

Binary Classification, Part 4

① AdaBoost (Freund-Schapire, 1997)

\mathcal{G} : base classifiers $g: \mathcal{X} \rightarrow \{-1, +1\}$ (weak learners)

$(X_1, Y_1), \dots, (X_n, Y_n)$ iid on $\mathcal{X} \times \{-1, +1\}$

Goal: return $\hat{f}_n \in \text{conv}(\mathcal{G})$ [use $\text{sgn} \hat{f}_n$ to classify]

$$\hat{f}_n(x) = \frac{\sum_{k=1}^K \alpha_k g_k(x)}{\sum_{k=1}^K \alpha_k}$$

where K is fixed (# iterations) and $\alpha_k \geq 0$ and $g_k \in \mathcal{G}$ are determined iteratively from data.

Note: can use $\text{sgn} \left(\sum_{k=1}^K \alpha_k g_k(x) \right)$ to classify

AdaBoost

• init: $w^{(1)} = (w_1^{(1)}, \dots, w_n^{(1)})$, $w_i^{(1)} = 1/n \quad \forall i$

• for $k=1, \dots, K$ do:

– $g_k := \underset{g \in \mathcal{G}}{\text{argmin}} e_k(g)$ where $e_k(g) := \sum_{i=1}^n w_i^{(k)} \mathbb{1}_{\{Y_i \neq g(X_i)\}}$

$$e_k := e_k(g_k)$$

// standing assumption: $e_k \leq 1/2$

– $w^{(k)} \rightarrow w^{(k+1)}$ update

$$\alpha_k := \frac{1}{2} \log \frac{1-e_k}{e_k}$$

// $e_k \leq 1/2 \Leftrightarrow 1-e_k \geq e_k \quad \alpha_k \geq 0$

$$w_i^{(k+1)} = \frac{w_i^{(k)} e^{-\alpha_k Y_i g_k(X_i)}}{\sum_k}$$

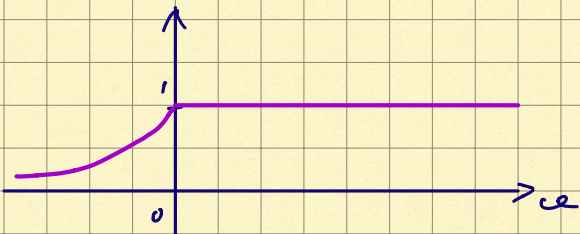
$\forall k$: norm.

$$\| e^{-\alpha_k Y_i g_k(x_i)} = \begin{cases} e^{\alpha_k} & \text{if } Y_i \neq g_k(x_i) \\ e^{-\alpha_k} & \text{if } Y_i = g_k(x_i) \end{cases}$$

• return $\hat{f}_n(x) := \frac{\sum_{k=1}^K \alpha_k g_k(x)}{\sum_{k=1}^K \alpha_k}$

Preview: • w.l.p., $L(\hat{f}_n) \leq \prod_{k=1}^K 2\sqrt{e_k(1-e_k)}$

• key analysis tool: $\varphi(u) = 1 \wedge e^u \text{ smin} \{1, e^u\}$



1-Lip., $0 \leq \varphi(\cdot) \leq 1$

$\gamma \geq \gamma' > 0 \Rightarrow$

$\varphi(\cdot/\gamma) \geq \varphi(\cdot/\gamma')$

• result due to Koltchinskii - Panchenko (2002)

Lemma $\frac{1}{n} \sum_{i=1}^n \exp\left(-Y_i \sum_{k=1}^K \alpha_k g_k(x_i)\right) = \prod_{k=1}^K 2\sqrt{e_k(1-e_k)}$

