

25 Oct 2021

CS492 (F) - Computational Learning Theory

Chap 5 . Support Vector Machine

1. Overview

(in short ~ SVM)

(1) Support Vector Machine is one of the most elegant machine learning algorithms, which performs well in practice and comes with a strong and somewhat surprising theoretical guarantee. Studying it will help you to improve your toolset for developing or selecting a learning algorithm. It will also help you to see how the learning theory that we studied can be used to analyse realistic learning algorithms.

(2) SVM is a learning algo. for classification. Given a labelled sample $S = ((x_1, y_1), \dots, (x_m, y_m)) \in (\mathcal{X} \times \mathcal{Y})^m$ with $y \in \{-1, 1\}$, it computes a linear classifier h , i.e., a function from \mathcal{X} to \mathcal{Y} that has the following form:

$$h(x) = \text{sgn}(f(x)), \quad f(x) = \langle w, x \rangle + b \quad \text{for some } w \in \mathbb{R}^n, b \in \mathbb{R}$$

of course, here we are assuming that $\mathcal{X} \subseteq \mathbb{R}^n$ so that the inner product is well-defined.

(3) The key idea behind SVM is to maximize so called margin when computing h . We will soon explain what this means. One important consequence of margin maximization is that we get generalization bound independent of the dimension n of the input x .
 $(\mathcal{X} \subseteq \mathbb{R}^n)$.

Contrast this with the generalisation bound for the set of hyperplanes. $\mathcal{H}_{PL} = \{x \mapsto \text{Sgn}(\langle w, x \rangle + b) \mid w \in \mathbb{R}^n, b \in \mathbb{R}\}$, that we derived in Chap 3. We recall the bound below:

$$\textcircled{1} \quad \text{VCDim}(\mathcal{H}_{PL}) = n+2 \quad \dots \quad \text{See Example 3.12 in the textbook.}$$

\textcircled{2} By Cor 3.19, we have the following:

For all $D \in \mathcal{P}(X \times Y)$, $m \in \mathbb{N}$, $\delta > 0$,

$$\Pr_{\text{S} \sim D^m} [\forall h \in \mathcal{H}_{PL} \quad R(h) \leq \hat{R}_S(h) + \sqrt{\frac{2(n+1) \log \frac{em}{\delta}}{m}} + \sqrt{\frac{\log 1/\delta}{2^m}}] \geq 1 - \delta.$$

\textcircled{3} Note that the bound becomes loose as m increases.

(4) The independence on the input dimension n indicates that SVM might work well when n is very large or even infinite. Such very large or infinite n arises when we use a large number of features to express inputs to our classifier h (or its linear function f).

(5) Another important and appealing property of SVM is that (the dual version of)
the algorithm and its result (i.e., learnt classifier) both process inputs x only via inner products. This lets SVM be combined nicely with kernel method, which we will study after this chapter. This combination enables SVM to work on inputs of infinite dimension.

2 Support Vector Machine ... Separable case.

(1) We explain SVM in a simpler setting first where labelled inputs can be separated by a hyperplane.

(2) Assumption: about the distribution $D \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$.

$\exists w, b$ s.t.

$$\mathbb{P}_{(x,y) \sim D} [(w \cdot x + b) y > 0] = 1.$$

Intuitively, this means that $f(x) = \langle w, x \rangle + b$ classifies every example from D correctly. We will also assume that $\mathcal{X} = \mathbb{R}^n$.

(3) SVM algorithm. (Primal form). minimization of the size of w .

Input: $S = ((x_1, y_1), \dots, (x_m, y_m)) \sim D^m$

$$(w^*, b)^* \stackrel{\text{def}}{=} \underset{\substack{w \in \mathbb{R}^m \\ b \in \mathbb{R}}}{\operatorname{argmin}} \frac{1}{2} \|w\|^2$$

$$\text{subject to } y_i (\langle w, x_i \rangle + b) \geq 1$$

for all $i \in [m]$.

$$f^*(x) \stackrel{\text{def}}{=} \langle w^*, x \rangle + b^* ; h^*(x) \stackrel{\text{def}}{=} \operatorname{sgn}(f^*(x))$$

1, not 0.

return h^* .

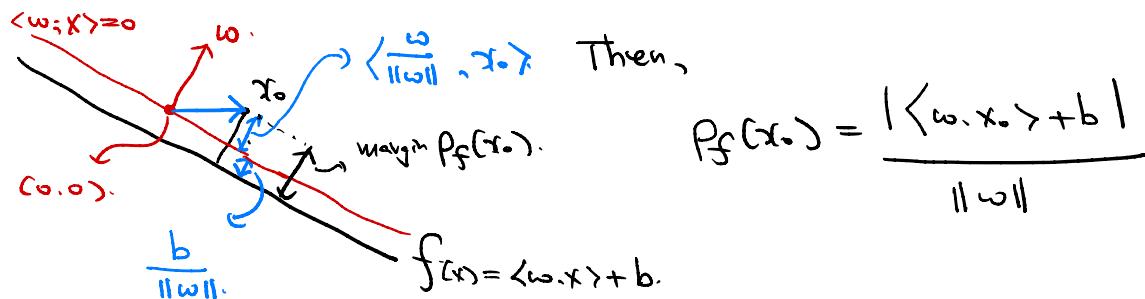
(4) By our assumption, the constraint of the SVM optimization is satisfiable. Note that b influences only the constraint, not the optimization objective. Furthermore, it can be eliminated in the constraint: when $J = \{i \mid y_i = +1\}$,

$$(\exists b, \forall i \in [m], y_i (\langle w, x_i \rangle + b) \geq 1) \Leftrightarrow \forall j \in J \forall i \in [m] \setminus J, -1 - \langle w, x_j \rangle \geq 1 - \langle w, x_i \rangle.$$

If we do so, we get an optimization problem over w only whose objective is strictly convex and whose constraint is the conjunction of linear inequalities. Since the constraint is satisfiable, there is a unique sol. w^* .

