# Neural Net Classifiers and their Rademacher Complexities

$$x \in \mathcal{X} \qquad g_f(x) = \text{sgn} f(x)$$
$$f : \mathcal{X} \to \mathbb{R}, \text{ in } \mathcal{F}$$

- linear classifiers: $\mathcal{X} \subseteq \mathbb{R}^d$

$$f(x) = \langle w, x \rangle \qquad w \in \mathbb{R}^d$$

(can cover affine $\quad f(x) = \langle w, x \rangle + b$
by adding an all-1 coordinate to $x$:

$$x \mapsto \binom{x}{1} \in \mathbb{R}^{d+1})$$

$$\mathcal{F} := \left\{ x \mapsto \langle w, x \rangle : \quad \|w\| \le B \right\}$$
$$R_n(\mathcal{F}(x^n)) = \frac{1}{n} \mathbb{E}_\varepsilon \left[ \sup_{\|w\| \le B} \left| \sum_{i=1}^n \varepsilon_i \langle w, x_i \rangle \right| \right]$$
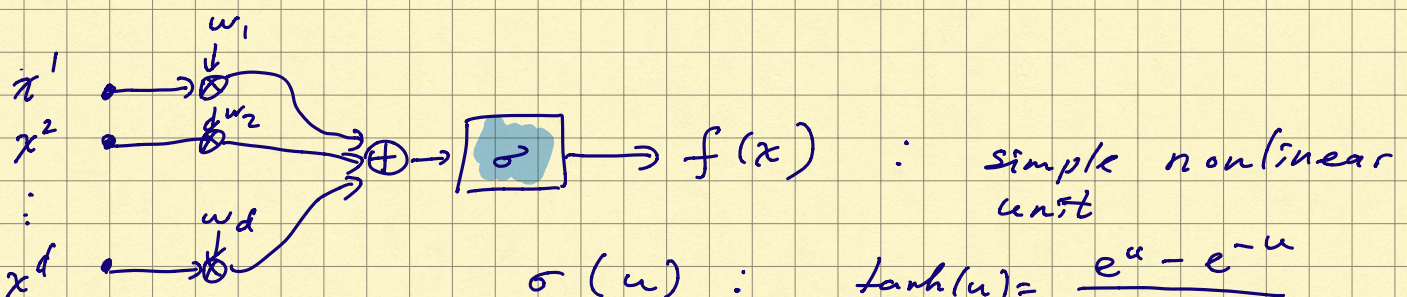$$\le \frac{B}{n} \sqrt{\sum_{i=1}^n \|x_i\|^2}$$

$$\mathcal{X} = \left\{ x \in \mathbb{R}^d : \|x\| \le R \right\} \implies R_n(\mathcal{F}(x^n)) \le \frac{BR}{\sqrt{n}}.$$

- simple nonlinearity : single neuron

$$f(x) = \sigma \left( \langle w, x \rangle \right) \qquad x, w \in \mathbb{R}^d$$

where $\sigma : \mathbb{R} \to \mathbb{R}$ is a continuous fcn,
$\sigma(0) = 0$, Lipschitz continuous: $|\sigma(u) - \sigma(v)| \le L |u - v|$.

$$x = (x^1, x^2, \dots, x^d)^T$$



$\sigma(u)$ : Simple nonlinear unit

$$\tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$
$$(u)_+ = \max(0, u), \text{ ReLU}$$

$$\mathcal{F}_\sigma := \left\{ x \mapsto \sigma\left( \sum_{j=1}^{d} w_j x^j \right) : \|w\| \leq B \right\}$$

