

# Support Vector Machine (SVM)

Oct 126 / 2023

§1

## Overview

(1) Classification Pb

Classification problem

$$S \in (\mathcal{X} \times \mathcal{Y})^m, \mathcal{Y} = \{-1, 1\}$$

SVM  $h \in [\mathcal{X} \rightarrow \mathcal{Y}]$ ,  $h(x) = \text{sgn}(\langle w, x \rangle + b)$ .

(2) Theoretical Analysis

$$\text{① } \Pr[\forall h \in H, R(h) \leq \underbrace{\dots}_{\substack{\text{Independent of} \\ \text{input dimension}}} \geq 1 - \delta.$$

②  $H_{PL} = \{x \mapsto \text{sgn}(\langle w, x \rangle + b) \mid w \in \mathbb{R}^n, b \in \mathbb{R}\}$ , in comparison using VCDim,  
 $\text{VCDim}(H_{PL}) = n+1$

$$\Pr[\forall h \in H_{PL}, R(h) \leq \hat{R}_S(h) + \underbrace{\sqrt{\frac{2(n+1)\log(\frac{em}{\delta})}{m}} + \sqrt{\frac{\log(1/\delta)}{2m}}}_{\hookrightarrow n\text{-dependent.}} \geq 1 - \delta.$$

③ Dualisation via Lagrangian.

§2.

## Support Vector Machine

(1)

Linearly separable assumption

$D \in \text{Pr}(\mathcal{X} \times \mathcal{Y})$ ,  $\mathcal{X} = \mathbb{R}^n$ . s.t.  $\exists w \in \mathbb{R}^n, b \in \mathbb{R}$  with

$$\Pr_{(x,y) \sim D} [(\langle w, x \rangle + b)y > 0] = 1.$$

(2)

SVM.

input  $S = ((x_1, y_1), \dots, (x_m, y_m))$

output  $(w^*, b^*) \in \arg \min \frac{1}{2} \|w\|^2$  subject to  $y_i(\langle w, x_i \rangle + b) \geq 1 \quad \forall i \in [m]$

return  $f^*(x) = \langle w^*, x \rangle + b^*$

(3)

### Margin Maximization

$$\textcircled{1} f(x) = \langle w, x \rangle + b, f(x_0) = 0.$$

$$x_0 \in \mathbb{R}^n, x_0 = x + \gamma w$$

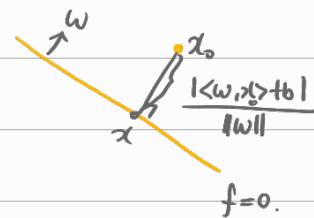
$$\begin{aligned} \textcircled{1} |f(x_0)| &= |f(x) - f(x_0)| = |\langle w, x_0 - x \rangle| \\ &= |\gamma| \|w\|^2 \end{aligned}$$

$$\Rightarrow \|x_0 - x\| = \|\gamma w\| = \frac{|f(x_0)|}{\|w\|}$$

(ii) For any  $x$  with  $f(x) = 0$ ,

$$|f(x) - f(x_0)| = |\langle w, x_0 - x \rangle| \geq \|w\| \cdot \|x_0 - x\|$$

$$\Rightarrow \|x_0 - x\| \geq \frac{|f(x) - f(x_0)|}{\|w\|} = \frac{|f(x_0)|}{\|w\|}$$



$$\min_x \|x_0 - x\|^2 \text{ s.t. } f(x) = 0.$$

By Lagrangian multiplier method,  $x - x_0 = \lambda \nabla f(x) = \lambda w$  for some  $\lambda$   
 $\Rightarrow \|x - x_0\| = \lambda \|w\|$ .

$$\textcircled{2} \sup_{w,b} \min_{i \in [m]} P_{w,b}(x_i) = \frac{|\langle w, x_i \rangle + b|}{\|w\|} \text{ subject to } y_i (\langle w, x_i \rangle + b) \geq 0 \quad \forall i \in [m].$$

$$= \frac{y_i (\langle w, x_i \rangle + b)}{\|w\|},$$

$$= \quad " \quad \text{subject to} \quad \min_{i \in [m]} y_i (\langle w, x_i \rangle + b) = 1$$