(5) To understand the meaning of the optimization of SVM,
 we define the margin $\rho_f(x_0)$ of x_0 to a hyperplane $f(x) = \langle w, x \rangle + b$.

It is the distance of x_0 to f .



For a sample $S_x = (x_1, \dots, x_m)$, the geometric margin ρ_f over S_x is $\min_{j \in [m]} \rho_f(x_j)$.

(6) The next derivation shows that the SVM computes a hyperplane f with the maximum geometric margin ρ_f over S_x .

$$\text{Let } \rho_{w,b}(x_0) = \rho_{x_0 \rightarrow \langle w, x \rangle + b}$$

$$\begin{aligned} & \sup_{w,b} \min_{j \in [m]} \rho_{w,b}(x_j) \\ & \forall i, y_i (\langle w, x_i \rangle + b) \geq 0 \\ &= \sup_{w,b} \min_{j \in [m]} \frac{|\langle w, x_j \rangle + b|}{\|w\|} \\ & \forall i, y_i (\langle w, x_i \rangle + b) \geq 0. \\ &= \sup_{w,b} \min_{j \in [m]} \frac{y_j (\langle w, x_j \rangle + b)}{\|w\|} \\ & \forall i, y_i (\langle w, x_i \rangle + b) \geq 0. \\ &= \sup_{w,b} \min_{j \in [m]} \frac{y_j (\langle w, x_j \rangle + b)}{\|w\|} = \sup_{w,b} \frac{1}{\|w\|} \min_{j \in [m]} y_j (\langle w, x_j \rangle + b) = 1 \\ & \min_{j \in [m]} y_j (\langle w, x_j \rangle + b) = 1 \\ &= \boxed{\sup_{w,b} \frac{1}{\|w\|}.} \\ & \forall i \in [m], y_i (\langle w, x_i \rangle + b) \geq 1 \end{aligned}$$

(7) If the maximum of the above supremum exists, it also solves the following optimization problem:

1

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

w.b.

$$\text{subject to } y_i(\langle w, x_i \rangle + b) \geq 1 \text{ for all } i \in [m]$$

SVM

Also, if this optimization problem has a solution, so does the supremum in the previous page. By the explanation in (4), there exists a unique w such that for some b , the SVM optimization problem is solved by (w, b) .

(8) Dual optimization problem and support vectors.

① The SVM optimization problem is convex, and can be dualized via Lagrangian L and Lagrange variables.

$d = (d_1, d_2, \dots, d_m)$ one for each conjunct $y_i(\langle w, x_i \rangle + b) \geq 1$.
with $d_i \geq 0$,

Lagrange RHS of \geq
variable i

$$L(w, b, d) \stackrel{\text{def}}{=} \frac{1}{2} \|w\|^2 + \sum_{i=1}^m d_i (1 - y_i(\langle w, x_i \rangle + b))$$

\uparrow each conjunct $0 \geq 1 - y_i(\langle w, x_i \rangle + b)$

Then,

$$\min_{w,b} \frac{1}{2} \|w\|^2 = \inf_{w,b} \sup_d L(w,b,d)$$

$$\text{s.t. } \forall i \in [m]. y_i(\langle w, x_i \rangle + b) \geq 1 \quad \text{s.t. } \forall i \in [m]. d_i \geq 0.$$

② If we swap \inf and \sup , we get something potentially smaller:

$$\sup_d \inf_{w,b} L(w,b,d) \leq \inf_{w,b} \sup_d L(w,b,d).$$

$\forall i \in [m]. d_i \geq 0.$

But for this specific SVM optimization (i.e., ① above), the equality holds. Furthermore, the so called KKT condition applies in this case, which implies that it suffices to consider

(ω, b, α) 's satisfying the following conditions:

$$\text{(i)} \quad \nabla_{\omega} L(\omega, b, \alpha) = 0 \quad \wedge \quad \nabla_b L(\omega, b, \alpha) = 0$$

$$\text{(ii)} \quad \wedge \quad \forall i \in [m]. \quad \alpha_i (1 - y_i (\langle \omega, x_i \rangle + b)) = 0.$$

Note:

$$\text{(i)} \quad \nabla_{\omega} L(\omega, b, \alpha) = 0 \quad \Rightarrow \quad \omega + \sum_{i=1}^m \alpha_i y_i x_i = 0.$$

$$\text{(ii)} \quad \nabla_b L(\omega, b, \alpha) = 0 \quad \Rightarrow \quad \sum_{i=1}^m \alpha_i y_i = 0.$$

$$\forall i \in [m]. \quad \alpha_i (1 - y_i (\langle \omega, x_i \rangle + b)) = 0 \quad \Rightarrow \quad \forall i \in [m]. \quad \alpha_i = 0 \quad \text{or}$$

$$y_i (\langle \omega, x_i \rangle + b) = 1.$$

When these conditions hold,

$$\begin{aligned} L(\omega, b, \alpha) &= \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i (\langle \omega, x_i \rangle + b)) \\ &= \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^m \alpha_i - \sum_{i,j=1}^m \alpha_i y_i \alpha_j y_j \langle x_i, x_j \rangle \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i y_i \alpha_j y_j \langle x_i, x_j \rangle. \end{aligned}$$

Thus, the dual-equivalent optimization problem is:

$$\boxed{\begin{aligned} \text{3} \quad \sup_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle. \\ \text{s.t.} \quad & \alpha_i \geq 0 \quad \wedge \quad \sum_{i=1}^m \alpha_i = 0 \end{aligned}}$$

③ In typical use cases, there exists α that attains the supremum in 3 and has some $\alpha_i \neq 0$. Assume such an α .