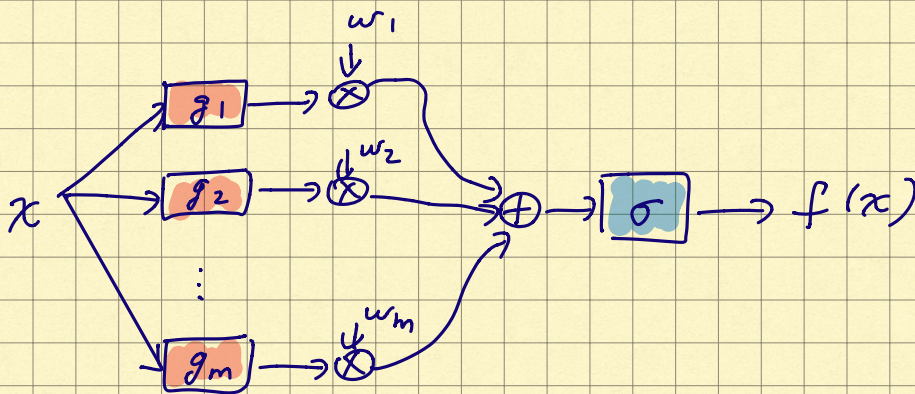$$\mathcal{F}_\sigma = \sigma \circ \mathcal{F} \qquad (\mathcal{F} : \text{lin. classifiers})$$

Contraction principle:

$$R_n(\mathcal{F}_\sigma) \leq 2L\, R_n(\mathcal{F}) \leq \frac{2LBR}{\sqrt{n}}$$

$$\text{if} \quad \mathcal{X} = \{ x \in \mathbb{R}^d : \|x\| \leq R \}$$

— neural nets

$$\mathcal{G} : \quad \text{base classifiers} \quad g : \mathcal{X} \to \mathbb{R}$$



$$\mathcal{F}_1 := \left\{ x \mapsto \sigma\left( \sum_{j=1}^{m} w_j\, g_j(x) \right) : \begin{array}{l} m \in \mathbb{N} \\ \|w\|_1 = \sum_{j=1}^{m} |w_j| \leq B \\ g_1, \ldots, g_m \in \mathcal{G} \end{array} \right\}$$

$$\mathcal{F}_1 = \sigma \circ \left( B \cdot \mathrm{absconv}(\mathcal{G}) \right)$$

