

FOSML: Lecture #15

①

- We now develop learning bounds for SVM. These provide good mathematical justification of the SVM method, even in the non-separable case.
- Recall that the VC-dimension of a family of linear hypotheses on \mathbb{R}^D is $D+1$.
- Apply our learning bound based on VC-dimension immediately yields that for any linear hypothesis h , and $\forall \delta > 0$, then with probability exceeding $1-\delta$,

$$R(h) \leq \hat{R}_S(h) + \sqrt{\frac{2(D+1) \log\left(\frac{em}{D+1}\right)}{m}} + \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{2m}},$$

where $m = |S|$.

- Note that when D is large, this bound is poor, unless m is very large.
- We shall develop learning bounds that are independent of D , i.e. are not cursed by dimensionality.
- Our results will hold for general \mathbb{R} -valued functions, as opposed to classifiers outputting ± 1 .

Defn (margin loss function): For any $p > 0$, the p -margin loss is the function $L_p: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ defined as

$$L_p(y, y') = \Phi_p(y y')$$

where $\Phi_p(x) = \min\left\{1, \max\left\{0, 1 - \frac{x}{p}\right\}\right\}$

$$= \begin{cases} 1, & \text{if } x \leq 0 \\ 1 - \frac{x}{p}, & \text{if } x \in (0, p) \\ 0, & \text{if } x \geq p \end{cases}$$

Defn (empirical margin loss): Given a sample $S = \{x_i\}_{i=1}^m$ and a hypothesis h ,
 the empirical margin loss is

$$\hat{R}_{S,p}(h) = \frac{1}{m} \sum_{i=1}^m \Phi_p(y_i h(x_i)).$$

Remark: $\forall i \in \{1, \dots, m\}$, $\Phi_p(y_i h(x_i)) \leq \mathbb{1}_{y_i h(x_i) \leq p}$
 $\Rightarrow \hat{R}_{S,p}(h) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i h(x_i) \leq p}$.

Note that Φ_p is $\frac{1}{p}$ Lipschitz, unlike the 0-1 loss which isn't even continuous.

Lemma (Talagrand): Let $\{\Phi_i\}_{i=1}^m$ be L -Lipschitz functions from $\mathbb{R} \rightarrow \mathbb{R}$. Let $\{\epsilon_i\}_{i=1}^m$ be Rademacher r.v. Then for any hypothesis class, \mathcal{H} , consisting of \mathbb{R} -valued functions, the following inequality holds:

$$\frac{1}{m} \mathbb{E}_{\epsilon} \left(\sup_{h \in \mathcal{H}} \sum_{i=1}^m \epsilon_i (\Phi_i \circ h)(x_i) \right) \leq \frac{L}{m} \mathbb{E}_{\epsilon} \left(\sup_{h \in \mathcal{H}} \sum_{i=1}^m \epsilon_i h(x_i) \right) \\ = L \cdot \hat{R}_S(\mathcal{H})$$

In particular, if $\Phi_i \equiv \Phi \forall i$, then $\hat{R}_S(\Phi \circ \mathcal{H}) \leq L \hat{R}_S(\mathcal{H})$.

Proof: Fix a sample $S = \{x_i\}_{i=1}^m$. Then

$$\frac{1}{m} \mathbb{E}_{\epsilon} \left(\sup_{h \in \mathcal{H}} \sum_{i=1}^m \epsilon_i (\Phi_i \circ h)(x_i) \right) = \frac{1}{m} \mathbb{E}_{\epsilon_1, \dots, \epsilon_m} \left(\mathbb{E}_{\epsilon_m} \left(\sup_{h \in \mathcal{H}} (u_{m-1}(h) + \epsilon_m (\Phi_m \circ h)(x_m)) \right) \right)$$

where $u_{m-1}(h) = \sum_{i=1}^{m-1} \epsilon_i (\Phi_i \circ h)(x_i)$. By the definition of the supremum,

$\forall \varepsilon > 0$, $\exists h_1, h_2 \in \mathcal{H}$ s.t.

$$u_{m-1}(h_1) + \Phi_m \circ h_1(x_1) \geq (1-\varepsilon) \left(\sup_{h \in \mathcal{H}} u_{m-1}(h) + (\Phi_m \circ h)(x_1) \right)$$

$$u_{m-1}(h_2) - \Phi_m \circ h_2(x_n) \geq (1-\varepsilon) \left(\sup_{h \in \mathcal{H}} u_{m-1}(h) - (\Phi_m \circ h)(x_n) \right).$$

