

LECTURE 2: GENERALIZATION AND UNIFORM CONVERGENCE THEORY

SUPERVISED LEARNING PROBLEM:

- Data: $\{(y_i, \alpha_i)\}_{i \leq n}^n$ $y_i \in \mathbb{R}$ $\alpha_i \in \mathbb{R}^d$ iid $\sim P$
- Loss fct: $l: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$
- Goal: Fit a predictor $\hat{f}: \mathbb{R}^d \rightarrow \mathbb{R}$ that has small
Expected/population/test risk: $R(\hat{f}) := E[l(y, f(x))]$
 $(y, x) \sim P$

Empirical risk: $\hat{R}_m(f) := \frac{1}{m} \sum_{i \leq m} l(y_i, f(x_i))$

class of
models

ERM: $\hat{f}_m := \underset{f}{\operatorname{argmin}} \left\{ \hat{R}_m(f) \text{ s.t. } f \in F \right\}$

Generalization question: \rightarrow what is $R(\hat{f}_m)$?

\rightarrow how to guarantee that we do well
on new unseen data?

"Classical" idea: Uniform convergence:

If $\varepsilon_m(F) := \sup_{f \in F} |\hat{R}_m(f) - R(f)|$ is small ②

then $R(\hat{f}_m) \approx \inf_{f \in F} R(f)$



"Proof": $f_* = \underset{f \in F}{\operatorname{argmin}} R(f)$

$$R(\hat{f}_m) = R(f_*) + [R(\hat{f}_m) - \hat{R}_m(\hat{f}_m)] + [\hat{R}_m(\hat{f}_m) - \hat{R}_m(f_*)]$$

$$+ [\hat{R}_m(f_*) - R(f_*)] \leq 0$$

$$\leq R(f_*) + 2 \sup_{f \in F} |R(f) - \hat{R}_m(f)| \quad \square$$

$$\equiv: \varepsilon_m(F)$$

Rmk 1: $\hat{R}_m(f) - R(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)) - \mathbb{E}[l(y, f(x))]$

For fixed f $= O(1/\sqrt{n})$ CLT

→ here we want a bound that holds uniformly over all $f \in F$

(3)

Rmk 2: In classical statistical learning theory,
 "a model learnt from data is effective at predicting on new data" is a consequence
 of U.C.

→ It is the mechanism that guarantees generalization

→ Highly influential idea: research then consists
 in carefully crafting fit classes such that $\mathbb{E}_m(F)$ small
 \hookrightarrow constrain number of parameters OR regularize weights
 \hookrightarrow in classical SLT: models are underparametrized or
 overconstrained

→ In Deep Learning, it seems good generalization
 happens thanks to a completely different principle

⇒ More in 2 next lectures

→ UC Theory still highly informative

Also help to compare and contrast what happens in DL
 \curvearrowright dependency on different parameters

GOAL rest of lecture

* Basic definitions + tools to bound $\mathbb{E}_m(F)$

* Ex 1: $F = 2\text{-layer NNs}$

* Ex 2: $F = \text{Multi-layer NNs}$

4

Def [Rademacher complexity] • $(\mathcal{Z}, \mathbb{P})$ prob space
 • $\mathcal{G} \subseteq \mathbb{R}^{\mathcal{Z}}$ class of fits $g: \mathcal{Z} \rightarrow \mathbb{R}$

$$Rd_m(\mathcal{G}) := \mathbb{E} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{i=1}^m z_i g(z_i) \right| \right]$$

where $z_1, \dots, z_m \stackrel{iid}{\sim} \mathbb{P}$

$z_1, \dots, z_m \stackrel{iid}{\sim} \text{Unif}(\{-1, 1\})$

Rmk: In our case, $z_i = (y_i, x_i)$ $g(z_i) = l(y_i, f(x_i))$
 $\mathcal{G} = l \circ F$

Lemma 1: $\left[\frac{1}{2} Rd_m(\mathcal{G}) \leq \mathbb{E} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{i=1}^m g(z_i) - \mathbb{E} g \right| \right] \leq 2 Rd_m(\mathcal{G}) \right]$