Contraction principle again:  $\sigma(0) = 0$
 $\sigma \quad L - \text{Lip.}$

$$R_n(\mathcal{F}_1) \leq 2L \cdot R_n(B \cdot \mathrm{absconv}(\mathcal{G}))$$

$$= 2LB \cdot R_n(\mathrm{absconv}(\mathcal{G}))$$

$$= 2LB \cdot R_n(\mathcal{G}).$$
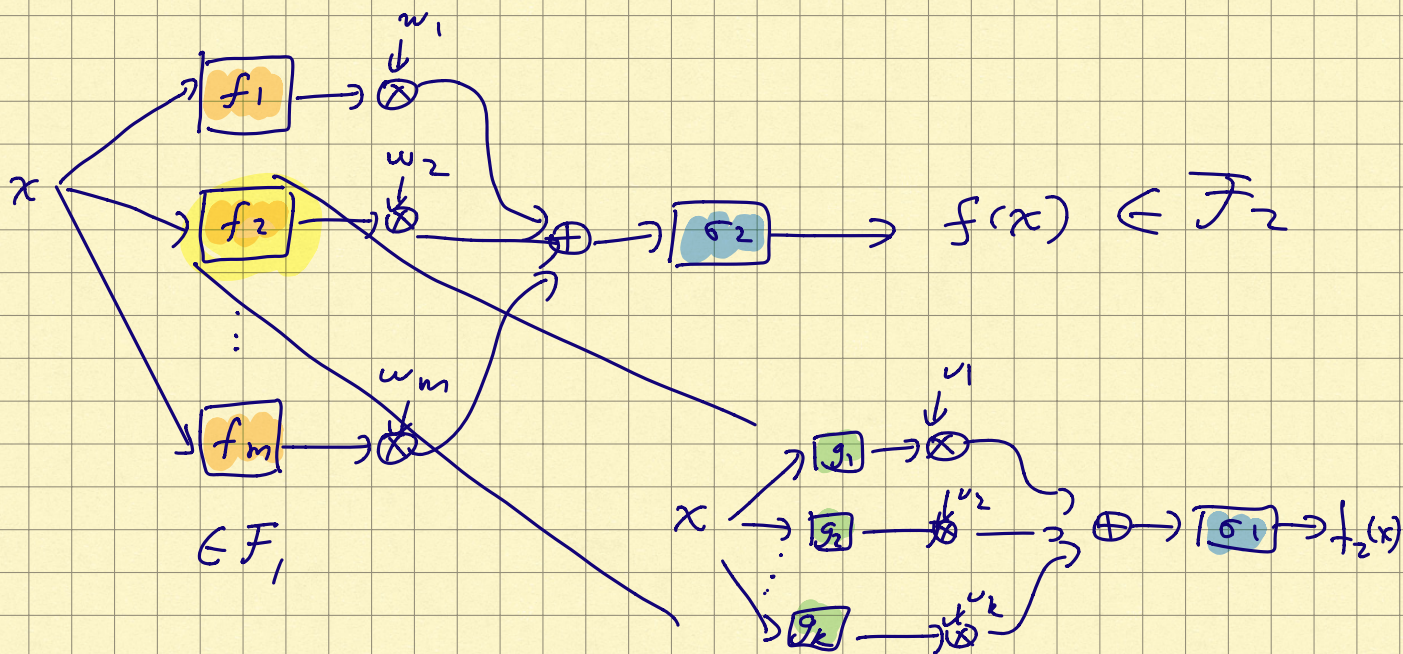
— adding layers  (making neural nets deeper)

$f_1, \ldots, f_m \in F_1$ above

$$f_j(x) = \sigma_1 \left( \sum_{k=1}^{m_j} w_k^j g_k^j(x) \right)$$

where $m_j \in \mathbb{N}$; $\|w^j\|_1 = \sum_{k=1}^{m_j} |w_k^j| \leq B_1$

$$g_1^j, \ldots, g_{m_j}^j \in G$$

$$f(x) = \sigma_2 \left( \sum_{j=1}^{m} w_j f_j(x) \right) \qquad \|w\|_1 \leq B_2$$



$$\mathcal{F}_2 = \sigma_2 \circ B_2 \, \text{absconv} (\mathcal{F}_1)$$

$$R_n(\mathcal{F}_2) \leq 2 L_2 B_2 \cdot R_n(\mathcal{F}_1)$$

$$\leq 2 L_2 B_2 \cdot 2 L_1 B_1 \cdot R_n(G)$$

$$= 2^2 L_2 L_1 B_2 B_1 \cdot R_n(G)$$

– add more layers: $\quad \mathcal{F}_j = \sigma_j \circ B_j \, \text{absconv}(\mathcal{F}_{j-1})$

$$\cdots \qquad \sigma_0 = G$$

$$j = 1, \ldots, \ell$$

$$\ell - \# \text{ layers}$$

$$R_n(\mathcal{F}_\ell) \leq 2^\ell \cdot \prod_{j=1}^{\ell}(L_j B_j) \cdot R_n(\mathcal{G})$$

grows as $\exp(\ell)$

(recursive use of contraction principle: K-P. 2002)

Can we do better? $2^\ell \rightarrow \sqrt{\ell} \Leftarrow$

$\prod_{j=1}^{\ell}(L_j B_j)$ still present

- Bartlett - Foster - Telgarsky (2017)
- Golowich - Rakhlin - Shamir (2017)

GRS : peeling layer-by-layer + log exp trick

$$\mathcal{F}_j = \sigma_j \circ B_j \text{ absconv}(\mathcal{F}_{j-1}) \qquad j=1,\ldots,\ell$$
$$\mathcal{F}_0 = \mathcal{G}$$

