

CS861 Notes

Young Wu

April 13, 2019

1 Lecture 1

output, label space Y (classification: finite discrete) (regression \mathbb{R})

input, item, instance, object, point space X (e.g. \mathbb{R}^d)

training set: $(x_i \in X, y_i \in Y)_{i=1:n} \stackrel{iid}{\sim} P_{X \times Y}$

test data $\stackrel{iid}{\sim} P_{X \times Y}$

finding "best" $h : X \rightarrow Y$

$h^*(x) \in \arg \max_y P(y|x)$, unknown, cannot compute

loss function: $\ell : Y \times Y \rightarrow \mathbb{R}_{\geq 0}$

e.g. 0 - 1 loss:

$$\ell(y_1, y_2) = \begin{cases} 1 & \text{if } y_1 \neq y_2 \\ 0 & \text{otherwise} \end{cases}$$

squared loss:

$$\ell(y_1, y_2) = \frac{1}{2} (y_1 - y_2)^2$$

Risk: $R(h) := \mathbb{E}_{(x,y) \sim P} [\ell(h(x), y)]$

Empirical Risk: $\hat{R}(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$

$h^* = \arg \min_{\{h: X \rightarrow Y\}} R(h)$, cannot compute

Hypothesis space: $\mathcal{H} = \{h\}$

What ML does: Given $(x_i, y_i)_{i=1:n} := S$ or $D \stackrel{iid}{\sim} P_{XY}^n$

$$\begin{aligned} \hat{h} &\in \arg \min_{h \in \mathcal{H}} \hat{R}(h) \\ \Rightarrow \hat{R}(\hat{h}) &= 0 \end{aligned}$$

Assumption: Realizability

$$\begin{aligned} \exists h^* \in \mathcal{H}, R(h^*) &= 0 \\ \Rightarrow \hat{R}(h^*) &= 0 \end{aligned}$$

Probably Approximately Correct

$$R(\hat{h}) \leq \varepsilon$$

2 Lecture 2

1. Get dataset $S = (x_i, y_i)_{i=1:n}$
2. Run ML $\hat{h}_s = \text{ML}(S), \hat{R}_s(\hat{h}_s) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{h}_s(x_i), y_i)$
3. Test set error $\hat{R}(\hat{h}_s) = \frac{1}{m} \sum_{j=1}^m \ell(\hat{h}_s(x_j), y_j), e.g. 0.01, e.g. 0$

$$T = (x_j, y_j)_{j=1:m}$$

$$R(h_s) = \mathbb{E}_{(x,y) \sim P_{XY}} \left[\ell(\hat{h}_s(x), y) \right], \text{ head prob}$$

Given $\hat{R} = 0$

Suppose $R > \varepsilon$

event prob

Prob (m trials all tails) $< (1 - \varepsilon)^m$

Prob (T has $\hat{R}(\hat{h}_s) = 0$) $\leq (1 - \varepsilon)^m$

Suppose $\hat{R}(\hat{h}_s) = 0$

want: statement about $R(\hat{h}_s)$ being large (bad)

$$\hat{h}_s \in \arg \min_{h \in \mathcal{H}} \hat{R}(h)$$

$\in \leftarrow$ we picked ANY of them!

Empirical Risk Minimization

$R(\hat{h}_s)$: head prob of coin \hat{h}_s

$$\left\{ S : \exists h \in \mathcal{H} : R(h) < \varepsilon \wedge \hat{R}_s(h) = 0 \right\}$$

$$(X \times Y)^n \setminus \left\{ S : R(\hat{h}_s) > \varepsilon \right\}$$

$$\mathbb{P}_{S \sim P^n} \left[\underbrace{\left\{ S : R(\hat{h}_s) > \varepsilon \right\}}_S \right] < \delta$$

fix $h \in \mathcal{H}_\varepsilon, \mathcal{H}_\varepsilon = \{h : R(h) > \varepsilon\} \subset \mathcal{H}, \mathcal{S}_h = \{S : \hat{R}_s(h) = 0\}$