Using it, we can define ω and b :

$$\omega = \sum_{i=1}^m \alpha_i y_i x_i, \quad b = y_{j_0} - \langle \omega, x_{j_0} \rangle \quad \text{for } j_0 \text{ with } \alpha_{j_0} \neq 0.$$

Thus, the hyperplane built by SVM is $\{x | y_{j_0} - \sum_{i=1}^m \alpha_i y_i \langle x_i, x_i \rangle = y_{j_0} - \langle \omega, x_{j_0} \rangle\}$

$$f(x) = \sum_{i=1}^m \alpha_i y_i \langle x_i, x \rangle - \langle \omega, x_{j_0} \rangle + y_{j_0} = \sum_{i=1}^m \alpha_i y_i \langle x_i, x \rangle - \langle x_i, x_i \rangle + y_{j_0}.$$

④ Note that only those x_i with $\alpha_i \neq 0$ contribute to the definition of f . Such x_i 's are called support vectors. Another important fact is that x, x_i are used in f only through their inner products.

⑤ Finally we can derive the geometric margin of f as follows:

$$\begin{aligned} \text{d}_{\text{SVM}} b &= \text{d}_{\text{SVM}} (y_i - \sum_{j=1}^m \alpha_j y_j \langle x_i, x_j \rangle) \quad \text{by the first and third conditions of KKT} \\ \therefore \sum_{i=1}^m \alpha_i y_i b &= \sum_{i=1}^m \alpha_i - \|w\|^2 \\ \therefore 0 &= \sum_{i=1}^m \alpha_i - \|w\|^2 \quad \text{by the second condition of KKT} \\ \therefore \rho^2 &= \frac{1}{\|w\|^2} = \frac{1}{\sum_{i=1}^m \alpha_i}. \end{aligned}$$

3. Support Vector Machine · Non-separable Case.

(1) In practice, often S_{SVM} cannot be separated by a hyperplane. How should we extend SVM to handle this non-separable case?

One answer that we will look at is to use slack variables.

For each constraint $y_i(\langle w, x_i \rangle + b) \geq 1$ in the original SVM optimisation \square , we introduce a slack variable $\xi_i \geq 0$, which measures the amount of violation of this constraint, and replace the constraint by

$$y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i$$

and include ξ_i as a part of optimisation objective.

(2) Thus, the SVM for the non-sep. case works as follows:

Input $S = ((x_1, y_1), \dots, (x_m, y_m)) \subset D^m$ hyperparameters

$$(\hat{w}^*, \hat{b}^*, \hat{\zeta}^*) \stackrel{\text{def}}{=} \arg \min_{w, b, \zeta} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \hat{\zeta}_i^p$$

subject to $y_i (\langle w, x_i \rangle + b) \geq 1 - \hat{\zeta}_i$

$\hat{\zeta}_i \geq 0$

for all $i \in [m]$.

new part. penalize incorrect or small-margin predictions.

+ 4

$$f^*(x) \stackrel{\text{def}}{=} \langle \hat{w}^*, x \rangle + \hat{b}^*$$

$$h^*(x) \stackrel{\text{def}}{=} \operatorname{sgn}(f^*(x))$$

return h^* .

- ① The new $C \sum_{i=1}^m \hat{\zeta}_i^p$ corresponds to empirical risk, and the old $\frac{1}{2} \|w\|^2$ corresponds to regularizer or measure of model complexity. We didn't have something like $C \sum_{i=1}^m \hat{\zeta}_i^p$ in the separable case, because in a sense we can achieve zero empirical risk at that time.

- ② Note that the minimization forces

$$\hat{\zeta}_i = \max(0, 1 - y_i (\langle w, x_i \rangle + b))$$

If we let $\Phi_p(u) = \max(0, (1+u)^p)$, then the optimization from above is equivalent to:

$$\arg \min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \Phi_p(-y_i (\langle w, x_i \rangle + b))$$

which is closely related to what we looked at in Sect. 7, since Φ_p is convex and non-decreasing. In fact, Φ_1 is precisely the hinge loss in Sect. 7.

(3) The new optimisation problem can also be dualised via Lagrangian and Lagrange variables. In this case, we can also switch inf and sup, and the KKT conditions can be imposed to transform the optimisation problem and analyse its solution. In so doing, we can characterise support vectors as well.