→ will only prove this upper bound

Proof: [In this lecture, we are only interested in upper-bound]

Decoupling argument:

$$\mathbb{E}_Z \sup_g \left| \frac{1}{m} \sum_{i=1}^m g(z_i) - \mathbb{E} g \right|$$

$$= \mathbb{E}_Z \sup_g \left| \mathbb{E}_{Z'} \left[\frac{1}{m} \sum_{i=1}^m g(z_i) - g(z'_i) \right] \right|$$

$$\begin{aligned}
 &\leq \underset{\text{(Sensen's)}}{\mathbb{E}_{\substack{Z, Z' \\ \mathcal{Z}}} \sup_g \left| \frac{1}{m} \sum_{i=1}^m [g(Z_i) - g(Z'_i)] \right|} \quad (5) \\
 &= \mathbb{E}_{\substack{Z, Z' \\ \mathcal{Z}}} \sup_g \left| \frac{1}{m} \sum_{i=1}^m g(Z_i) Z_i - \frac{1}{m} \sum_{i=1}^m g(Z'_i) Z_i \right| \\
 &\leq 2 R_{d_m}(g) \quad \square \\
 &\text{(triangular inequality)}
 \end{aligned}$$

Remark: This implies $\mathbb{E}[R(\hat{f}_m)] \leq \inf_{f \in F} R(f) + 2 R_{d_m}(l \circ F)$

To show without expectation, need to show

$$\sup_{f \in F} |\hat{R}_m(f) - R(f)| \text{ concentrates around its } \mathbb{E}.$$

E.g. Markov's inequality: w.p. $\geq 1 - \delta$

$$R(\hat{f}_m) \leq \inf_{f \in F} R(f) + \frac{2}{\delta} R_{d_m}(l \circ F)$$

Can improve using concentration: e.g., $\sqrt{\log(\frac{1}{\delta})}$

if l is bounded by McDiarmid's bounded difference inequality

(6)

Right now: $Rd_m(\ell \circ F) \rightarrow$ want $Rd_m(F)$

Lemma 2: [Contraction inequality]

$\phi_1, \dots, \phi_m: \mathbb{R} \rightarrow \mathbb{R}$ λ -Lipschitz $\phi_i(0)=0$

$$\mathbb{E}_{\mathcal{Z}} \left[\sup_{f \in F} \left| \frac{1}{m} \sum_{i=1}^m \phi_i(f(z_i)) z_i \right| \right] \leq \lambda \mathbb{E}_{\mathcal{Z}} \left[\sup_{f \in F} \left| \frac{1}{m} \sum_{i=1}^m f(z_i) z_i \right| \right]$$

Proof: (skipped during class) Add in lecture notes

you can try to do it by yourself

\Rightarrow can reduce to ourselves to $\lambda = 1$

and $\phi_2 = \dots = \phi_m(t) = t$

\Rightarrow can do $\phi_1(f(z_1)) \rightarrow f(z_1)$ "by hand"

(Add it in the lecture notes)

(7)

Bound on $\mathbb{E}_m(F)$: $l(y, \cdot)$ L-lipschitz

$$\mathbb{E}[\mathbb{E}_m(F)] \leq 2 \text{Rd}_m(l \circ F)$$

$$\text{Rd}_m(l \circ F) \leq \mathbb{E} \left| \frac{1}{m} \sum_{i=1}^m l(y_i, 0) z_i \right|$$

$$+ \mathbb{E} \sup_{f \in F} \left| \frac{1}{m} \sum_{i=1}^m z_i [l(y_i, f(x_i)) - l(y_i, 0)] \right|$$

$$\leq \mathbb{E} \left[\frac{1}{m^2} \sum_{ij} z_i z_j l(y_i, 0) l(y_j, 0) \right]^{\frac{1}{2}} \stackrel{\text{≤ Jensen's}}{\leq}$$

$$+ L \mathbb{E} \sup_{f \in F} \left| \frac{1}{m} \sum_i z_i f(x_i) \right|$$

$$\leq \frac{1}{\sqrt{m}} \mathbb{E}[l(y, 0)^2]^{\frac{1}{2}} + L \cdot \text{Rd}_m(F)$$