∴ Objective is invariant to the scaling of  $(w, b)$ .

$$w' = \frac{w}{\eta}, b' = \frac{b}{\eta}, \eta = \min_{i \in [m]} y_i (\langle w, x_i \rangle + b) > 0.$$

$$\Rightarrow \min_{i \in [m]} y_i (\langle w', x_i \rangle + b') = 1$$

and objective value is the same.

$$= \sup_{w,b} \frac{1}{\|w\|} \text{ s.t. } \min_{i \in [m]} y_i(\langle w, x_i \rangle + b) = 1$$

$$= \sup_{w, b} \frac{1}{\|w\|} \text{ s.t. } \min_{i \in [m]} y_i (\langle w, x_i \rangle + b) \geq 1.$$

7<sup>o</sup>) w,b with  $\min_{i \in [m]} y_i (\langle w, x_i \rangle + b) = \mu \geq 1,$   
 $w' = \frac{w}{\mu}, b' = \frac{b}{\mu}$  has  $\min_{i \in [m]} y_i (\langle w', x_i \rangle + b') = 1$   
 and  $\frac{1}{\mu} = \frac{\mu}{\mu} > \frac{1}{\mu}$

Thus, sup. achieved in smaller domain.

$$= \inf_{w,b} \frac{1}{2} \|w\|^2 \text{ s.t. } y_i(\langle w, x_i \rangle + b) \geq 1 \quad \forall i \in [m].$$

(Question)

Formulation without  $b$ ?

$$\min \frac{1}{2} \|W\|^2 \text{ s.t. } \forall j \in J \quad \forall i \in [m] \setminus J$$

$$-1 - \langle w, \gamma_i \rangle \geq 1 - \langle w, \gamma_j \rangle,$$

where  $J = \{j | y_j = +1\}$ .

( $\Rightarrow$ ) Take  $b \in [\max_{j \in J} 1 - \langle w, x_j \rangle, \min_{\{i \in [m]\} \setminus J} -1 - \langle w, x_i \rangle]$ .

Then, for  $j \in J$ ,  $b \geq 1 - \langle w, x_j \rangle \Rightarrow \langle w, x_j \rangle + b \leq 1$

$$i \notin J, b \leq -\langle w, x_i \rangle \Rightarrow \langle w, x_i \rangle + b \leq -1.$$

$$(\Leftarrow) -(\langle \omega, \lambda_i \rangle + b) \geq 1$$

$$(\Leftarrow) \quad \langle w, x_j \rangle + b \geq 1 \quad \forall j \in J \quad \langle w, x_i \rangle + b \leq -1 \quad \forall i \in [m] \setminus J$$

$$\text{Thus, } \max_{j \in J} |-\langle w, x_j \rangle| \leq b \leq \min_{i \in [m] \setminus J} -|-\langle w, x_i \rangle|$$

{}

## SUM for Separable Case

$$(1) \quad \min_{w,b} \frac{1}{2} \|w\|^2 \quad \text{s.t. } y_i(\langle w, x_i \rangle + b) \geq 1 \quad \forall i \in [n]$$

$\Downarrow F(w, b)$

(2) Dualisation

① Dual variables  $\alpha_i \geq 0$ ,  $0 \geq 1 - y_i(\langle w, x_i \rangle + b) =: g_i(w, b)$   
 (Lagrangian Variables)  
 $\alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}_{\geq 0}^m$

$$② L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i g_i(w, b)$$

Fact

For  $\forall \alpha \in \mathbb{R}^m$ ,

$$(1) \quad \inf_{w,b} L(w, b, \alpha) \leq F(w, b) \quad \forall w, b \text{ s.t. } 0 \geq g_i(w, b) \quad \forall i \in [m].$$

$$\therefore L(w, b, \alpha) = F(w, b) + \sum_{i=1}^m \alpha_i g_i(w, b) \\ \leq F(w, b)$$

Corollary

$$\sup_{\alpha} \inf_{w,b} L(w, b, \alpha) \leq \inf_{w,b} F(w, b). \\ \text{s.t. } 0 \geq g_i(w, b) \quad \forall i \in [m]$$