① Lagrange variables ... α_i for $0 \geq 1 - \hat{y}_i - y_i(\langle w, x_i \rangle + b)$
 β_i for $0 \geq -\hat{y}_i$

② Lagrangian \mathcal{L} .

$$\begin{aligned} \mathcal{L}(w, b, \hat{\gamma}, \alpha, \beta) \\ = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \hat{\gamma}_i + \sum_{i=1}^m \alpha_i (1 - \hat{\gamma}_i - y_i(\langle w, x_i \rangle + b)) + \sum_{i=1}^m \beta_i (-\hat{\gamma}_i) \end{aligned}$$

③ KKT conditions.

$$\nabla_w \mathcal{L} = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i$$

$$\nabla_b \mathcal{L} = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0$$

$$\nabla_{\hat{\gamma}_i} \mathcal{L} = 0 \Rightarrow C - \alpha_i - \beta_i = 0 \quad (\text{i.e., } \alpha_i + \beta_i = C)$$

$$\forall i. \alpha_i (1 - \hat{\gamma}_i - y_i(\langle w, x_i \rangle + b)) = 0 \Rightarrow \forall i. \alpha_i = 0 \text{ or.}$$

$$y_i (\langle w, x_i \rangle + b) = 1 - \hat{\gamma}_i$$

$$\forall i. \beta_i (-\hat{\gamma}_i) = 0 \Rightarrow \forall i. \beta_i = 0 \text{ or. } \hat{\gamma}_i = 0.$$

④ By what we explained,

$$\begin{aligned} \text{optimization of SVM} &= \inf_{w, b, \hat{\gamma}} \sup_{\alpha, \beta} \mathcal{L}(w, b, \hat{\gamma}, \alpha, \beta) \\ &= \sup_{\alpha, \beta} \inf_{w, b, \hat{\gamma}} \mathcal{L}(w, b, \hat{\gamma}, \alpha, \beta) \quad \text{for all } \hat{\gamma} \\ &\text{s.t. } \alpha_i \geq 0, \beta_i \geq 0, C = \alpha_i + \beta_i, \\ &\quad \sum_{i=1}^m \alpha_i y_i = 0, \quad w = \sum_{i=1}^m \alpha_i y_i x_i \end{aligned}$$

5

We simplify ⑤ by rewriting L and $\beta_i \geq 0$:

$$\beta_i \geq 0 \quad \rightarrow \quad d_i \leq C.$$

$$\begin{aligned} L(w, b, \beta, \rho) &= \sum_{i,j=1}^m \frac{1}{2} d_i d_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^m d_i - \sum_{i,j=1}^m d_i d_j y_i y_j \langle x_i, x_j \rangle \\ &= \sum_{i=1}^m d_i - \frac{1}{2} \sum_{i,j=1}^m d_i d_j y_i y_j \langle x_i, x_j \rangle. \end{aligned}$$

Thus, ⑤ can be simplified as follows.

$$\boxed{6} \quad \begin{aligned} &\sup_{d_i} \sum_{i=1}^m d_i - \frac{1}{2} \sum_{i,j=1}^m d_i d_j y_i y_j \langle x_i, x_j \rangle \\ \text{subj. to } &0 \leq d_i \leq C \quad \text{and} \quad \sum_{i=1}^m d_i y_i = 0 \quad \text{for all } j \in [m]. \end{aligned}$$

Note that the only difference wrt. the objective for the step case is the presence of the upper bound C. By the way, this upper bound ensures that the constraint of ⑥ defines a compact set, so that the optimisation objective (as a continuous function on d) can be maximised by some d.

So, we can replace sup by max

⑤ If a maximizing α of ⑥ has some d_{j_0} with $0 < d_{j_0} < C$, we can derive w and b as before:

$$w = \sum_{i=1}^m d_i y_i w_i \quad b = y_{j_0} - \sum_{j=1}^m d_j y_j \langle x_j, x_{j_0} \rangle.$$

Support vectors are, however, slightly different. x_i with d_{j_0} is still called support vector, but $d_{j_0} \neq 0$ doesn't imply that $y_i (\langle w, x_i \rangle + b) = 1$

Instead, it implies that $y_i (\langle w, x_i \rangle + b) = 1 - \beta_{j_0}$. The traditional support vectors are x_i 's with $d_i \neq 0$ and $C - d_i = \beta_i \neq 0$.

4. Margin Theory

(1) As we said in the beginning of this chapter, maximizing the margin ensures that SVM has a generalization bound independent of the input dimension. Specifically, we have the following bound: for all $r, r' > 0$, $D \in \Pr(\mathcal{X} \times \mathcal{Y})$, $m \in \mathbb{N}$, and $S > 0$,

$$\text{if } \mathbb{P}_{\substack{(x,y) \sim D}} [\|x\| \leq r] = 1 \quad \text{and} \quad f_f = \{x \mapsto \text{sgn}(\langle \omega, x \rangle) \mid \omega \in \mathbb{R}^m \text{ and } \|\omega\| \leq \Lambda\}$$

then

$$\mathbb{P}_{\substack{S \sim \mathcal{D}^m}} [\forall h_\omega \in \mathcal{H}, \forall p \in [0, r']]$$

$$R(h_\omega) \leq \frac{1}{m} \sum_{i=1}^m \max(0, 1 - \frac{y_i (\langle \omega, x_i \rangle)}{p}) + 4 \sqrt{\frac{r^2 \Lambda^2}{p^2 m}} + \sqrt{\frac{\log \log(2r'/p)}{m}} + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \geq 1 - \delta.$$

large margin p makes these small.

* bound not exactly for the formula used by SVM, but for closely related one

Contrast it with the VC-dimension bound for

the hypo. set with hyperplanes \mathcal{H}_{PL} (which is very closely related to \mathcal{H}_f from above).

$$\mathbb{P}_{\substack{S \sim \mathcal{D}^m}} [\forall h \in \mathcal{H}_{PL}, R(h) \leq \hat{R}_S(h)] + \sqrt{\frac{2(m+1) \log \frac{m}{\delta}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \geq 1 - \delta.$$

dependence on n , the dimension of the inputs.