Proof $\exp\left(-Y_i \sum_{k=1}^K \alpha_k g_k(x_i)\right)$

$= \prod_{k=1}^K \exp(-\alpha_k Y_i g_k(x_i))$

$= \prod_{k=1}^K \frac{w_i^{(k+1)}}{w_i^{(k)}} \stackrel{\text{M}_k}{\sim}$

$= \prod_{k=1}^K \frac{w_i^{(k+1)}}{w_i^{(k)}} \cdot \prod_{k=1}^K \mathbb{E}_k$

$= \frac{w_i^{(K+1)}}{w_i^{(1)}} \cdot \prod_{k=1}^K \mathbb{E}_k$

$= n w_i^{(K+1)} \cdot \prod_{k=1}^K \mathbb{E}_k$

$$\frac{1}{n} \sum_{i=1}^n \exp\left(-\gamma_i \sum_{k=1}^K \alpha_k g_k(x_i)\right) = \underbrace{\sum_{i=1}^n w_i^{(k+1)}}_{=1} \cdot \prod_{k=1}^K \zeta_k$$

where $\zeta_k = 2\sqrt{e_k(1-e_k)}$ (can be shown from defn.) □

Per fcn $\varphi(u) = 1 \wedge e^u \leq e^u$

By K.-P. adaptive margin bound, w.p. $\geq 1-\delta$,

$$L(\hat{f}_n) \leq \inf_{\gamma \in (0,1]} \left\{ A_{\varphi(\cdot/\gamma),n}(\hat{f}_n) + \frac{C}{\gamma} \mathbb{E} R_n(G(X^n)) + C \sqrt{\frac{\log \log(1/\gamma)}{n}} \right\} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

\downarrow
 take any $0 < \gamma \leq 1$

Fix $\gamma \in (0,1]$

$$A_{\varphi(\cdot/\gamma),n}(\hat{f}_n) = \frac{1}{n} \sum_{i=1}^n \varphi\left(-\frac{1}{\gamma} \gamma_i \hat{f}_n(x_i)\right)$$

$$= \frac{1}{n} \sum_{i=1}^n \varphi\left(-\frac{1}{\gamma} \gamma_i \frac{\sum_{k=1}^K \alpha_k g_k(x_i)}{\sum_{k=1}^K \alpha_k}\right)$$

$$\leq \frac{1}{n} \sum_{i=1}^n \varphi\left(-\gamma_i \sum_{k=1}^K \alpha_k g_k(x_i)\right) \quad \text{for a properly chosen } \gamma$$

$$\left(\text{take } \gamma = 1 \wedge \frac{1}{\sum_{k=1}^K \alpha_k} \leq \frac{1}{\sum_{k=1}^K \alpha_k} \right)$$

$$\leq \frac{1}{n} \sum_{i=1}^n \exp\left(-\gamma_i \sum_{k=1}^K \alpha_k g_k(x_i)\right) = \prod_{k=1}^K 2\sqrt{e_k(1-e_k)}$$

$$\therefore L(\hat{f}_n) \leq \prod_{k=1}^K 2\sqrt{e_k(1-e_k)} + C \underbrace{\left(1 \vee \sum_{k=1}^K \alpha_k\right)}_{1/\gamma} \mathbb{E} R_n(G(x^n))$$

$$+ C \sqrt{\frac{\log \log \left(1 \vee \sum_{k=1}^K \alpha_k\right)}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}} \quad \text{w.p.} \geq 1-\delta$$

for $\gamma = 1 \wedge \frac{1}{\sum_{k=1}^K \alpha_k}$. □

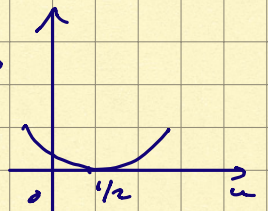
Note: • $2\sqrt{e_k(1-e_k)} \leq 1$

$$e_k(1-e_k) \leq \frac{1}{4}$$

$$e_k - e_k^2 \leq \frac{1}{4}$$

$$e_k^2 - e_k + 1/4 \geq 0$$

$$(e_k - 1/2)^2 \geq 0$$



$$\bullet \sum_{k=1}^K \alpha_k = \sum_{k=1}^K \frac{1}{2} \log \frac{1-e_k}{e_k} = \sum_{k=1}^K \log \sqrt{\frac{1-e_k}{e_k}}$$

$$= \log \prod_{k=1}^K \left(\frac{1-e_k}{e_k}\right)^{1/2}$$

$$\bullet e_k < 1/2 \quad \forall k \quad \Rightarrow \quad \prod_{k=1}^K 2\sqrt{e_k(1-e_k)} \xrightarrow{K \rightarrow \infty} 0$$