<u>Step 1</u>

$$R_n(\mathcal{F}_\ell) = R_n(\mathcal{F}_\ell(x^n)), \quad x^n \text{ fixed}$$

$$R_n(\mathcal{F}_\ell) = \frac{1}{n} \mathbb{E}_\varepsilon\left[\sup_{f \in \mathcal{F}_\ell} \left| \sum_{i=1}^{n} \varepsilon_i f(x_i)\right|\right]$$

$$= \frac{1}{\lambda n} \mathbb{E}_\varepsilon\left[\log \exp\left(\lambda \sup_{f \in \mathcal{F}_\ell} \left| \sum_{i=1}^{n} \varepsilon_i f(x_i)\right|\right)\right]$$

$\lambda > 0$:
to be tuned

$$\leq \frac{1}{\lambda n} \log \mathbb{E}_\varepsilon\left[\sup_{f \in \mathcal{F}_\ell} \exp\left(\lambda \left| \sum_{i=1}^{n} \varepsilon_i f(x_i)\right|\right)\right]$$

Note: $G(u) := e^{\lambda u}$ $(\lambda > 0)$ is convex, nondecreasing

Step 2

Lemma 1 (GRS) Let $G$ be a convex, nondecreasing fcn $\mathbb{R} \to \mathbb{R}$; let $\bar{F}$ be of the form

$$\bar{F} = \sigma \circ B \text{ absconv}(\bar{F}')$$

for $\sigma: \mathbb{R} \to \mathbb{R}$, $\sigma(0) = 0$, $L$-Lip. Then

$$\mathbb{E}_\varepsilon \left[ G\left( \sup_{f \in \bar{F}} \left| \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right) \right]$$

$$\leq 2 \, \mathbb{E}_\varepsilon \left[ G\left( L B \cdot \sup_{f' \in \bar{F}'} \left| \sum_{i=1}^n \varepsilon_i f'(x_i) \right| \right) \right]$$

Let's apply lemma 1 to $\bar{F}_\ell = \sigma_\ell \circ B_\ell \, \text{absconv}(\bar{F}_{\ell-1})$

$$G(u) = e^{\lambda u}$$

$$\mathbb{E}_\varepsilon \left[ \exp\left( \lambda \sup_{f \in \bar{F}_\ell} \left| \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right) \right]$$

$$\leq 2 \, \mathbb{E}_\varepsilon \left[ \exp\left( \lambda L_\ell B_\ell \cdot \sup_{f \in \bar{F}_{\ell-1}} \left| \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right) \right]$$

$$\vdots$$

$$\leq 2^\ell \, \mathbb{E}_\varepsilon \left[ \exp\left( \lambda \prod_{j=1}^\ell (L_j B_j) \cdot \sup_{g \in \mathcal{G}} \left| \sum_{i=1}^n \varepsilon_i g(x_i) \right| \right) \right]$$

Step 3 $\qquad M := \prod_{j=1}^\ell (L_j B_j)$

$$\mathbb{E}_\varepsilon \left[ \exp\left( \lambda M \sup_{g \in \mathcal{G}} \left| \sum_{i=1}^n \varepsilon_i g(x_i) \right| \right) \right] \leq \, ?$$

$$A := \{ (g(x_1), \ldots, g(x_n)) : g \in \mathcal{G} \}$$

<u>Lemma 2</u> (GRS) Let $A \subseteq \mathbb{R}^n$ be a bdd set. Then, $\forall \lambda > 0$,

$$\mathbb{E}_\varepsilon \left[ \exp\left( \lambda \sup_{a \in A} \left| \sum_{i=1}^n \varepsilon_i a_i \right| \right) \right]$$

$$\leq \exp\left( \frac{\lambda^2}{2} \sum_{i=1}^n \sup_{a \in A} |a_i|^2 \right) \exp\left( \lambda n R_n(A) \right)$$