(2)

$(w^*, b^*, \alpha^*)$  is a saddle of  $L \Rightarrow (w^*, b^*)$  is a soln. of the

original problem.



$$L(w^*, b^*, \alpha) \leq L(w^*, b^*, \alpha^*) \leq L(w, b, \alpha^*) \quad \forall \alpha \geq 0, \forall w, b.$$

$$\therefore L(w^*, b^*, \alpha) \leq L(w^*, b^*, \alpha^*).$$

$$\sum_{i=1}^m \alpha_i^* g_i(w^*, b^*) \leq \sum_{i=1}^m \alpha_i^* g_i(w^*, b^*)$$

$$\alpha_i^* g_i(w^*, b^*) = 0$$

①  $(w^*, b^*)$  is feasible. i.e.  $g_i(w^*, b^*) \leq 0 \quad \forall i \in [m]$

Otherwise, we can choose  $\alpha_i > 0$  large s.t.  $L(w^*, b^*, \alpha) > L(w^*, b^*, 0)$

②  $0 \leq \sum_{i=1}^m \alpha_i^* g_i(w^*, b^*)$

As  $\alpha_i^* g_i(w^*, b^*) \leq 0 \quad \forall i \in [m]$ ,  $\alpha_i^* g_i(w^*, b^*) = 0 \quad \forall i \in [m]$

③  $F(w^*, b^*) = L(w^*, b^*, \alpha^*)$

$\leq L(w', b', \alpha^*)$

$\leq F(w', b') \quad (\text{If } (w', b') \text{ feasible}).$

Thus,  $(w^*, b^*)$  is the optimal solution of primal

$$\nabla_{w,b} \inf L(w, b, \alpha) = L(w^*, b^*, \alpha^*) \geq L(w^*, b^*, \alpha)$$

④  $(w^*, b^*, \alpha^*)$  is the optimal solution of dual :

$$\sup_{\alpha \geq 0} \inf_{w,b} L(w, b, \alpha) = L(w^*, b^*, \alpha^*).$$

$$\circlearrowleft) L(w^*, b^*, \alpha^*) \leq L(w, b, \alpha^*) \Rightarrow L(w^*, b^*, \alpha^*) \leq \inf_{w,b} L(w, b, \alpha^*)$$

( $\leq$ ) by weak duality.

$$\leq \sup_{\alpha \geq 0} \inf_{w,b} L(w, b, \alpha)$$

Fact 3

$(w^*, b^*)$  is a soln of the original problem

if and only if  $\exists \alpha^*$  s.t.  $(w^*, b^*, \alpha^*)$  is a saddle of the associated Lagrangian.

$$\sup_{\alpha} \inf_{w,b} L(w, b, \alpha) \quad \text{s.t.} \quad \nabla_{w,b} L(w, b, \alpha) = 0 \\ \alpha \geq 0.$$

Note

(Lagrangian of SVM).

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i (\langle w, x_i \rangle + b))$$

$$\nabla_w L(w, b, \alpha) = w + \left( \sum_{i=1}^m \alpha_i (-y_i x_i) \right) = 0 \Rightarrow w = \sum \alpha_i y_i x_i$$

$$\nabla_b L(w, b, \alpha) = \sum_{i=1}^m \alpha_i (-y_i) = 0 \Rightarrow \sum \alpha_i y_i = 0.$$

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

Solve  $\sup_{\alpha \geq 0} L(\alpha)$

Then  $w^* = \sum_i \alpha_i^* y_i x_i$ .

If  $\exists i_0$  s.t.  $\alpha_{i_0} > 0$ .