② Neural Nets

- linear classifier: $x \mapsto \text{sgn} \langle w, x \rangle$ $w, x \in \mathbb{R}^d$

$$\mathcal{F} := \left\{ x \mapsto \langle w, x \rangle : \|w\| \leq B \right\}$$

$$\|w\| := \left(\sum_{j=1}^d |w_j|^2 \right)^{1/2}$$

w/o restriction on w , $R_n \lesssim \sqrt{d/n}$

$$\begin{aligned}
R_n(\mathcal{F}(X^n)) &= \frac{1}{n} \mathbb{E}_{\mathcal{E}} \left[\sup_{\|w\| \leq B} \left| \sum_{i=1}^n \varepsilon_i \langle w, X_i \rangle \right| \right] \\
&= \frac{1}{n} \mathbb{E}_{\mathcal{E}} \left[\sup_{\|w\| \leq B} \left| \langle w, \sum_{i=1}^n \varepsilon_i X_i \rangle \right| \right] \\
&= \frac{B}{n} \mathbb{E}_{\mathcal{E}} \left\| \sum_{i=1}^n \varepsilon_i X_i \right\|
\end{aligned}$$

[Cauchy-Schwarz: $\|v\| = \sup_{\|w\| \leq 1} |\langle w, v \rangle|$]

$$\begin{aligned}
\mathbb{E}_{\mathcal{E}} \left\| \sum_{i=1}^n \varepsilon_i X_i \right\| &= \mathbb{E}_{\mathcal{E}} \sqrt{\sum_{i=1}^n \sum_{j=1}^n \varepsilon_i \varepsilon_j \langle X_i, X_j \rangle} \\
&\leq \sqrt{\sum_{i=1}^n \sum_{j=1}^n \langle X_i, X_j \rangle \underbrace{\mathbb{E}[\varepsilon_i \varepsilon_j]}_{\mathbb{1}_{\{i=j\}}}} \\
&= \sqrt{\sum_{i=1}^n \|X_i\|^2}
\end{aligned}$$

$$\therefore R_n(\mathcal{F}(X^n)) \leq \frac{B}{n} \sqrt{\sum_{i=1}^n \|X_i\|^2}$$

E.g. if X_1, \dots, X_n are elements of a ball of rad. R centered at 0 , then

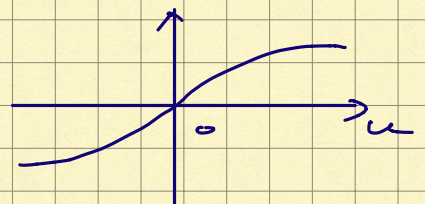
$$R_n(\mathcal{F}(X^n)) \leq \frac{BR}{\sqrt{n}} \leftarrow \text{dimension-free}$$

- add a nonlinearity

$$\sigma: \mathbb{R} \rightarrow \mathbb{R}$$

$$\sigma(0) = 0, \quad L\text{-Lipschitz}$$

e.g. $\sigma(u) = \tanh u = \frac{e^u - e^{-u}}{e^u + e^{-u}}$



$$x \mapsto \text{sgn } \sigma(\langle w, x \rangle)$$

$$\mathcal{F}_\sigma := \left\{ x \mapsto \sigma(\langle w, x \rangle) : \|w\| \leq B \right\}$$