$$\mathcal{S} \subseteq \{\mathcal{S}_h : h \in \mathcal{H}_\varepsilon\} := \bigcup_{h \in \mathcal{H}_\varepsilon} \mathcal{S}_h$$

$$\begin{aligned}
\mathbb{P}(\mathcal{S}) &\leq \mathbb{P}\left(\bigcup_{h \in \mathcal{H}_\varepsilon} \mathcal{S}_h\right) \\
&\stackrel{Union B}{\leq} \sum_{h \in \mathcal{H}_\varepsilon} \mathbb{P}(\mathcal{S}_h) \\
&= \sum_{h \in \mathcal{H}_\varepsilon} \mathbb{P}(n \text{ tails given } \mathbb{P}(\text{head}) > \varepsilon) \\
&\leq \sum_{h \in \mathcal{H}_\varepsilon}^{R(h) > \varepsilon} (1 - \varepsilon)^n \\
&\leq \sum_{h \in \boxed{\mathcal{H}}} (1 - \varepsilon)^n \\
&\stackrel{\boxed{\mathcal{H} \text{ finite}}}{=} |\mathcal{H}| (1 - \varepsilon)^n
\end{aligned}$$

3 Lecture 3

- "task, world, environment, population" fixed unknown P_{XY}
- training data $S \stackrel{iid}{\sim} P^n$
- ERM: $\hat{h}_s = \arg \min_{h \in \mathcal{H}} \hat{R}(h) := \frac{1}{n} \sum_{(x,y) \in S} \ell(h(x), y)$
- Hypo space \mathcal{H}

"ideally" want:

$$h^* \in \arg \min_{h \in \mathcal{H}} R(h) := \mathbb{E}_{(x,y) \sim P} \ell(h(x), y)$$

A1: $\ell 0-1$ loss

A2: $R(h^*) = 0$ "realizable case"

A3: $|\mathcal{H}| < \infty$

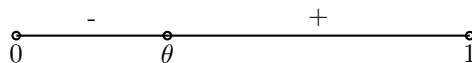
- "Wolf" $\mathcal{H}_\varepsilon := \{h \in \mathcal{H} : R(h) > \varepsilon\}$

A2: $\hat{R}(\hat{h}_s) = 0$ "looks like a sheep"

- "Sheep" $\mathcal{H} \setminus \mathcal{H}_\varepsilon$

Bad event:

- Want: $\mathbb{P}_{S \sim P^n} [\{S : S \text{ enables a wolf to look like the best sheep}\}] < \delta$



$$P_X = U[0, 1]$$

$$P\left(y=+|x\right)=\left\{\begin{array}{ll}1 & \text{if } x\geqslant \theta \\ 0 & \text{otherwise}\end{array}\right.$$

$$y\in\{-,+\}$$

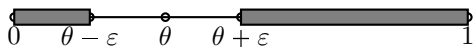
$$h_a=\left\{\begin{array}{ll}+ & \text{if } x\geqslant a \\ - & \text{otherwise}\end{array}\right.$$

$$\mathcal{H}=\{h_a:a\in[0,1]\}$$

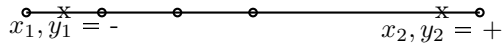
$$\ell:0-1$$

$$R\left(h_a\right)=\left|a-\theta\right|$$

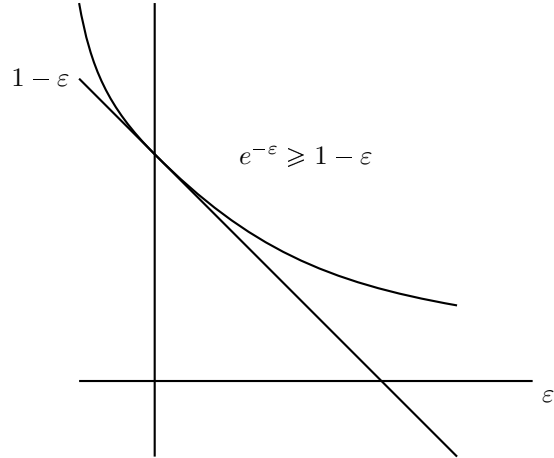
”wolves”