By KKT,  $g_{i_0}(w^*, b^*) = 1 - y_{i_0} (\langle w^*, x_{i_0} \rangle + b^*) = 0$ .

$$\begin{aligned} \Rightarrow b^* &= y_{i_0} - \langle w^*, x_{i_0} \rangle \\ &= y_{i_0} - \sum_{i=1}^m \alpha_i^* y_i \langle x_i, x_{i_0} \rangle \end{aligned}$$

$$f(x) = \langle w^*, x \rangle + b^*$$

$$= \left( \sum_i \alpha_i^* y_i \langle x_i, x - x_{i_0} \rangle \right) + y_{i_0}$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

$$\rho_{w^*, b^*} = \frac{\min_{i \in [m]} y_i (\langle w^*, x_i \rangle + b)}{\|w^*\|} = \frac{1}{\|w^*\|}$$

$$\rho_{w^*, b^*}^2 = \frac{1}{\|w^*\|^2} = \frac{1}{\sum_{i=1}^m \alpha_i^*}$$

$$\alpha_i^* (1 - y_i (\langle w^*, x_i \rangle + b)) = 0.$$

$$\alpha_i^* - \alpha_i^* y_i (\langle w^*, x_i \rangle) = \alpha_i^* y_i b$$

# §1 Remainder/Motivation

Nov 12 / 2023

(1)

## SVM for Separable Case

$$\underset{w,b}{\operatorname{argmin}} \frac{1}{2} \|w\|^2 \quad \text{s.t. } y_i(\langle w, x_i \rangle + b) \geq 1 \quad \forall i \in [m].$$

(2)

## Non-Separable Case

- ① Relax the constraint by slack variable
- ② Additional opt objective that encourage minimum violation of constraint.

§2

## SVM for Non-separable Cases

(1)

### (Slack Variable)

$\xi_i \geq 0$  for each data point  $(x_i, y_i)$ ,  $\xi = (\xi_1, \dots, \xi_m)$ .

(2)

$p \geq 1$  ( $\geq 0$ ,

$$(w^*, b^*, \xi^*) = \underset{w,b,\xi}{\operatorname{argmin}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i^p \quad \text{s.t. } y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \quad \forall i.$$

return  $x \mapsto \operatorname{sgn}(\langle w, x \rangle + b)$

(3)

### (Removing $\xi$ ?)

Given  $(w, b)$ ,  $\xi_i^* = \max(0, 1 - \underbrace{y_i(\langle w, x_i \rangle + b)}_{f_{w,b}(x_i)})$

$$(\xi_i^*)^p = \max(0, (1 - f_{w,b}(x_i))^p) = \bar{\Phi}_p(-y_i f_{w,b}(x_i)),$$

where  $\bar{\Phi}(a) = \max(0, (1+a)^p)$ .

For  $p \geq 0$ ,  $\bar{\Phi}$  is nondecreasing convex.

$$\underset{w,b}{\operatorname{argmin}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \bar{\Phi}_p(-y_i f_{w,b}(x_i))$$

(4)

Dualisation - P=1

$$0 \geq 1 - \xi_i - y_i (\langle w, x_i \rangle + b) \Rightarrow \alpha_i \geq 0 \quad \alpha = (\alpha_1, \dots, \alpha_m)$$

$$0 \geq -\xi_i \quad \beta_i \geq 0 \quad \beta = (\beta_1, \dots, \beta_m)$$

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (\langle w, x_i \rangle + b)) + \sum_{i=1}^m \beta_i (-\xi_i)$$

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

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

$$\frac{\partial}{\partial \xi_i} L = C - \alpha_i - \beta_i = 0 \Rightarrow \alpha_i + \beta_i = C.$$

$$\text{Thus, } w = \sum_{i=1}^m \alpha_i y_i x_i, \sum_{i=1}^m \alpha_i y_i = 0, \alpha_i + \beta_i = C.$$

$$\begin{aligned} \inf_{w, b, \xi} L(w, b, \xi, \alpha, \beta) &= \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \xi_i \xi_j \langle x_i, x_j \rangle \\ &\quad + \sum_{i=1}^m (C - \alpha_i - \beta_i) \xi_i \\ &\quad + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i y_i \langle \sum_{j=1}^m \alpha_j y_j x_j, x_i \rangle \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \xi_i \xi_j \langle x_i, x_j \rangle \end{aligned}$$

$$\text{with } \sum \alpha_i y_i = 0, 0 \leq \alpha_i = C - \beta_i \leq C.$$

KKT Condition:

$$\alpha_i^* (1 - \xi_i^* - y_i (\langle w^*, x_i \rangle + b^*)) = 0 \quad \forall i$$

$$\beta_i^* (-\xi_i^*) = 0 \quad \forall i$$

$$\alpha_i^* \neq 0 \Rightarrow y_i (\langle w^*, x_i \rangle + b^*) = 1 - \xi_i^*$$

$$\textcircled{1} \xi_i^* = 0 \Rightarrow y_i (\langle w^*, x_i \rangle + b^*) = 1. \quad (\text{on marginal plane})$$