$$\underline{\text{Summary}}: \mathbb{E}[R(f_m)] \leq \inf_{f \in F} R(f) + \frac{C}{\sqrt{m}} \mathbb{E}[l(y, 0)^2]^{\frac{1}{2}}$$

$$+ C L \cdot \text{Rd}_m(F)$$

Example 1: $F :=$ 2-layer neural nets

$$f(x; \theta) = \frac{1}{M} \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle)$$

$$x \in \mathbb{R}^d$$

$$\theta = \{(a_j, w_j)\}_{j \leq m}$$

Assume: $x \in B_2^d(C_x \sqrt{d})$ = ball of radius $C_x \sqrt{d}$

$$F = \left\{ \begin{array}{l} w_j \in B_2^d(C_w) \\ a = (a_1, \dots, a_M) \in S_{p, M} = \left\{ \frac{1}{M} \sum_{i=1}^M |a_i|^p \leq r_0^p \right\} \\ = \left\{ \|a\|_p \leq r_0 M^{1/p} \right\} \\ = B_p^M(r_0 M^{1/p}) \end{array} \right.$$

By Jensen's $S_{p, M} \subseteq S_{p', M}$ $p' < p$

$$Rd_m(F) = \mathbb{E} \sup_{a, w} \left| \frac{1}{M} \sum_{i=1}^M z_i \cdot \frac{1}{M} \sum_{\ell=1}^M a_\ell \sigma(\langle w_\ell, x_i \rangle) \right|$$

$$= \frac{1}{M} \mathbb{E} \sup_{W} \sup_{a \in S_{p, M}} \left\langle a, \frac{1}{M} \sum_{i=1}^M z_i \cdot \sigma(W x_i) \right\rangle$$

$$\sup_{\|a\|_p \leq C} \langle a, v \rangle = C \|v\|_q \quad q = \frac{p}{p-1}$$

$$= n_0 M^{\frac{1}{p}-1} \mathbb{E} \sup_{\omega} \left\| \frac{1}{m} \sum_{i=1}^m z_i \sigma(\langle \omega_{\alpha_i}, \cdot \rangle) \right\|_q$$

$$= n_0 M^{\frac{1}{p}-1} \mathbb{E} \sup_{\omega_1, \dots, \omega_M} \left(\sum_{\ell=1}^M \left| \frac{1}{m} \sum_{i=1}^m z_i \sigma(\langle \omega_{\ell}, \alpha_i \rangle) \right|^q \right)^{\frac{1}{q}}$$

$$= n_0 M^{\frac{1}{p} + \frac{1}{q} - 1} \mathbb{E} \sup_{\omega} \left| \frac{1}{m} \sum_{i=1}^m z_i \sigma(\langle \omega, \alpha_i \rangle) \right|$$

n_0 Rd of one neuron \rightarrow indep of M and p

Contraction

$$\leq n_0 L_6 \mathbb{E} \sup_{\omega} \left| \langle \omega, \frac{1}{m} \sum_i z_i \alpha_i \rangle \right|$$

$$= n_0 L_6 C_{\omega} \mathbb{E} \left\| \frac{1}{m} \sum_i z_i \alpha_i \right\|_2$$

$$\leq n_0 L_6 C_{\omega} \mathbb{E} \left[\left\| \frac{1}{m} \sum_i z_i \alpha_i \right\|_2^2 \right]^{\frac{1}{2}} \leq \frac{1}{m^2} \sum_i \mathbb{E} \|\alpha_i\|_2^2$$

$$\leq n_0 L_6 C_{\omega} \frac{1}{\sqrt{m}} \mathbb{E} [\|x_i\|^2]^{\frac{1}{2}}$$

$$\leq n_0 L_6 C_{\omega} C_x \sqrt{\frac{d}{m}}$$

10

Conclusion: $Rd_m(l \circ F) \leq \frac{C}{\sqrt{m}} + L \cdot L_g \cdot n_0 C_\omega C_\alpha \sqrt{\frac{d}{m}}$

Rank: ① Bound independent of M

↳ NNs can generalize even with $M = \infty$

as long as $\|\alpha\|_p \ll \sqrt{\frac{m}{d}}$ [Bartlett, 1996]

1st realized: # params not the right "complexity" measure ^(VCdim) → norm of the parameters.