S



$$\begin{aligned} & \mathbb{P}\left[\left\{S:\exists h\in\mathcal{H}_\varepsilon,\hat{R}\left(h\right)=0\right\}\right] \\ &= \mathbb{P}\left[\bigcup_{h\in\mathcal{H}_\varepsilon}\left\{S:\hat{R}\left(h\right)=0\right\}\right] \\ &\stackrel{Unionb.}{\leqslant} \sum_{h\in\mathcal{H}_\varepsilon}\mathbb{P}\left[\left\{S:\hat{R}\left(h\right)=0\wedge R\left(h\right)>\varepsilon\right\}\right] \\ &\stackrel{\mathcal{H}_\varepsilon}{\leqslant} \sum_{h\in\mathcal{H}_\varepsilon}\left(1-\varepsilon\right)^n \\ &\stackrel{\mathcal{H}_\varepsilon\subseteq\mathcal{H}}{\leqslant} \sum_{h\in\mathcal{H}}\left(1-\varepsilon\right)^n \\ &\stackrel{A3}{=} |\mathcal{H}|\left(1-\varepsilon\right)^n \\ &\stackrel{e^{-\varepsilon}\geqslant 1-\varepsilon}{\leqslant} |\mathcal{H}|e^{-\varepsilon n}:=\delta \end{aligned}$$



$$-\epsilon n = \log \left(\frac{\delta}{|\mathcal{H}|} \right)$$

$$\epsilon = \frac{1}{n} \log \left(\frac{|\mathcal{H}|}{\delta} \right)$$

$$\mathbb{P}(\{\text{Bad } S\}) \leq \delta$$

$$\mathbb{P}((X \times Y)^n \setminus \{\text{Bad } S\}) > 1 - \delta$$

$$\mathbb{P}\left[R(\hat{h}_s) \leq \epsilon\right] \geq 1 - \delta$$

With prob at least $1 - \delta$,

$$R(\hat{h}_s) \leq \epsilon := \frac{\log(|\mathcal{H}|) - \log(\delta)}{n}$$

$$n := \frac{\log(|\mathcal{H}|) - \log(\delta)}{\epsilon}$$

4 Lecture 4

A1: ℓ is 0 – 1 loss

$$\mathbb{P}_{S \sim P^n} \left(\left\{ S : \exists h \in \mathcal{H}_\epsilon, \underbrace{\hat{R}_s(h) = 0}_{\text{A2: Realizable } \min_{h \in \mathcal{H}} R(h) = 0} \right\} \right)$$

$$\mathcal{H}_\epsilon = \{h \in \mathcal{H} : R(h) > \epsilon\}$$

$$R(h) = \mathbb{E}_{(x,y) \sim P} \ell(h(x), y)$$

$$\hat{R}_s(h) = \frac{1}{n} \sum_{x,y \in S} \ell(h(x), y)$$

$$\begin{aligned}
& \mathbb{P} \left(\bigcup_{h \in \mathcal{H}_\varepsilon} \left\{ S : \hat{R}_s(h) = 0 \right\} \right) \\
& \stackrel{\text{Union}}{\leq} \sum_{h \in \mathcal{H}_\varepsilon} \mathbb{P} \left(S : \hat{R}_s(h) = 0 \wedge R(h) > \varepsilon \right) \\
& \stackrel{\text{Realizability}}{\leq} \sum_{h \in \mathcal{H}_\varepsilon} (1 - \varepsilon)^n \\
& \stackrel{A3: |\mathcal{H}| < \infty}{\leq} |\mathcal{H}| (1 - \varepsilon)^n \\
& \leq |\mathcal{H}| e^{-\varepsilon n} := \delta
\end{aligned}$$

