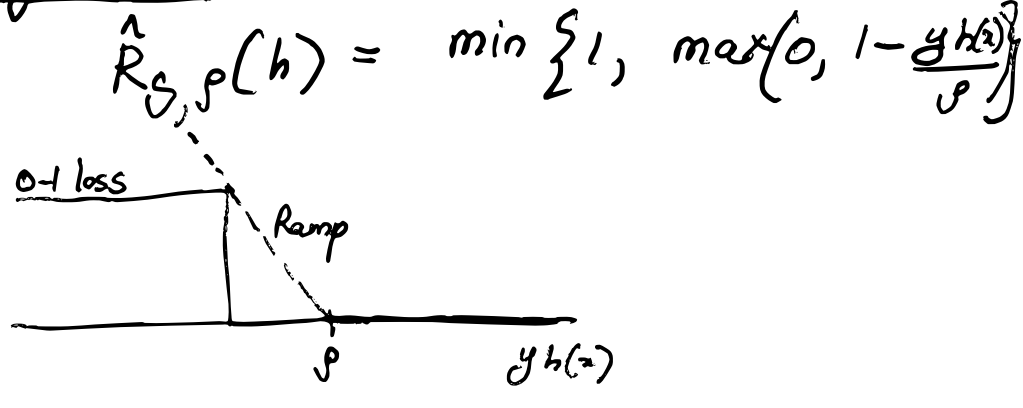
$$S = \{(x_i, y_i)\}_{i=1}^m$$

$$\mathbb{E}_{S \sim \mathcal{D}^m} \hat{R}_S(h) = R(h)$$

Suppose we proved

$$R(h) < \hat{R}_S(h) + \frac{2}{\rho} \text{Rad}_S(\mathcal{H}) + 3 \sqrt{\frac{\log^2 \delta}{2m}} \rightarrow \textcircled{A}$$

Ramp/ Margin loss



If \textcircled{A} , we are done

$$R(h) < \hat{R}_S(h) + \dots < \hat{R}_{S,\rho}(h) + \dots$$

(In practice, we solve ERM using ramp/hinge loss)

Proof sketch

Today

↓

Want: $R(h) < \hat{R}_S(h) + 2 \text{Rad}_S(\mathcal{H}) + 3 \sqrt{\frac{\log^2 \delta}{2m}}$

and use

$$\text{Rad}_S(\mathcal{H}) < \frac{\gamma \Lambda}{\sqrt{m}}$$

$$\gamma = \sup_x \|x\| \quad \Lambda = \sup_w \|w\|$$

$$\ell(z, h) = \mathbb{1}_{\{h(x) \neq y\}}$$

$$0 \leq \ell(z, h) \leq 1$$

$$0 \leq \hat{R}_S(h) \leq 1$$

$$0 \leq R_S(h) \leq 1 \quad R_S(h) = \mathbb{E}_{z \sim \mathcal{D}} \ell(z, h)$$

In more generality, we will prove a Rademacher Complexity-based generalization bound when $\ell \in [0, 1]$.

(Chapter 3 of Mohri + Chapter 5 of Mohri).

(Also for algorithmic-stability - different idea, depending on learning algorithm for generalization - useful inequality)

McDiarmid's inequality)

$$1) \Phi(S) = \sup_{h \in \mathcal{H}} \hat{R}_S(h) - R_S(h)$$

$$2) \text{ Apply M. inequality to } \Phi(S)$$

McDiarmid's inequality

$$S = z_1, \dots, z_m$$

If S' differs from S in one element such that

$$|\Phi(S) - \Phi(S')| \leq c \text{ holds}$$

for every $S, S' \sim D^m$,

Then, for any $t > 0$,

$$\Pr(\Phi(S) - \mathbb{E}_S \Phi(S) > t)$$

$$\leq e^{-\frac{2t^2}{mc^2}}$$

ie with probability $> 1 - \delta$,

$$\Phi(S) - \mathbb{E}_S \Phi(S) < \sqrt{\frac{mc^2 \log 1/\delta}{2}}$$

$$\Phi(S) = \sup_{h \in \mathcal{H}} \hat{R}_S(h) - R_S(h)$$

$$\begin{aligned} |\Phi(S) - \Phi(S')| &= \left| \sup_{h \in \mathcal{H}} (\hat{R}_S(h) - R_S(h)) - \sup_{h \in \mathcal{H}} (\hat{R}_{S'}(h) - R_{S'}(h)) \right| \\ &\leq \left| \sup_{h \in \mathcal{H}} (\hat{R}_S(h) - \hat{R}_{S'}(h)) \right| \quad \left(\sup A - \sup B \leq \sup(A-B) \right) \\ &= \frac{1}{m} \left| \sup_{h \in \mathcal{H}} \sum_{i=1}^m (\ell(z_i, h) - \ell(z'_i, h)) \right| \\ &< \frac{1}{m} \end{aligned}$$

(Φ satisfies B.D. property from Mc. ineq).