The SVM margin bound says that with high probability, the generalization error of any h is bounded by the empirical slack-variable-based error plus terms related to margin.

(2) We will now discuss margin theory, which shows the SVM margin bound in the previous page.

Def 5.5 Let $p > 0$.



$$\textcircled{1} \quad \Phi_p: \mathbb{R} \rightarrow \mathbb{R}_+ = [0, \infty)$$

$$\textcircled{2} \quad \Phi_p(x) \stackrel{\text{def}}{=} \min(1, \max(1 - \frac{|x|}{p})) = \begin{cases} 1 & \text{if } |x| \leq p \\ 1 - \frac{|x|}{p} & \text{if } 0 \leq |x| \leq p \\ 0 & \text{if } |x| \geq p \end{cases}$$

\textcircled{3} The p -margin loss function is $L_p: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ defined by

$$L_p(y, y') \stackrel{\text{def}}{=} \Phi_p(y y')$$

\textcircled{1}

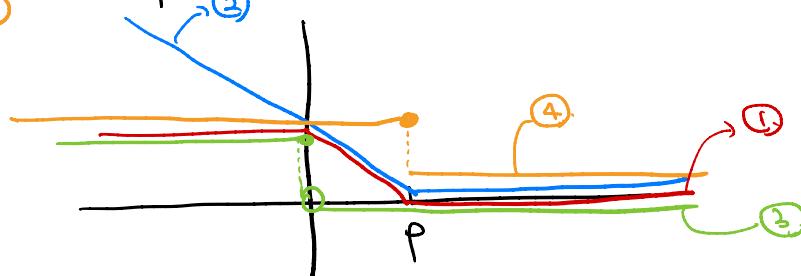
\textcircled{2}

for p .

\textcircled{3}

Compare Φ_p with the hinge loss and the functions $x \mapsto \mathbf{1}_{\{x \leq 0\}}$ and $x \mapsto \mathbf{1}_{\{x \leq p\}}$.

\textcircled{4}



Only \textcircled{2} is convex.

Def 5.6 $S = ((x_1, y_1), \dots, (x_m, y_m))$... sample $f: \mathcal{X} \rightarrow \mathbb{R}$.
or just empirical margin loss.

The empirical p -margin loss of f with respect to S is:

$$\widehat{R}_{S,p}(f) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m L_p(f(x_i), y_i) = \frac{1}{m} \sum_{i=1}^m \Phi_p(y_i f(x_i)).$$

almost = in practice.

Note that

$$\widehat{R}_S(\text{sgn } f) \leq \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{f(x_i) y_i \leq 0\}} \leq \widehat{R}_{S,p}(f) \leq \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{f(x_i) y_i \leq p\}}.$$

(3) Let $\mathcal{F} \subseteq [\mathbb{X} \rightarrow \mathbb{R}]$ measurable

We consider two induced.

Sets of functions from $\mathbb{X} \times \mathbb{Y}$ to \mathbb{R} :

$$\mathcal{G}_0 = \{ z = (x, y) \mapsto y f(x) \mid f \in \mathcal{F} \}.$$

$$\mathcal{G} = \{ \underbrace{\Phi_p}_{\text{def}} \circ g_0 \mid g_0 \in \mathcal{G}_0 \}.$$

$$= (x, y) \mapsto \Phi_p(y f(x)) \quad \text{where } g_0(x, y) = y f(x)$$

$$= (x, y) \mapsto L_p(f(x), y)$$

Note that the range of each $g \in \mathcal{G}$ is $[0, 1]$. Thus, we can apply Thm 3.3., and get the following:

$$\forall D \in \mathcal{P}(\mathbb{X} \times \mathbb{Y}) \quad \forall n \in \mathbb{N} \quad \forall \delta > 0.$$

$$\Pr_{S \sim D^m} \left[\forall g \in \mathcal{G}, \mathbb{E}_{(x, y) \sim D} [g(x, y)] \leq \frac{1}{m} \sum_{i=1}^m g(x_i, y_i) + 2R_m(g) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

$$\text{and} \quad \Pr_{S \sim D^m} \left[\forall g \in \mathcal{G}, \mathbb{E}_{(x, y) \sim D} [g(x, y)] \leq \frac{1}{m} \sum_{i=1}^m g(x_i, y_i) + 2\hat{R}_S(g) + 3\sqrt{\frac{\log \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

But the sum on the RHS is $\hat{R}_{S, p}(f)$ where $g(x, y) = L_p(f(x), y)$.

Also, the expectation on the LHS is larger than or equal to the risk of ($\text{sgn } f$):

$$\begin{aligned} R(\text{sgn } f) &= \Pr_{(x, y) \sim D} [\text{sgn}(f(x)) \neq y] \\ &\stackrel{\text{hypothesis}}{\leq} \Pr_{(x, y) \sim D} [1_{\{f(x)y \leq 0\}}] \leq \Pr_{(x, y) \sim D} [\Phi_p(f(x), y)] \\ &\quad \text{and } L_p(f(x), y) = g(x, y). \end{aligned}$$

By Talagrand's lemma that we will explain soon,

$$\hat{R}_S(g) \leq \frac{1}{p} \hat{R}_S(g_0) \quad \text{and } R_m(g) \leq \frac{1}{p} R_m(g_0).$$