\textcircled{2} \xi\_i^\* \neq 0 \Rightarrow x\_i \text{ is outlier}

# Soft Margin SVM



## §3 Margin Theory

(1)

Motivation

$$\text{① } H_{PL} = \{x \mapsto \text{sgn}(\langle w, x \rangle + b) \mid w \in \mathbb{R}^n, b \in \mathbb{R}\}$$

$$VC\text{Dim}(H_{PL}) = n+1$$

$$\underset{S \sim D^m}{\mathbb{P}} \left[ \forall h \in H_{PL} \quad R(h) \leq \hat{R}_S(h) + \sqrt{\frac{2 \log(\epsilon m / \delta)}{m}} + \sqrt{\frac{\log(1/\delta)}{2m}} \right] \geq 1 - \delta.$$

Bound useless if  $n \gg m$

(2)

Goal

$$r, r' > 0 \quad \text{Def} \Pr(x|y), m \in \mathbb{N}, \delta > 0$$

Assume  $\Pr[\|x\| \leq r] = 1$  and  $H = \{x \mapsto \text{sgn}(\langle w, x \rangle) \mid w \in \mathbb{R}^n, \|w\| \leq 1\}$

$$\text{Then, } \underset{S \sim D^m}{\mathbb{P}} \left[ \forall h_w \in H, \forall p \in (0, r'] , R(h_w) \leq \frac{1}{m} \sum_{i=1}^m \max(0, 1 - \frac{y_i \langle w, x_i \rangle}{p}) + 4 \sqrt{\frac{r^2 \Lambda^2}{p^2 m}} + \sqrt{\frac{\log \log(2r'/p)}{m}} + \sqrt{\frac{\log(2/\delta)}{cm}} \right] \geq 1 - \delta$$

upper /

/

Q.

How can we manipulate  $w, \rho$  to make bound small

$$w' = \frac{w}{\rho} \quad \|w'\| = \frac{\|w\|}{\rho} \leq 1 \Rightarrow \rho \leq \frac{1}{\|w\|}.$$

Theorem 5.8

$D \in \Pr(\mathcal{X} \times \mathcal{Y})$ ,  $m \in \mathbb{N}$ ,  $\delta > 0$ ,  $\mathcal{F} \subseteq [\mathcal{X} \rightarrow \mathbb{R}]$  s.t.

$\forall x \in \mathcal{X}, \exists M > 0$  with  $H = \{ \underbrace{\text{sgn of } h_f}_{h_f} \mid f \in \mathcal{F} \}$ ,  $\sup_{f \in \mathcal{F}} |f(x)| < M$

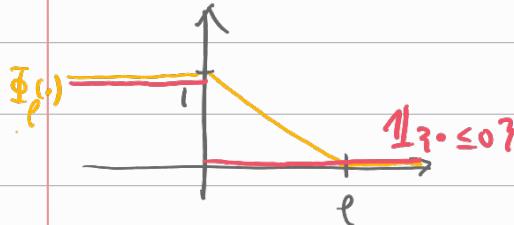
$$\Rightarrow \underset{S \sim D^m}{\mathbb{P}} \left[ \forall f \in \mathcal{F} \quad R(h_f) \leq \hat{R}_{Sp}(f) + \frac{2}{\rho} R_m(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2m}} \right] \geq 1 - \delta$$

Definition

(Empirical  $\rho$ -margin Loss)

$$\hat{R}_{Sp}(f) := \frac{1}{m} \sum_{i=1}^m L_\rho(f(x_i), y_i),$$

$$L_\rho(f(x_i), y_i) := \Phi_\rho(y_i f(x_i)) = \min(1, \max(0, 1 - \frac{f(x_i) y_i}{\rho}))$$