Apply Mc. in. : with probability $1-\delta$

$$\underbrace{\sup_{h \in \mathcal{H}} R_S(h) - \hat{R}_S(h)}_{\Phi(S)} \leq \underbrace{\mathbb{E}[\Phi(S)]}_{\sqrt{\frac{\log 1/\delta}{2m}}}$$

$$\begin{aligned} \mathbb{E}_S[\Phi(S)] &= \mathbb{E}_{S \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} (R_S(h) - \hat{R}_S(h)) \\ &= \mathbb{E}_S \sup_{h \in \mathcal{H}} (\mathbb{E}_{S'} \hat{R}_{S'}(h) - \hat{R}_S(h)) \\ &\leq \mathbb{E}_{S, S'} \sup_{h \in \mathcal{H}} (\hat{R}_{S'}(h) - \hat{R}_S(h)) \quad (\because \sup(A+B) \leq \sup A + \sup B) \\ &= \mathbb{E}_{S, S'} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m (\ell(z'_i, h) - \ell(z_i, h)) \\ &= \mathbb{E}_{\sigma, S, S'} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(z'_i, h) - \ell(z_i, h)) \quad (\because \mathbb{E}_{S, S'} \text{ symmetric diff}) \\ &\leq 2 \mathbb{E}_{\sigma, S} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \ell(z_i, h) \rightarrow \textcircled{B} \end{aligned}$$

$$(\sup(A+B) \leq \sup A + \sup B)$$

Recall

$$\text{Rad}_S(\mathcal{H}) = \mathbb{E}_{\sigma} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i)$$

$$\text{Rad}_m(\mathcal{H}) = \mathbb{E}_{S \sim \mathcal{D}^m} \text{Rad}_S(\mathcal{H})$$

$$\begin{aligned} \mathbb{E}_S \Phi(S) &\leq 2 \mathbb{E}_{\sigma, S} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \ell(z_i, h) \\ &= 2 \mathbb{E}_{\sigma, S} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \left(\frac{1 - y_i h(x_i)}{2} \right) \end{aligned}$$

(0-1 loss)

$$= \mathbb{E}_{\sigma, S} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i y_i h(x_i)$$

($y_i = \pm 1$)

$$\leq \mathbb{E}_{\sigma, S} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i)$$

(h belongs to linear fns.)

$$= \text{Rad}_m(\mathcal{H})$$

Implications:

Bias-complexity tradeoff

General form

$$R(h) \leq \hat{R}_S(h) + \text{Rad}(\text{sgn} \mathcal{H})$$

with high probability

A class \mathcal{H} is PAC-learnable with a "sample complexity" $m_{\mathcal{H}}: \mathbb{R}^{2^+} \rightarrow \mathbb{N}$ if $m \geq m_{\mathcal{H}}(\epsilon, \delta) \Rightarrow$
 $|R(h) - \hat{R}_S(h)| < \epsilon$
 with probability $1-\delta$,
 for any \mathcal{D} .

For $h_S \in \text{ERM}(S, \mathcal{H})$

$$R(h_S) = \underbrace{R(h_S) - \min_{h \in \mathcal{H}} R(h)}_{\text{estimation error}} + \underbrace{\min_{h \in \mathcal{H}} R(h)}_{\text{approximation error}}$$

depends on \mathcal{H}

↑ comp of \mathcal{H} ,

↓ approx error

Sample complexity: how representative is empirical error of $R(h)$

Computational complexity

Boosting

Freund Shapire 1995 , 1999

If you have a bunch of
"rough" hypotheses that are "easy
to learn" (small computational complexity)
can you combine ("boost") them
to get better generalization/
training error

→ Approximation error BAD
→ Estimation error good ✓

Practical algorithm: Adaboost
"adaptive"

→ Iterative algorithm I/p: $S = \{(x_i, y_i)\}_{i=1}^m$
algorithm for "weak" learning
→ how do you combine hypotheses?
"weak rules"

"weak" learning returns $f_t(x) = \pm 1$

$$\underline{h_T}(x) = \text{sgn} \left(\sum_{t=1}^T w_t \underline{f_t}(x) \right)$$

(after T rounds of boosting)

o/p: h_T

There are no hyperparameters!
except for T