Natural $M = \infty$: μ Prob meas on $B_2^d(C_\omega)$

$$F_p(n_0) := \left\{ f(x) = \int \sigma(\langle w, x \rangle) \alpha(w) \mu(dw) \right\}$$

$$\|\alpha\|_{L_p(\mu)} = \left(\int |\alpha(w)|^p \mu(dw) \right)^{\frac{1}{p}} \leq n_0$$

$$Rd_m(F_p(n_0)) \leq L_g n_0 C_\alpha C_\omega \sqrt{\frac{d}{m}}$$

② Independent of p : $F_p(n_0) \subset F_{p'}(n_0) \subset F_1(n_0)$

$$1 < p' < p$$

↳ strongest core for $p=1$ $\int |\alpha(w)| \mu(dw) \leq n_0$

"Convex Neural networks"

(11)

$$\overline{\mathcal{F}_1}(n_0) = \left\{ f(x) = \int \sigma(\langle w, x \rangle) v(dw) : \|v\|_{TV} \leq n_0 \right\}$$

$$v = v_+ - v_- \quad \|v\|_{TV} = v_+(J) + v_-(J) \leftarrow \begin{matrix} \text{Total variation} \\ \text{Signed measure} \end{matrix}$$

→ Will go back to this later this semester
model

Rank: p=2 \mathcal{F}_2 is a RKHS

(also talk about it later this semester)

(12)

Ex 2: Multilayer NNs

[Bartlett, Foster, Telgarsky, 2017]

(^{"spectral-based"} generalization bound for multi-layer)
(Lipschitz norm of the NN)

Here slightly simplified setting

$$f(x; W) = \sigma \circ W_L \circ \sigma \circ W_{L-1} \circ \dots \circ \sigma \circ W_1 x$$

* $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ 1-Lipschitz $\sigma(0) = 0$

* $W = (W_L, W_{L-1}, \dots, W_1)$

$$W_L \in \mathbb{R}^{1 \times d} \quad W_{L-1}, \dots, W_2 \in \mathbb{R}^{M \times M} \quad W_1 \in \mathbb{R}^{M \times d}$$

$$\textcircled{H}(b_1, \dots, b_L, \delta_1, \dots, \delta_L) = \left\{ \begin{array}{l} \|W_\ell\|_{op} \leq \delta_\ell \\ \|W_\ell\|_{1,2} \leq b_\ell \end{array} \right\}$$

$$\|W\|_{1,2} = \left\| \underbrace{\left(\|W_{\cdot,1}\|_1, \dots, \|W_{\cdot,M}\|_1 \right)}_{L_1 \text{ norm of columns.}} \right\|_2$$

$$F := \left\{ f(\cdot, W) : W \in \textcircled{H}(b_1, \dots, b_L, \delta_1, \dots, \delta_L) \right\}$$

Goal: bound $Rd_m(F)$

(13)

$X = [x_1, \dots, x_d] \in \mathbb{R}^{d \times m}$ data

$\mathcal{F}_m(X) := \{f_m(W) : W \in \mathbb{W}\} \subseteq \mathbb{R}^m$

$$f_m(W) = (f(x_1; W), \dots, f(x_m; W))$$

Approach:

Def [Conditional Rademacher complexity]

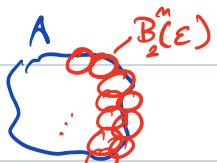
$$Rd_m(\mathcal{F}_m(X)) := \mathbb{E}_z \sup_{W \in \mathbb{W}} \frac{1}{m} \sum_{i=1}^m z_i f(x_i; W)$$

$$= \mathbb{E}_z \sup_{W \in \mathbb{W}} \frac{1}{m} \langle f_m(W), z \rangle$$

Rmk: $Rd_m(\mathcal{F}) = \mathbb{E}_X Rd_m(\mathcal{F}_m(X))$

Def: [Covering number] $A \subseteq \mathbb{R}^m$

$\mathcal{N}(A; \varepsilon) := \min \# \text{ of } l_2 \text{ balls of radius } \varepsilon \text{ to cover } A$