(pf of thm 5.8)  $\exists g \subseteq [Z \rightarrow [0, 1]]$  -  $D \in \Pr(Z)$

$$\underset{S \sim D^m}{\mathbb{P}} \left[ \forall g \in \mathcal{G} \quad \underset{z \sim D}{\mathbb{E}} [g(z)] \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2 R_m(g) + \sqrt{\frac{\log(1/\delta)}{2m}} \right] \geq 1 - \delta.$$

$$g_0 := \{ (x, y) \in \mathcal{X} \times \mathcal{Y} \mapsto f(x)y \mid f \in \mathcal{F} \}$$

$$g_* := \{ (x, y) \mapsto \Phi_\rho(g_0(x, y)) \mid g_0 \in g_0 \}.$$

$$\underset{S \sim D^m}{\mathbb{E}} \left[ \forall g \in \mathcal{G} \quad \underset{(x, y) \sim D}{\mathbb{E}} [g_0(x, y)] \leq \frac{1}{m} \sum_{i=1}^m g_0(x_i, y_i) + 2 R_m(g) + \sqrt{\frac{\log(1/\delta)}{2m}} \right] \geq 1 - \delta$$

$$\forall f \in \mathcal{F} \quad \underset{(x, y) \sim D}{\mathbb{E}} [\Phi_\rho(f(x)y)] = \frac{1}{m} \sum_{i=1}^m \Phi_\rho(f(x_i)y_i)$$

$$R(h_f) = \mathbb{E}_{(x,y) \sim D} [ \mathbb{1}_{\{\text{sgn}(f(x)) \neq y\}} ] \leq \mathbb{E}_{(x,y) \sim D} [ \Phi_p(f(x)y) ]$$

Note .  $\Phi_p$  is  $\frac{1}{p}$ -Lipschitz.

By Telagrand's lemma,  $\hat{R}_s(g) \leq \frac{1}{p} \hat{R}_s(g_0) = \frac{1}{p} \hat{R}_s(\tilde{f})$ .

Exercise

$$\hat{R}_s(g_0) = \hat{R}_s(\tilde{f})$$

(pf)

$$\hat{R}_s(\tilde{f}) = \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(x_i) \right]$$

$$\begin{aligned} \hat{R}_s(g) &= \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{g \in g_0} \sum_{i=1}^m \sigma_i g_0(x_i, y_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i y_i f(x_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{y: y_i \sim \text{Unif}(-1,1)} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i y_i f(x_i) \right] \right] \\ &= \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(x_i) \right] \\ &= \hat{R}_s(\tilde{f}) \end{aligned}$$

Hence,  $R_m(g_0) = R_m(\tilde{f})$ , and  $R_m(g) \leq \frac{1}{p} R_m(\tilde{f})$ .

Lemma

(Telagrand)

$g_0 \subseteq [Z \rightarrow \mathbb{R}]$  s.t.  $\forall z \in Z, \exists M > 0$  with  $\sup_{g \in g_0} |g_0(z)| < M$

and  $\Phi: \mathbb{R} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz. and Bounded

Then,  $g_1 := \Phi(g_0) = \{ \Phi(g_0) \mid g_0 \in g_0 \}$  has

$$\hat{R}_s(g_1) \leq L \hat{R}_s(g_0)$$

(pf)

$$\hat{R}_s(g_1) = \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{g \in g_0} \sum_{i=1}^m \sigma_i \Phi(g_0(z_i)) \right]$$

It suffices to prove

$$\mathbb{E} \left[ \sup_{g_0} \left( \sum_{i=1}^k \sigma_i \Phi(g_0(z_i)) + L \sum_{j=k+1}^m \sigma_j g_0(z_j) \right) \right]$$

$$\leq \mathbb{E} \left[ \sup_{g_0} \left( \sum_{i=1}^k \sigma_i \Phi(g_0(z_i)) + L \sum_{j=k+1}^m \sigma_j g_0(z_j) \right) \right]$$