Furthermore -

$$\hat{R}_S(g_0) = \frac{1}{m} \mathbb{E}_{\substack{g_0 \sim \mathcal{G}_0 \\ \text{Unif}[0, 1]^m}} \left[\sup_{g \in \mathcal{G}_0} \sum_{i=1}^m g_i g_0(x_i, y_i) \right] = \frac{1}{m} \mathbb{E}_{f \in \mathcal{F}} \left[\sup_{g \in \mathcal{G}_0} \sum_{i=1}^m g_i y_i f(x_i) \right] = \frac{1}{m} \mathbb{E}_{f \in \mathcal{F}} \left[\sup_{g \in \mathcal{G}_0} \sum_{i=1}^m g_i f(x_i) \right] = \hat{R}_S(F)$$

By taking the expectation over S , we can also get

$$R_m(\mathcal{G}_0) = R_m(\mathcal{F}).$$

Thus, we get the following generalization bound:

Theorem 5.8. [Margin Bound for Binary Classification].

$\forall D \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}) \quad \forall n \in \mathbb{N} \quad \forall S \geq 0$.

$$\Pr_{\text{S} \sim D^m} [\forall f \in \mathcal{F} \quad R(\text{sgn} \circ f) \leq \hat{R}_{S,p}(f) + \frac{2}{p} R_m(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2m}}] \geq 1 - \delta.$$

and

$$\Pr_{\text{S} \sim D^m} [\forall f \in \mathcal{F} \quad R(\text{sgn} \circ f) \leq \hat{R}_{S,p}(f) + \frac{2}{p} \hat{R}_S(\mathcal{F}) + 2\sqrt{\frac{\log^2(1/\delta)}{2m}}] \geq 1 - \delta.$$

Note that using large p reduces the model-complexity part of the bound, but increases the margin error. So, there is a trade-off here.

To complete the proof of the theorem, we need to show the next lemma.

Lemma 5.9 [Talagrand's lemma].

$\bar{\mathbb{E}}_1, \dots, \bar{\mathbb{E}}_m$ --- l -Lipschitz functions from \mathbb{R} to \mathbb{R} .
(bounded)

\mathcal{G}_1 --- any set of measurable functions from \mathbb{Z}^m to \mathbb{R} .

$S = (z_1, \dots, z_m) \in \mathbb{Z}^m$.

means
 $\{\bar{\mathbb{E}}_i \circ g \mid g \in \mathcal{G}_1\}$.

\Rightarrow

$$\frac{1}{m} \mathbb{E}_{\vec{\sigma} \in \text{Unif}(-1, 1)^m} \left[\sup_{g \in \mathcal{G}_1} \sum_{i=1}^m \bar{\mathbb{E}}_i \times (\bar{\mathbb{E}}_i \circ g(z_i)) \right]$$

$$\leq \frac{l}{m} \mathbb{E}_{\vec{\sigma}} \left[\sup_{g \in \mathcal{G}_1} \sum_{i=1}^m \bar{\sigma}_i g(z_i) \right] = l \hat{R}_S(\mathcal{G}_1).$$

In particular, if $\bar{\mathbb{E}}_i = \bar{\mathbb{E}}$ for all $i \in [m]$,

$$\hat{R}_S(\bar{\mathbb{E}} \circ \mathcal{G}_1) \leq l \hat{R}_S(\mathcal{G}_1).$$

Proof We will show that

$$\mathbb{E} \left[\sup_{g \in \mathcal{G}_1} \sum_{i=1}^m b_i \times \Phi_i(g(z_i)) \right] \leq l \mathbb{E} \left[\sup_{g \in \mathcal{G}_1} \sum_{i=1}^m b_i g(z_i) \right].$$

The proof is by induction on m .

The base case is $m=0$. Then, both sides of the inequality are $-\infty$. So, the inequality holds.

For the inductive case $m > 0$, we will show that

$$\begin{aligned} \mathbb{E} \left[\sup_{\substack{b_m \\ g}} u_{m-1}(g) + b_m \Phi_m(g(z_m)) \right] \\ \leq \mathbb{E} \left[\sup_{\substack{b_m \\ g}} u_{m-1}(g) + b_m l g(z_m) \right] \end{aligned}$$

where $u_{m-1}(g) = \sum_{i=1}^{m-1} b_i \Phi_i(g(z_i))$

Pick any $\varepsilon > 0$. Then, for all (b_1, \dots, b_{m-1}) , $\exists g_1, g_2 \in \mathcal{G}_1$ s.t.

$$\begin{aligned} u_{m-1}(g_1) + \Phi_m(g_1(z_m)) &\geq (1-\varepsilon) \sup_{g \in \mathcal{G}_1} u_{m-1}(g) + \Phi_m(g(z_m)) \\ u_{m-1}(g_2) - \Phi_m(g_2(z_m)) &\geq (1-\varepsilon) \sup_{g \in \mathcal{G}_1} u_{m-1}(g) - \Phi_m(g(z_m)), \end{aligned}$$

because the supremums in both are bounded. Note that $u_{m-1}(g)$ depends on (b_1, \dots, b_{m-1}) , and so do the choices of g_1 and g_2 .