Lemma 3: Assume $0 \in F_m(X) \subseteq [0, 1]^n$

$$Rd_m(F_m(X)) \leq \inf_{\alpha \geq 0} \left\{ \frac{2\alpha}{\sqrt{m}} + \frac{24}{m} \int_{\alpha}^{\sqrt{m}} \sqrt{\log N(F_m(X), \varepsilon)} d\varepsilon \right\}$$

↑ "Dudley integral" ↑ "metric entropy"

Proof: Classical "chaining" argument [Vershynin]
 → will add in the notes

Goal: Bound $N(F_m(X), \varepsilon)$, i.e., construct ε -covering

$$F_m(X) := \{ f_m(W) : W \in \mathbb{W} = \mathbb{H}_1 \times \dots \times \mathbb{H}_L \}$$

$$\mathbb{H}_l = \{ W_l : \|W_l\|_{op} \leq s_l, \|W_l\|_{1,2} \leq b_l \}$$

$$f_m(W) = \sigma \circ W_L \circ \dots \circ \sigma \circ A_1 X \in \mathbb{R}^n$$

↙ output of intermediate
layers

$$f_m^{(l)}(W^l) = \sigma \circ W_l \circ \dots \circ \sigma \circ A_1 X \in \mathbb{R}^{M \times n}$$

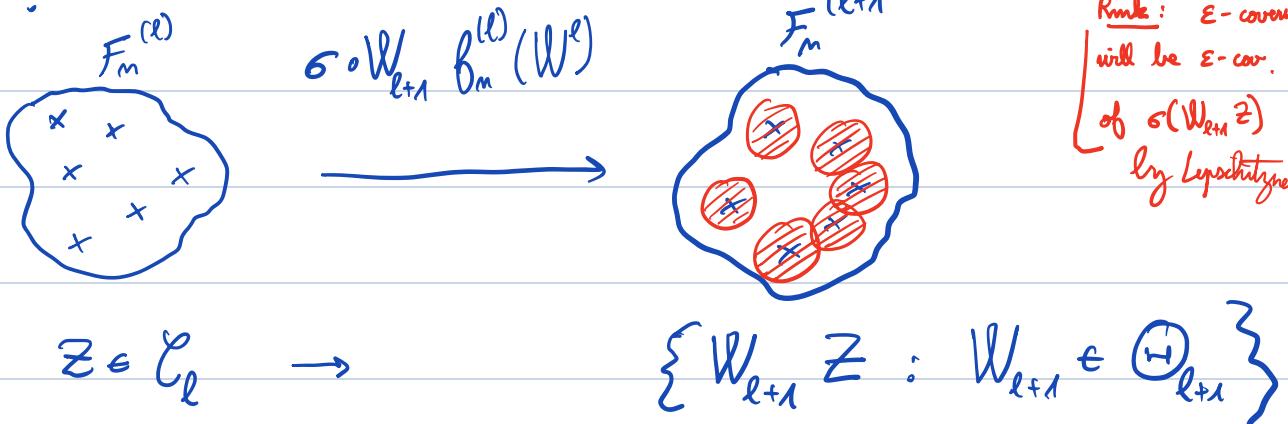
$$f_m^{(l+1)}(W^{l+1}) = \sigma \circ W_{l+1} f_m^{(l)}(W^l)$$

(15)

Idea: Construct covering \mathcal{C}_L by constructing sequentially coverings \mathcal{C}_l δ_l covering of

$$\mathcal{F}_m^{(l)}(X) := \left\{ f_m^{(l)}(W^l) : W^l \in \mathbb{H}^l \right\} \subseteq \mathbb{R}^{M \times m}$$