Which reduces to  $\mathbb{E} \left[ \sup_{g_0} \left( \sum_{i=1}^k \sigma_i \Phi(g_0(z_i)) + L \sum_{j=k+1}^m \sigma_j g_0(z_j) \right) \right] \leq \mathbb{E} \left[ \sup_{g_0} \left( \sum_{i=1}^k \sigma_i \Phi(g_0(z_i)) + L \sum_{j=k+1}^m \sigma_j g_0(z_j) \right) \right]$

$$\sup_{g_0} \left( \sum_{i=1}^k \sigma_i \Phi(g_0(z_i)) + \Phi(g_0(z_k)) + L \sum_{j=k+1}^m \sigma_j g_0(z_j) \right)$$

$$\sup_{g_0} \left( " - \Phi(g_0(z_k)) " \right) \Psi(g_0)$$

for  $\varepsilon \in (0,1)$ ,  $\exists g_1, g_2 \in \mathcal{G}_0$  st.  $\Psi(g_1) > (1-\varepsilon) \sup_{g_0} \Psi(g_0)$   
 $\Psi(g_2) > (1-\varepsilon) \sup_{g_0} \Psi(g_0)$

Let us denote  $\sum_{i=1}^k \sigma_i \Phi(g_0(z_i)) + \sum_{j=k+1}^m \sigma_j g_0(z_j) =: u(g_0) \Rightarrow \Psi(g_0) = u(g_0) + \Phi(g_0)$   
 $\Psi(g_0) = " - "$

$$(1-\varepsilon) \times (*) = (1-\varepsilon) \left[ \frac{1}{2} \times \sup_{g_0} (u(g_0) + \Phi(g_0(z_k))) + \frac{1}{2} \sup_{g_0} (u(g_0) - \Phi(g_0(z_k))) \right]$$

$$\leq \frac{1}{2} (u(g_1) + \Phi(g_1(z_k)) + u(g_2) - \Phi(g_2(z_k)))$$

$$\leq \frac{1}{2} (u(g_1) + u(g_2) + L |g_1(z_k) - g_2(z_k)|)$$

$$= \frac{1}{2} (u(g_1) + u(g_2) + L s(g_1(z_k) - g_2(z_k))),$$

with  $s := \text{sgn}(g_1(z_k) - g_2(z_k))$

$$= \frac{1}{2} (u(g_1) + L s g_1(z_k)) + \frac{1}{2} (u(g_2) - L s g_2(z_k))$$

$$\leq \frac{1}{2} \sup_g (u(g) + L s g(z_k)) + \frac{1}{2} \sup_g (u(g) - L s g(z_k))$$

$$= \mathbb{E}_{\sigma_k} \left[ \sup_g (u(g) + L \sigma_k g(z_k)) \right]$$

$$= (***)$$

As  $\varepsilon > 0$  is arbitrary,  $(*) \leq (***)$

Theorem 5.10  $S_{\mathcal{F}} \in \{\chi \in \mathcal{X} \mid \|x\| \leq r\}$ .  $\mathcal{F} = \{\chi \mapsto \langle w, \chi \rangle \mid \|w\| \leq \Lambda\}$ . Then,

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

(pf)

$$\hat{R}_{S_{\mathcal{F}}}(\mathcal{F}) = \frac{1}{m} \mathbb{E} \left[ \sup_{\substack{w \in \mathbb{R}^n \\ \|w\| \leq \Lambda}} \sum_{i=1}^m \sigma_i \langle w, x_i \rangle \right]$$

$$\leq \frac{1}{m} \mathbb{E} \left[ \sup_{\sigma} \sum |\sigma_i \langle w, x_i \rangle| \right]$$

$$\leq \frac{1}{m} \mathbb{E} \left[ \sup_{\sigma} \|w\| \|\sum \sigma_i x_i\| \right] \quad (\text{-S})$$

$$\leq \frac{1}{m} \mathbb{E} \left[ \|\sum \sigma_i x_i\| \right] \quad (\|w\| \leq \Lambda)$$

$$= \frac{1}{m} \mathbb{E} \left[ \|\sum \sigma_i x_i\| \right]$$

$$\leq \frac{1}{m} (\mathbb{E} [\|\sum \sigma_i x_i\|^2])^{\frac{1}{2}} \quad (\text{Jensen})$$