Thus, $(1-\varepsilon) \mathbb{E} \left[\sup_{\substack{b_m \\ g \in \mathcal{G}_1}} u_{m-1}(g) + b_m \Phi_m(g(z_m)) \right]$

$$\begin{aligned} &= (1-\varepsilon) \left(\frac{1}{2} \left[\sup_g u_{m-1}(g) + \Phi_m(g(z_m)) \right] + \frac{1}{2} \left[\sup_g u_{m-1}(g) - \Phi_m(g(z_m)) \right] \right) \\ &\leq \frac{1}{2} (u_{m-1}(g_1) + \Phi_m(g_1(z_m))) + \frac{1}{2} (u_{m-1}(g_2) - \Phi_m(g_2(z_m))) \\ &\leq \frac{1}{2} (u_{m-1}(g_1) + u_{m-1}(g_2)) + \frac{1}{2} (\Phi_m(g_1(z_m)) - \Phi_m(g_2(z_m))). \end{aligned}$$

Let $s = \text{sgn}(g_1(z_m) - g_2(z_m))$.

$$\begin{aligned} &\leq \frac{1}{2} \left[(u_{m-1}(g_1) + u_{m-1}(g_2)) + s l (g_1(z_m) - g_2(z_m)) \right] \quad (\because l \text{-Lipschitz}) \\ &= \frac{1}{2} (u_{m-1}(g_1) + s l g_1(z_m)) + \frac{1}{2} (u_{m-1}(g_2) - s l g_2(z_m)) \\ &\leq \frac{1}{2} \left(\sup_{g \in \mathcal{G}_1} (u_{m-1}(g) + s l g(z_m)) \right) + \frac{1}{2} \left(\sup_{g \in \mathcal{G}_1} (u_{m-1}(g) - s l g(z_m)) \right) \end{aligned}$$

$$= \mathbb{E} \left[\sup_{g \in \mathcal{G}} u_m(g) + 6m \log(2m) \right].$$

symmetrization by the Rademacher variable ξ_m . D.

(4). To get a meaningful generalization bound for SVM, we have to adjust two aspects of Thm 5.8. First, we need to make it uniform with respect to the margin ρ . Doing so would let us say something about an algorithm that computes ρ based on a sample S and so uses a random ρ . Second, we need to replace $R_m(F)$ by something easier to deal with.

The next results do these adjustments.

Thm 5.9

$$F \subseteq [X \rightarrow \underset{\text{measurable}}{\mathbb{R}}].$$

$$\Omega \in \Pr(X \times Y) \quad \delta > 0 \quad r > 0 \quad m \in \mathbb{N}.$$

\Rightarrow

$$\Pr \left[\sup_{f \in F} R(f) - \widehat{R}_{S,p}(f) \geq \frac{4}{\rho} R_m(F) + \sqrt{\frac{\log \log_2 2^r / \rho}{m}} + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \right] \geq 1 - \delta.$$

and

$$\Pr \left[\sup_{f \in F} R(f) - \widehat{R}_{S,p}(f) \geq \frac{4}{\rho} \widehat{R}_S(F) + \sqrt{\frac{\log \log_2 2^r / \rho}{m}} + 3 \sqrt{\frac{\log \frac{4}{\delta}}{m}} \right] \geq 1 - \delta.$$

We prove the first bound only. The proof of the second is similar.
 Proof. $(p_k \in (0, r])_{k \geq 1}$ and $(\varepsilon_k \in (0, \infty))_{k \geq 1}$... two sequences.

Then, by Thm 5.8,

$$\Pr \left[\sup_{f \in F} R(f) - \widehat{R}_{S,p}(f) > \frac{2}{p_k} R_m(F) + \sqrt{\frac{\log \frac{1}{\delta p_k}}{2m}} \right] \leq \underline{\delta}_k = \varepsilon_k = \exp(-2m \varepsilon_k^2).$$

$$\text{Choose } \varepsilon_k = \varepsilon + \sqrt{\frac{\log k}{m}}$$

$$\begin{aligned} & \mathbb{P} \left[\sup_{\substack{f \in \mathcal{F} \\ k \geq 1}} R(\text{sgn} \circ f) - \widehat{R}_{S, p_k}(f) - \frac{2}{p_k} R_m(\mathcal{F}) - \varepsilon_k \geq 0 \right] \\ & \leq \sum_{k \geq 1} \mathbb{P} \left[\sup_f R(\text{sgn} \circ f) - \widehat{R}_{S, p_k}(f) - \frac{2}{p_k} R_m(\mathcal{F}) - \varepsilon_k \geq 0 \right] \\ & \leq \sum_{k \geq 1} \exp(-2m\varepsilon_k^2) \leq \sum_{k \geq 1} \exp \left(-2m \left(\varepsilon^2 + \frac{\log k}{m} \right) \right) \\ & = \exp(-2m\varepsilon^2) \sum_{k \geq 1} \frac{1}{k^2} = \exp(-2m\varepsilon^2) \times \frac{\pi^2}{6} = 2 \exp(-2m\varepsilon^2). \end{aligned}$$

Choose $p_k = r/2^k$. Let $p_0 = r$.