<u>Proof idea</u>   $U(\varepsilon_1, \dots, \varepsilon_n) := \sup_{a \in A} \left| \sum_{i=1}^n \varepsilon_i a_i \right|$

$U(\varepsilon_1, \dots, \varepsilon_i, \dots, \varepsilon_n) - U(\varepsilon_1, \dots, -\varepsilon_i, \dots, \varepsilon_n) \leq \sup_{a \in A} |a_i|$

— mimic proof of McDiarmid, to bound

$$\mathbb{E}\left[ e^{\lambda U} \right] . \qquad\qquad \boxdot$$

Let $A := G_0(x^n)$ $\qquad\qquad M := \prod_{j=1}^\ell (L_j B_j)$

$$\mathbb{E}_\varepsilon \left[ \exp\left( \lambda M \cdot \sup_{g \in G} \left| \sum_{i=1}^n \varepsilon_i g(x_i) \right| \right) \right]$$

$$\leq \exp\left( \frac{\lambda^2 M^2}{2} \cdot \sum_{i=1}^n \sup_{g \in G} |g(x_i)|^2 \right) \cdot \exp\left( \lambda M n R_n(G_0(x^n)) \right)$$

<u>Step 4</u>

$$R_n(\mathcal{F}_\ell) \leq \frac{1}{\lambda n} \log \mathbb{E}_\varepsilon \left[ \exp\left( \lambda M \cdot \sup_{g \in G} \left| \sum_{i=1}^n \varepsilon_i g(x_i) \right| \right) \right]$$

$$\leq \frac{1}{\lambda n} \left( \ell \log 2 + \frac{\lambda^2 M^2}{2} \sum_{i=1}^n \sup_{g \in G} |g(x_i)|^2 + \lambda M n R_n(G) \right)$$

$$= M R_n(G) + \lambda \cdot \frac{M^2}{2n} \sum_{i=1}^n \sup_{g \in G} |g(x_i)|^2 + \frac{1}{\lambda} \cdot \frac{\ell \log 2}{n}$$

$$\inf_{\lambda \geq 0} \left\{ \frac{a}{\lambda} + b\lambda \right\} = 2\sqrt{ab} \qquad (a, b \geq 0)$$

min. over $\lambda \geq 0$ :

$$\min_{\lambda \geq 0} \quad = \quad \frac{1}{n} \cdot 2 \sqrt{\ell \log 2 \cdot \frac{M^2}{2} \left( \sum_{i=1}^{n} \sup_{g \in \mathcal{G}} |g(x_i)|^2 \right)}$$

$$= \frac{M}{n} \sqrt{2 \, \ell \cdot \log 2 \cdot \sum_{i=1}^{n} \sup_{g \in \mathcal{G}} |g(x_i)|^2}$$

$$\therefore \quad R_n(\mathcal{F}_\ell) \leq \prod_{j=1}^{\ell} (L_j B_j) \cdot R_n(\mathcal{G}(x^n))$$

$$+ \frac{1}{n} \prod_{j=1}^{\ell} (L_j B_j) \underbrace{\sqrt{\ell \log 4 \cdot \sum_{i=1}^{n} \sup_{g \in \mathcal{G}} |g(x_i)|^2}}_{O(\sqrt{\ell n})}$$

(check constant in front in $\sqrt{\ell \dots}$ term).

— take $L_j, B_j = 1 \quad \forall j$ ———

$$R_n(\mathcal{F}_\ell) \leq R_n(\mathcal{G}) + \frac{C}{n} \sqrt{\ell \cdot \sum_{i=1}^{\tilde{n}} \sup_{g \in \mathcal{G}} |g(x_i)|^2}$$

$$\lesssim R_n(\mathcal{G}) + C' \sqrt{\ell/n} .$$

· take $L_1 = \dots = L_\ell = 1$

· $B_j$ : largest $\ell_1$ norm of weights in layer $j$