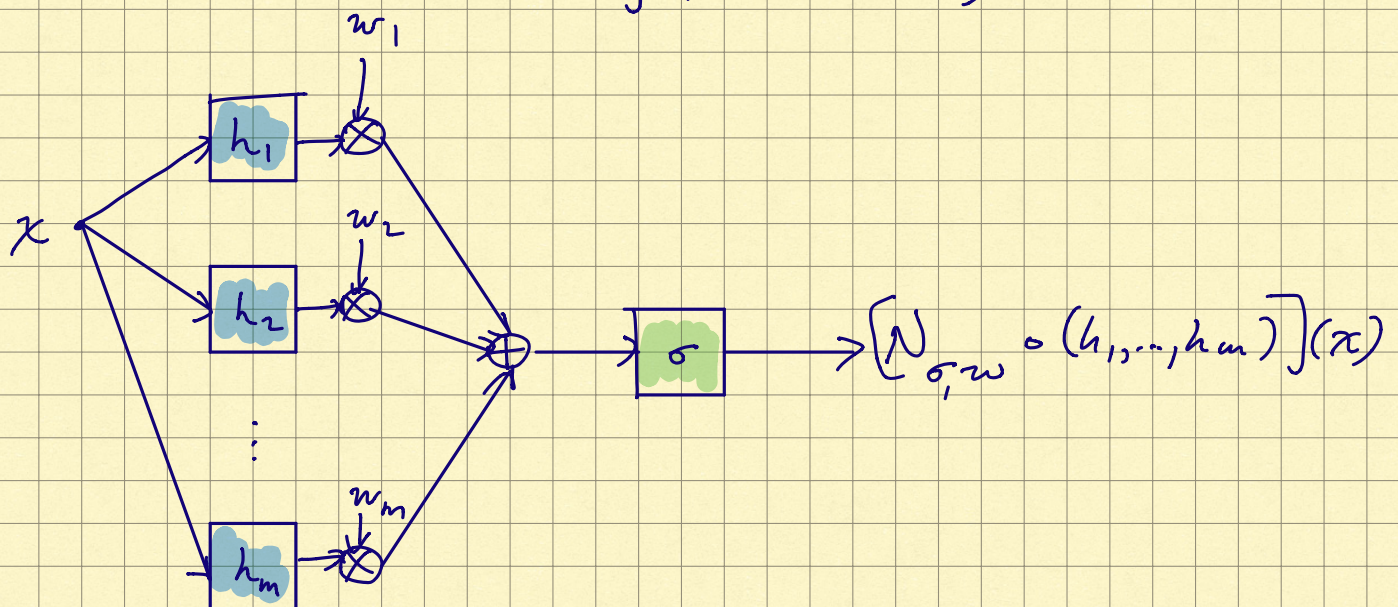
$$\begin{aligned} R_n(\mathcal{F}_\sigma(X^n)) &= R_n(\sigma \circ \mathcal{F}(X^n)) && \mathcal{F}: \text{linear} \\ &\leq 2L R_n(\mathcal{F}(X^n)) && (\text{Contraction principle}) \\ &\leq \frac{2LB}{n} \sqrt{\sum_{i=1}^n \|x_i\|^2} \end{aligned}$$

- neural nets

$$\left. \begin{array}{l} \sigma: \mathbb{R} \rightarrow \mathbb{R} \\ w \in \mathbb{R}^m \end{array} \right\} \begin{array}{l} N_{\sigma, w}: \mathbb{R}^m \rightarrow \mathbb{R} \\ N_{\sigma, w}(u) := \sigma\left(\sum_{j=1}^m w_j u_j\right) \end{array}$$

$$h_1, \dots, h_m: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\begin{aligned} [N_{\sigma, w} \circ (h_1, \dots, h_m)](x) &:= N_{\sigma, w}(h_1(x), \dots, h_m(x)) \\ &= \sigma\left(\sum_{j=1}^m w_j h_j(x)\right) \end{aligned}$$



deep neural nets - define recursively

• \mathcal{G}_j - base classifiers $g: \mathbb{R} \rightarrow \mathbb{R}, x \in \mathbb{R}^d$

• fix $l \in \mathbb{N}$ (# layers)

$\sigma_1, \dots, \sigma_l : \mathbb{R} \rightarrow \mathbb{R}$

$\sigma_l(0) = 0, \sigma_k$ L_k -Lip.