Given \mathcal{C}_l δ_l -covering of $\mathcal{F}_m^{(l)}$, we construct δ_{l+1} -covering \mathcal{C}_{l+1} as follows:



$B_{l+1}(z; \varepsilon_{l+1})$ ε_{l+1} covering for each of these sets \rightarrow union gives \mathcal{C}_{l+1}

$$|\mathcal{C}_{l+1}| \leq \sum_{z \in \mathcal{C}_l} |B_{l+1}(z; \varepsilon_{l+1})| \leq |\mathcal{C}_l| \max_{z \in \mathcal{C}_l} |B_{l+1}(z; \varepsilon_{l+1})|$$

$$\delta_{l+1} \leq \sup_{W_{l+1} \in \mathbb{H}_{l+1}} \|W_{l+1}\|_{op} \delta_l + \varepsilon_{l+1}$$

$$\boxed{\begin{aligned} f &= W_{l+1} \tilde{f} \rightarrow \delta_l \text{ close to } z \in \mathcal{C}_l \\ &= W_{l+1} (\tilde{f} - z) + W_{l+1} z \end{aligned}}$$

(16)

$$\text{Hence } \delta_{l+1} \leq \delta_{l+1} \delta_l + \varepsilon_{l+1}$$

$$\delta_L \leq \sum_{l=1}^L \prod_{j=l+1}^L \delta_j \cdot \varepsilon_l$$

$$\sup_{z \in C_\ell} |B_{\ell+1}(z; \varepsilon_{\ell+1})| \leq \sup_{f_m \in \mathcal{F}_m^{(\ell)}} N(T_{\ell+1}(f_m); \varepsilon_{\ell+1})$$

$$\text{where } T_\ell(f_m) := \{W_\ell f_m : W_\ell \in \mathbb{H}_\ell\}$$

$$\|W_\ell\|_{op} \leq \rho_\ell$$

$$\|W_\ell\|_{1,2} \leq b_\ell$$

$$\log N(\mathcal{F}_m; \varepsilon) \leq \sum_{l=1}^L \sup_{f_m \in \mathcal{F}_m^{(\ell)}} \log N(T_\ell(f_m); \varepsilon_\ell)$$

$$\varepsilon := \sum_{l=1}^L \prod_{j=l+1}^L \delta_j \cdot \varepsilon_l$$

Lemma 4: $T(z; b) := \{Az : \|A\|_{1,2} \leq b\}$

$$A \in \mathbb{R}^{M \times M}$$

$$z \in \mathbb{R}^{M \times N}$$

$$\log N(T(z; b); \varepsilon) \leq 2 \left(\frac{b \|z\|_F}{\varepsilon} \right)^2 \log(2M)$$

Will prove later if I have the time.

(This lemma is the reason for $\|W_\ell\|_{1,2} \leq b_\ell$ constraint)

$$R_\ell := \sup_{f_m \in F_m^{(\ell-1)}} \log N(T_\ell(f_m), \varepsilon_\ell) \leq 2 \left(\frac{b_\ell}{\varepsilon_\ell} \sup_{f_m \in F_m^{(\ell-1)}} \|f_m\|_F \right)^2 \log(2M)$$

$$f_m^{(\ell-1)}(W^{(\ell-1)}) = \sigma \circ W_{\ell-1} \circ \sigma \dots \circ \sigma \circ W_1 X \in \mathbb{R}^{M \times m}$$

$$\begin{aligned} \|f_m^{(\ell-1)}\|_F &\leq \|W_{\ell-1}\|_{op} \dots \|W_1\|_{op} \|X\|_F \\ &\leq \|X\|_F \cdot \prod_{i=1}^{\ell-1} \delta_i \end{aligned}$$

$$X_\ell \leq 2 \left(\frac{b_\ell}{\varepsilon_\ell} \prod_{i=1}^{\ell-1} \delta_i \right)^2 \|X\|_F^2 \log(2M)$$

$$\log N(F_m(X), \varepsilon) \leq 2 \sum_{\ell=1}^L \left(\frac{b_\ell}{\varepsilon_\ell} \prod_{i=1}^{\ell-1} \delta_i \right)^2 \|X\|_F^2 \log(2M)$$

$$\begin{aligned} \varepsilon &= \sum_{\ell=1}^L \prod_{j=\ell+1}^L \delta_j \varepsilon_\ell \\ &=: \alpha_\ell \cdot \varepsilon \end{aligned}$$

$$\sum_{\ell=1}^L \alpha_\ell = 1$$

optimizing over them!