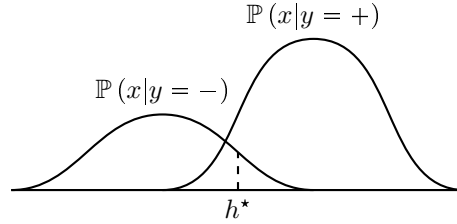
$$\begin{aligned}
n & \geq \frac{\log |\mathcal{H}| - \log \delta}{\varepsilon} \\
\varepsilon & \leq O \left(\frac{1}{n} \right)
\end{aligned}$$

Agnostic learning (wrt \mathcal{H}), $R(h^*) \geq 0$

$$h^* \in \arg \min_{h \in \mathcal{H}} R(h)$$

e.g.

$$\mathbb{P}(y = +) = \frac{1}{2}$$



Want Uniform Convergence:

$$\begin{aligned}
& \left\{ S : \exists h \in \mathcal{H}, \left| R(h) - \hat{R}_s(h) \right| > \varepsilon \right\} \\
& \forall h \in \mathcal{H}, \left| R(h) - \hat{R}_s(h) \right| \leq \varepsilon
\end{aligned}$$

" $R(\hat{h}_s) - R(h^*)$ small "

$$\begin{aligned}
& \mathbb{P} \left(\left\{ S : \exists h \in \mathcal{H}, \left| R(h) - \hat{R}_s(h) \right| > \varepsilon \right\} \right) \\
& \leq \sum_{h \in \mathcal{H}} \mathbb{P} \left(\left\{ S : \left| R(h) - \hat{R}_s(h) \right| > \varepsilon \right\} \right)
\end{aligned}$$

$$\stackrel{\text{Hoeffding's}}{\leq} 2|\mathcal{H}|e^{-\frac{2n\varepsilon^2}{(b-a)^2}} := \delta$$

new A1: $\ell \in [a, b]$ or $(x, y) \sim P \ell(h(x), y)$ is subGaussian

remove A2

retain A3

$$\text{Hoeffding's Ineq. } \left| \mu - \frac{1}{n} \sum_i^n \theta_i \right|$$

$$\mathbb{P}_{S \sim P^n} \left(\left\{ S : \left| R(h) - \hat{R}_s(h) \right| > \varepsilon \right\} \right) \leq 2e^{-\frac{2n\varepsilon^2}{(b-a)^2}}$$

$$\text{wp} \geq 1 - \delta, \forall h \in \mathcal{H}, \left| R(h) - \hat{R}_s(h) \right| \leq \varepsilon$$

$$\begin{aligned} & R(\hat{h}_s) - R(h^*) \\ &= R(\hat{h}_s) - \hat{R}_s(\hat{h}_s) + \hat{R}_s(\hat{h}_s) - \hat{R}_s(h^*) + \hat{R}_s(h^*) - R(h^*) \\ &\leq \varepsilon + \underbrace{\left(\hat{R}_s(\hat{h}_s) - \hat{R}_s(h^*) \right)}_{\text{ERM} \leq 0} + \varepsilon \\ &\leq 2\varepsilon \end{aligned}$$

$$\begin{aligned} \frac{2n\varepsilon^2}{(b-a)^2} &= \log \frac{2|\mathcal{H}|}{\delta} \\ n &= \frac{\log(2|\mathcal{H}|) - \log \delta}{2\varepsilon^2} (b-a)^2 \\ \varepsilon &\leq O\left(\frac{1}{\sqrt{n}}\right) \end{aligned}$$