Then for any $\varepsilon > 0$, by definition of $\mathbb{E}_{\mathcal{Z}_m}$, we have

$$(1-\varepsilon) \mathbb{E}_{\mathcal{Z}_m} \left(\sup_{h \in \mathcal{H}} u_{m-1}(h) + \mathcal{Z}_m(\Phi_m \circ h)(x_n) \right)$$

$$= (1-\varepsilon) \left[\frac{1}{2} \sup_{h \in \mathcal{H}} u_{m-1}(h) + (\Phi_m \circ h)(x_n) + \frac{1}{2} \sup_{h \in \mathcal{H}} u_{m-1}(h) - (\Phi_m \circ h)(x_n) \right]$$

$$\stackrel{\textcircled{\star}}{\leq} \frac{1}{2} \left[u_{m-1}(h_1) + (\Phi_m \circ h_1)(x_n) + u_{m-1}(h_2) - (\Phi_m \circ h_2)(x_n) \right]$$

Letting $s = \text{sign}(h_1(x_n) - h_2(x_n))$, $\textcircled{\star}$ implies

$$(1-\varepsilon) \mathbb{E}_{\mathcal{Z}_m} \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) + \mathcal{Z}_m(\Phi_m \circ h)(x_n) \right]$$

$$\leq \frac{1}{2} \left[u_{m-1}(h_1) + u_{m-1}(h_2) + sL(h_1(x_n) - h_2(x_n)) \right]$$

$$= \frac{1}{2} \left(u_{m-1}(h_1) + sL h_1(x_n) \right) + \frac{1}{2} \left(u_{m-1}(h_2) - sL h_2(x_n) \right)$$

$$\leq \frac{1}{2} \sup_{h \in \mathcal{H}} \left(u_{m-1}(h) + sL h(x_n) \right) + \frac{1}{2} \sup_{h \in \mathcal{H}} \left(u_{m-1}(h) - sL h(x_n) \right)$$

$$= \mathbb{E}_{\mathcal{Z}_m} \left(\sup_{h \in \mathcal{H}} u_{m-1}(h) + \mathcal{Z}_m L h(x_n) \right)$$

Since this holds $\forall \varepsilon > 0$, we get

$$\mathbb{E} \left(\sup_{h \in \mathcal{H}} U_{m-1}(h) + Z_m(\Phi_{m-1} h)(x_m) \right) \leq \mathbb{E} \left(\sup_{h \in \mathcal{H}} U_{m-1}(h) + Z_m L h(x_m) \right). \quad (9)$$

The result follows by doing the same thing for the expectations over Z_i , $i \leq m-1$.

We can use Talagrand's lemma to prove the following learning bound:

Theorem (margin bound for binary classification): Let \mathcal{H} be a set of \mathbb{R} -valued functions. Fix $\rho > 0$. Then $\forall \delta > 0$, with probability exceeding $1-\delta$, both of the following hold $\forall h \in \mathcal{H}$:

$$(a.) R(h) \leq \hat{R}_{S, \rho}(h) + \frac{2}{\rho} \tilde{R}_m(\mathcal{H}) + \sqrt{\frac{\log(\frac{1}{\delta})}{2m}}$$

$$(b.) R(h) \leq \hat{R}_{S, \rho}(h) + \frac{2}{\rho} \tilde{R}_m(\mathcal{H}) + 3 \sqrt{\frac{\log(\frac{1}{\delta})}{2m}}$$

Proof: We begin with (a); (b.) follows immediately by earlier results relating $R(\cdot)$ and empirical $R(\cdot)$.

Let $\tilde{\mathcal{H}} = \{z = (x, y) \mapsto y h(x)\}_{h \in \mathcal{H}}$. Consider the following family of functions taking values in $[0, 1]$:

$$\mathcal{H}_z = \{\Phi_{\rho} \circ f \mid f \in \tilde{\mathcal{H}}\}.$$

By our earlier results on RC , with probability exceeding $1-\delta$, $\forall g \in \mathcal{H}_z$,

$$\mathbb{E}(g(z)) \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2 \tilde{R}_m(\mathcal{H}_z) + \sqrt{\frac{\log(\frac{1}{\delta})}{2m}}$$

$$\Leftrightarrow \mathbb{E}(\Phi_{\rho}(y h(x))) \leq \hat{R}_{S, \rho}(h) + 2 \tilde{R}_m(\Phi_{\rho} \circ \tilde{\mathcal{H}}) + \sqrt{\frac{\log(\frac{1}{\delta})}{2m}}.$$

Since Φ_{ρ} is $\frac{1}{\rho}$ -Lipschitz, Talagrand's Lemma implies

$$\begin{aligned}
R_m(\mathbb{E}_p \tilde{H}) &\leq \frac{1}{p} R_m(\tilde{H}) \\
&= \frac{1}{p} \cdot \frac{1}{m} \mathbb{E}_{S, Z} \left(\sup_{h \in \mathcal{H}} \sum_{i=1}^m \delta_i y_i h(x_i) \right) \\
&= \frac{1}{p^m} \mathbb{E}_{S, Z} \left(\sup_{h \in \mathcal{H}} \sum_{i=1}^m \delta_i h(x_i) \right) \\
&= \frac{1}{p} R_m(\mathcal{H}).
\end{aligned}$$

This gives (4), as needed. ■