$$= \frac{1}{m} (\mathbb{E} [\langle \sum_i \sigma_i x_i, \sum_j \sigma_j x_j \rangle])^{\frac{1}{2}}$$

$$= \frac{1}{m} (\mathbb{E} \sum_i \sigma_i^2 \|x_i\|^2 + \sum_{i \neq j} \sigma_i \sigma_j \langle x_i, x_j \rangle)^{\frac{1}{2}}$$

$$\leq \frac{1}{m} (mr^2)^{\frac{1}{2}} = \sqrt{\frac{r^2 \Lambda^2}{m}}. \quad (\|x_i\| \leq r).$$

Theorem 5.9

$\mathcal{F}, D, \delta, m$  same as thm 5.8. Then for  $r > 0$ ,

$$\mathbb{P}_{S \sim D^m} \left[ \forall f \in \mathcal{F} \forall p \in (0, r] R(h_f) \leq \hat{R}_{Sp}(f) + \frac{4}{p} R_m(f) + \sqrt{\frac{\log \log(2r/p)}{m}} + \sqrt{\frac{\log(2f)}{2m}} \right] \geq 1 - \delta.$$

(proof)

$$\text{Let } \varepsilon = \sqrt{\frac{\log(2/\delta)}{2m}}, \text{ so that } \delta = 2e^{-2m\varepsilon^2}.$$

$$\text{Let } \varepsilon_k = \varepsilon + \sqrt{\frac{\log k}{m}}, \quad p_k = r'/2^k.$$

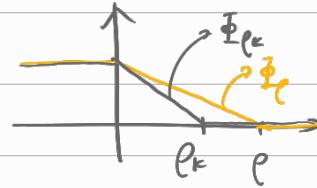
By thm 5.8,  $\mathbb{P}[\exists f \in \mathcal{F} \text{ s.t. } R(h_f) - \hat{R}_{Sp}(f) - \frac{2}{p_k} R_m(f) - \varepsilon_k > 0] \leq e^{-2m\varepsilon_k^2}$

$$\begin{aligned} \mathbb{P} \left[ \sup_{\substack{f \in \mathcal{F} \\ k \geq 1}} \dots \right] &\leq \sum_k e^{-2m\varepsilon_k^2} \\ &\leq \sum_k e^{-2m(\varepsilon^2 + \frac{\log k}{m})} = \sum_k \frac{1}{k^2} e^{-2m\varepsilon^2} \\ &= \frac{\pi^2}{6} e^{-2m\varepsilon^2} \leq 2e^{-2m\varepsilon^2} = \delta. \end{aligned}$$

For  $\rho \in (0, r']$ ,  $\rho \in [\rho_k, \rho_{k+1}]$  for some  $k \geq 1$ .

$$\frac{1}{\rho_k} = \frac{2}{\rho_{k+1}} \leq \frac{2}{\rho} \quad \sqrt{\log k} = \sqrt{\log \log_2 r' / \rho_k} \leq \sqrt{\log \log_2 2r' / \rho}$$

$$\hat{R}_{S, \rho_k}(f) \leq \hat{R}_{S, \rho}(f). \quad \therefore$$



$$\textcircled{1} = \sup_{\substack{f \in \mathcal{H} \\ \rho \in [0, r']}} R(h_f) - \hat{R}_{S, \rho}(f) - \frac{4}{\rho} R_m(\tilde{f}) - \sqrt{\frac{\log \log_2 (2r'/\rho)}{m}} - \varepsilon$$

$$\leq \sup_{\substack{f \in \mathcal{H} \\ k \geq 1}} R(h_f) - \hat{R}_{S, \rho_k}(f) - \frac{2}{\rho_k} R_m(f_i) - \sqrt{\frac{\log k}{m}} - \varepsilon = \textcircled{2}$$

$$\text{Thus, } \mathbb{P}[\textcircled{1} > 0] \leq \mathbb{P}[\textcircled{2} > 0] \leq 2e^{-2M\varepsilon^2} = \delta$$

$$\Rightarrow \mathbb{P}[\textcircled{1} \leq 0] \geq 1 - \delta.$$