$B_1, \dots, B_l > 0$

• $\mathcal{F}_0 := \mathcal{G}$

• given $\mathcal{F}_0, \dots, \mathcal{F}_{j-1},$ ($j = 1, \dots, l$)

$\mathcal{F}_j := \left\{ N_{\sigma_j, w} \circ (f_1, \dots, f_m) : \begin{array}{l} m \in \mathbb{N} \\ |w_1| + \dots + |w_m| \leq B_j \\ f_1, \dots, f_m \in \mathcal{F}_{j-1} \end{array} \right\}$

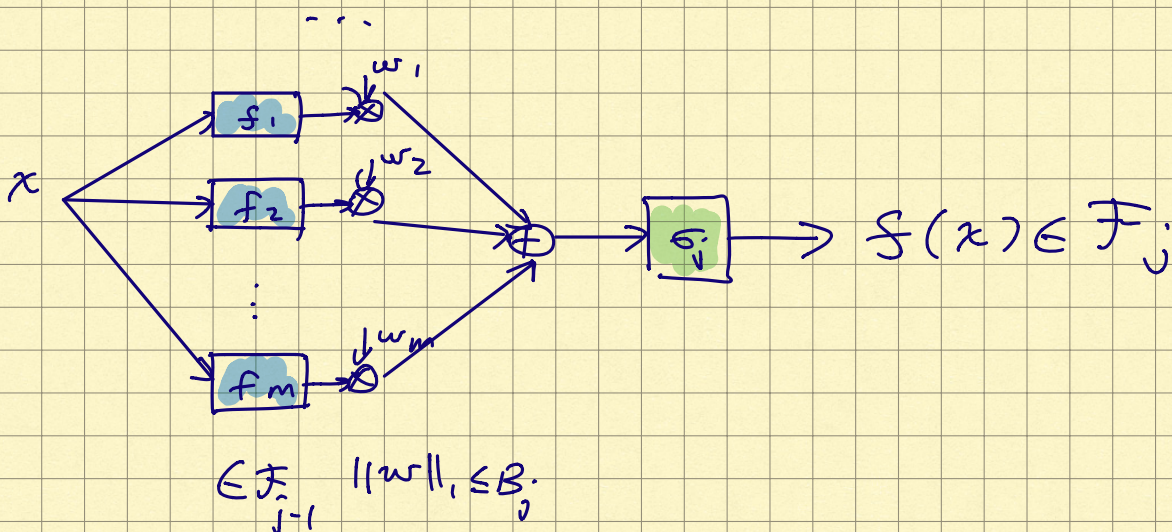
\mathcal{F}_l : class of all l -layer nets w/ activations $\sigma_1, \dots, \sigma_l,$ weight size constraints B_1, \dots, B_l

$R_n(\mathcal{F}_l) \leq ?$

$\mathcal{F}_0 = \mathcal{G}$

$\mathcal{F}_1 : f(x) = \sigma_1 \left(\sum_{k=1}^m w_k g_k(x) \right) \quad \begin{array}{l} m \in \mathbb{N} \\ |w_1| + \dots + |w_m| \leq B_1 \\ g_1, \dots, g_m \in \mathcal{G} \end{array}$

$\mathcal{F}_2 : f(x) = \sigma_2 \left(\sum_{k=1}^m w_k f_k(x) \right), f_1, \dots, f_m \in \mathcal{F}_1$



Naive bound: use contraction principle

$$F_j = \sigma_j \circ (B_j \cdot \text{absconv}(F_{j-1}))$$

$$R_n(F_\ell) \leq 2L_\ell B_\ell \cdot R_n(F_{j-1})$$

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

- say, $L_1, \dots, L_\ell \leq 1$

$$\prod_{j=1}^{\ell} (2B_j) = 2^\ell \prod_{j=1}^{\ell} B_j$$

needs to
be $\leq 2^{-\ell}$
for a bd
that does not
blow up w/l

Next lecture: $2^\ell \rightarrow \sqrt{\ell}$

- (Golowich-Rathlin-Shamir 2017)

- also see Bartlett-Foster-Talgarisky 2017