$$\log N(F_m, \varepsilon) \leq \frac{2 \|X\|_F^2 \log 2M}{\varepsilon^2} \sum_{l=1}^L \left(\frac{b_l}{\alpha_l} \prod_{j \neq l} \alpha_j \right)^2$$

$$\left(\prod_{j=1}^L \alpha_j \right)^2 \sum_{l=1}^L \left(\frac{b_l}{\alpha_l \delta_l} \right)^2$$

$$c_l := \frac{b_l}{\delta_l} \Rightarrow \min \sum_{l=1}^L \left(\frac{c_l}{\alpha_l} \right)^2 \text{ s.t. } \sum_{l=1}^L \alpha_l = 1$$

$$\alpha_l \geq 0$$

Convinc yourself: $\alpha_l = \frac{c_l^{2/3}}{\sum_{l'} c_{l'}^{2/3}}$ [Holder's inequality]

$$Q := \left(\sum_{l=1}^L \left(\frac{b_l}{\delta_l} \right)^{2/3} \right)^{3/2} \left(\prod_{i=1}^L \alpha_i \right) \|X\|_F \sqrt{\log 2M}$$

$$\Rightarrow \boxed{\log N(F_m; \varepsilon) \leq 2 \frac{Q^2}{\varepsilon^2}}$$

$$Rd_m(F_m) \leq \frac{2\alpha}{\Gamma_m} + \frac{2Q}{m} \int_{\alpha}^{\Gamma_m} \sqrt{\log N(F_m; \varepsilon)} d\varepsilon$$

$$\leq \frac{2\alpha}{\Gamma_m} + \frac{2Q}{m} Q \sqrt{2} \int_{\alpha}^{\Gamma_m} \frac{1}{\varepsilon} d\varepsilon \rightarrow \log \frac{\Gamma_m}{\varepsilon}$$

$$\alpha = \frac{1}{m}$$

$$\leq \frac{100}{m} Q \log(m)$$

$$\text{Conclusion: } R_{dm}(\mathcal{F}_m(X)) \leq C \cdot R_0 \cdot \|X\|_F \frac{\log m}{m} \cdot \sqrt{\log(2M)}$$

where "spectral complexity" of the network is given by

$$R_0 := \left(\sum_{l=1}^L \left(\frac{b_l}{\delta_l} \right)^{2/3} \right)^{3/2} \prod_{i=1}^L \delta_i$$

$$\|W_l\|_{op} \leq \delta_l \quad \|W_l\|_{1,2} \leq b_l$$

- RMK:
- Mild dependency on $\log(M)$
 - scales with "Lipschitz cste" $\prod_{i=1}^L \delta_i$ of network
 - dependency on # layers $\propto L^{3/2}$ (outside Lipschitz constant)

EX: $\|x_i\|_2 \asymp C\sqrt{d}$ W_l i.i.d. $\frac{1}{\sqrt{d}}$ sub-Gaussian

$$\|X\|_F \asymp \sqrt{md}$$

$$\|W_l\|_{op} \asymp C \quad \|W_l\|_{1,2} \asymp M$$

$$R_{dm}(\mathcal{F}_m) \asymp \sqrt{\frac{dM^2}{m}} L^{3/2} C^L$$

some issue I
talked last week

For $L = O(1)$: requires $m \gg M^2 d \gg \underbrace{M^2 \vee Md}_{\# \text{parameters}}$

Can these VC_D explain why NNs generalize?