Then, for any $p \in (0, r]$, $\exists k \geq 1$ s.t. $p \in (p_k, p_{k-1}]$.

$$\text{Note } \frac{1}{p_k} = \frac{2}{p_{k-1}} \leq \frac{2}{p} \quad \text{and so}$$

$$\sqrt{\log k} = \sqrt{\log \log_2(r/p_k)} \leq \sqrt{\log \log_2 2^r/p}$$

Also,

$$\widehat{R}_{S, p_k}(f) \leq \widehat{R}_{S, p}(f) \quad \text{for all } f \in \mathcal{F}.$$

As a result,

$$\begin{aligned} & \sup_{f \in \mathcal{F}} R(\text{sgn} \circ f) - \widehat{R}_{S, p}(f) - \frac{2}{p} R_m(\mathcal{F}) - \sqrt{\frac{\log \log_2 2^r/p}{m}} - \varepsilon \\ & \quad \leq \sup_{\substack{f \in \mathcal{F} \\ k \geq 1}} R(\text{sgn} \circ f) - \widehat{R}_{S, p_k}(f) - \frac{2}{p_k} R_m(\mathcal{F}) - \sqrt{\frac{\log k}{m}} - \varepsilon \\ & \quad = -\varepsilon_k \text{ from above.} \end{aligned}$$

Thus,

$$\begin{aligned} & \mathbb{P} \left[\sup_{f \in \mathcal{F}} R(\text{sgn} \circ f) - \widehat{R}_{S, p}(f) - \frac{2}{p} R_m(\mathcal{F}) - \varepsilon \geq 0 \right] \leq \mathbb{P} \left[\sup_{f \in \mathcal{F}} R(\text{sgn} \circ f) - \widehat{R}_{S, p_k}(f) - \frac{2}{p_k} R_m(\mathcal{F}) - \sqrt{\frac{\log k}{m}} - \varepsilon \geq 0 \right] \\ & \quad \leq 2 \exp(-2m\varepsilon^2) = \delta = \sqrt{\frac{\log 2/\delta}{2m}} \quad \square \end{aligned}$$

Thm 5.10

$$S_{\mathcal{X}} \in \{x \in \mathcal{X} \mid \|x\| \leq r\}^m$$

$$\mathcal{F} = \{x \mapsto \langle w, x \rangle \mid \|w\| \leq \Lambda\}.$$

\Rightarrow

$$\widehat{R}_{S_{\mathcal{X}}}(\mathcal{F}) \leq \sqrt{\frac{r^2 \Lambda^2}{m}}$$

.... bound on empirical Rademacher complexity.

Note that the thm implies that if $\mathbb{P} [\|x\| \leq r]$,
then $\widehat{R}_m(\mathcal{F}) \leq \sqrt{\frac{r^2 \Lambda^2}{m}}$ bound on Rademacher complexity.

Proof. Let $S_{\mathcal{X}} = (x_1, \dots, x_m)$.

$$\begin{aligned} \widehat{R}_{S_{\mathcal{X}}}(\mathcal{F}) &= \frac{1}{m} \mathbb{E}_{\substack{\omega \sim \text{Unif}\{-1, 1\}^m \\ \|\omega\| \leq \Lambda}} \left[\sup_{\omega} \sum_{i=1}^m \epsilon_i \langle \omega, x_i \rangle \right] \\ &= \frac{1}{m} \mathbb{E}_{\substack{\omega \\ \|\omega\| \leq \Lambda}} \left[\sup_{\omega} \left\langle \omega, \sum_{i=1}^m \epsilon_i x_i \right\rangle \right] \\ &\leq \frac{1}{m} \mathbb{E}_{\substack{\omega \\ \|\omega\| \leq \Lambda}} \left[\sup_{\omega} \|\omega\| \left\| \sum_{i=1}^m \epsilon_i x_i \right\| \right] = \frac{1}{m} \mathbb{E}_{\substack{\omega}} \left[\left\| \sum_{i=1}^m \epsilon_i x_i \right\| \right] \end{aligned}$$

Cauchy-Schwarz inequality.

Jensen's inequality.

$$\begin{aligned} &\leq \frac{1}{m} \mathbb{E}_{\substack{\omega}} \left[\left\| \sum_{i=1}^m \epsilon_i x_i \right\|^2 \right]^{\frac{1}{2}} \\ &\leq \frac{1}{m} \mathbb{E}_{\substack{\omega}} \left[\sum_{i=1}^m \epsilon_i^2 \|x_i\|^2 + \sum_{i \neq j} \epsilon_i \epsilon_j \langle x_i, x_j \rangle \right]^{\frac{1}{2}} \\ &= \frac{1}{m} \left(\sum_{i=1}^m \|x_i\|^2 \right)^{\frac{1}{2}} = \sqrt{\frac{r^2 \Lambda^2}{m}} \end{aligned}$$

We now bring Thm 5.9 and Thm 5.10 together.

Corollary

$$F = \{x \mapsto \langle w, x \rangle \mid \|w\| \leq \Delta\}.$$

$$\Pr_{D \in \mathcal{P}_r}(\exists x, y) \text{ s.t. } \Pr_{(x, y) \sim D}[\|x\| \leq r] = 1.$$

$$m \in \mathbb{N}, \quad \delta > 0, \quad r' > 0.$$

\Rightarrow

$$\Pr_{S \sim D^m} [\forall f \in F \forall p \in (0, r'] R(\text{sgn } f) \leq \widehat{R}_{S,p}(f) + \frac{4}{p} \sqrt{\frac{r^2 \Delta^2}{m}} + \sqrt{\frac{\log \log \frac{2m}{\delta}}{m}} + \sqrt{\frac{\log \frac{2}{\delta}}{2m}}] \geq 1 - \delta.$$

* The upper bound doesn't say anything about the dimension of the input. But the bound r on the size of the input x plays an important role here.

(S) From the above corollary, we can derive the upper bound mentioned when we started our discussion on margin theory. We just need to replace $\widehat{R}_{S,p}(f)$ by its upper bound:

$$\widehat{R}_{S,p}(f) \leq \frac{1}{m} \sum_{i=1}^m \max\left(0, 1 - y_i \langle \frac{w}{p}, x_i \rangle\right)$$

where $f(x) = \langle w, x \rangle$.