(very different opinions, especially about right complexity measure)

- If NNs verify these bounds \rightarrow explain why it generalizes
- Some correlate between real & gen bounds
 \hookrightarrow lots of empirical studies: overall more seem highly predictive of performance
- "Requires rethinking gene" \rightarrow when interpolate break down

\hookrightarrow above bound

- exponential dependency on L
- "cons" • no assumption on the distribution of the data
- need # samples \gg # parameters

Proof of Lemma 9

$$T(Z; b) := \{ AZ : \|A\|_{1,2} \leq b \}$$

$$AZ = \underbrace{(A \odot C)}_B \hat{Z}$$

has normalized
rows

$$Z = \begin{pmatrix} c_1 & \dots & c_m \end{pmatrix} \hat{Z}$$

$$c_i = \|z_{i,:}\|_2$$

$$C = \begin{pmatrix} c_1 & c_2 & \dots & c_m \\ \vdots & \vdots & & \vdots \\ c_1 & c_2 & & c_m \end{pmatrix} \quad M \downarrow M$$

\overbrace{C}^M

$$A \operatorname{diag}(c_1 \dots c_m) = A \odot C$$

$$\|B\|_1 = \sum_{ij} |B_{ij}| = \langle C, |A| \rangle \leq \|C\|_{\infty, 2} \|A\|_{1,2}$$

$$(\|M\|_{p,q} = \left\| (\|M_{:,1}\|_p, \dots, \|M_{:,n}\|_p) \right\|_q)$$

$$\langle X, Y \rangle \leq \|X\|_{p,\diamond} \|Y\|_{q,\wedge} \quad \frac{1}{p} + \frac{1}{q} = 1$$

$$\frac{1}{\rho} + \frac{1}{\eta} = 1$$

$$\|C\|_{\infty,2} = \left\| (|c_1|, \dots, |c_n|) \right\|_2 = \left(\sum_i \|z_i\|_2^2 \right)^{\frac{1}{2}} \quad (22)$$

$$= \|Z\|_F.$$

$$\text{Hence } \|B\|_1 \leq \|Z\|_F \cdot \|A\|_{1,2} \leq \|Z\|_F \cdot b =: t$$

$$T(z; b) \subseteq \{Bz : \|B\|_1 \leq t\}$$

Lemma 5: Let $u = \sum_{i=1}^p \lambda_i v_i$ vector is v_i , u , $\lambda_i \geq 0$

$$\Rightarrow \exists k_1, \dots, k_p \geq 0 \text{ integers } \sum_{i=1}^p k_i = k$$

s.t. $\left\| u - \frac{\|\lambda\|_1}{k} \sum_{i=1}^p k_i v_i \right\|_2^2 \leq \frac{\|\lambda\|_1^2}{k} \max_{i \leq p} \|v_i\|_2^2$

Remark: Result does not depend on the dimension of v_i 's
(hold for $\dim(v_i) = \infty$)

Proof: [Add a proof in the notes.]

Corollary: $S(v_1, \dots, v_p; r) := \{u = \sum \lambda_i v_i : \|\lambda\|_1 \leq r\}$

$$N(S; \varepsilon) \leq (2p)^k$$

$$\text{where } \varepsilon^2 = \frac{\tau^2}{k} \max_i \|w_i\|_2^2$$

We have $T(Z, b) \subseteq \{B\hat{Z} : \|B\|_1 \leq b\}$

$$B\hat{Z} = \sum_{i,j=1}^M B_{ij} [e_i e_j^\top \hat{Z}] = \sum_{i,j=1}^M B_{ij} v_{ij} \in \mathbb{R}^{M \times n}$$

Above corollary gives $\mathcal{N}(\{B\hat{Z} : \|B\|_1 \leq b\}; \varepsilon) \leq (2M^2)^k$

$$\begin{aligned} \varepsilon^2 &= \frac{b^2}{k} \max_{ij} \|e_i e_j^\top \hat{Z}\|_F^2 \\ &= \|\hat{Z}\|_F^2 = 1 \end{aligned}$$

$$\Rightarrow \log \mathcal{N}(T(Z, b), \varepsilon) \leq k \log (2M)^2$$

$$\leq 2 \left(\frac{b}{\varepsilon} \right)^2 \log (2M)$$

□