"learning alg" $A(S) = h$

5 Lecture 5

$$\mathbb{P}(\{S \text{ bad}\}) \leq \delta = |\mathcal{H}|e^{-n\varepsilon}$$

i.e. $\exists h, \left| R(h) - \hat{R}_s(h) \right| > \varepsilon$

Today's goal: $\mathbb{P}(\{S \text{ bad}\}) \geq \delta$

Any learning Algorithm $A : \{S\} \rightarrow \mathcal{H}$, ERM is a A

$$S \stackrel{iid}{\sim} P^n(x, y), n \leq \frac{|X|}{2}$$

Theorem 1. $\forall A, \exists P$

1. $\exists h : X \rightarrow Y, R_P(h) = 0$
2. $\mathbb{P} \left(\left\{ S : R_P(A(S)) \geq \frac{1}{8} \right\} \right) \geq \frac{1}{7}$

$$2n = |\{ \text{distinct } x' \text{'s in } S \} \cup \{ \text{some other distinct } x \text{ from } X \text{ not in } S \}|$$

$$\overbrace{x_1, \dots, x_n}^{S_n}, x_{n+1}, \dots, x_{2n}$$

Construct a family of $P(x, y) = P(x) \cdot P(y|x)$

1. $P(x)$ uniform on the $2n$ items
2. $\mathcal{C} := 2^{2n}$ labelings over $2n$ items

$$\mathcal{C} \text{ rows} \Rightarrow \mathcal{C} \text{ joint } P_{XY} = \begin{cases} 0 \dots 0 \dots 0 \dots 00, \mathbb{P}(y=1|x) = 0 \forall x \\ 0 \dots 0 \dots 0 \dots 01, \mathbb{P}(y=1|x_{2n}) = 1, \mathbb{P}(y=1|x \neq x_{2n}) = 0 \\ 1 \dots 1 \dots 1 \dots 11, \mathbb{P}(y=1|x) = 1 \forall x \end{cases}$$

Key idea: $\max_{c \in [\mathcal{C}]} \mathbb{E}_{S \sim P_c^n} R_{P_c}(A(S))$

$$\boxed{z \text{ rv } \geq 0, \mathbb{E}(z) \geq a \Rightarrow \mathbb{P}(z \geq b) \geq c}$$

$$\boxed{z \in [0, 1] \text{ rv }, \mathbb{E}(z) \geq u \Rightarrow \mathbb{P}(z \geq 1-a) \geq \frac{u - (1-a)}{a}}$$

Lemma B.1 Markov's ineq

$Q = (2n)^n$ distinct S sequences, S_1, S_2, \dots, S_Q

6 Lecture 6

Take $2n$ distinct item from X

$x_1 \dots x_{2n}$
 $0 \dots 00$
 $0 \dots 01$
 \dots

$\mathcal{C} := 2^{2n}$ different labelings of the $2n$ items

$$P_1(y|x) \text{ "world" } P_1(x, y) = P(x) P_1(y|x)$$

$$\begin{aligned}
& \dots \\
& P_C(y|x) \\
P(x) &= \frac{1}{2n} \\
R_C(h) &= \mathbb{E}_{(x,y) \sim P_C(x,y)} \underbrace{\ell}_{0-1 \text{ loss}}(h(x), y), c \in [C]
\end{aligned}$$

$$n \leq \frac{|X|}{2}$$

No-Free lunch Thm

$$\forall A, \exists P_c (c \in [C])$$

$$1. \exists h, R_c(h) = 0$$

$$2. \mathbb{P}_{S \sim P_c^n} \left[R_c(A(S)) \geq \frac{1}{8} \right] \geq \frac{1}{7}$$

$$\max_{c \sim [C]} \mathbb{E}_{S \sim P^n} R_c(A(S)) \geq \frac{1}{4} \xrightarrow{\text{Markov}} (2)$$

$$S = (x^1, x^2, \dots, x^n), x^{i \leftarrow \text{position in } S}$$

$$S_1 : x_1, x_1, \dots, x_1$$

$$S_2 : x_1, \dots, x_1, x_2$$

...

$$S_Q : x_{2n}, \dots, x_{2n}$$

$$Q := (2n)^n$$

$$\max_{c \in [C]} \mathbb{E}_{S \sim P_c^n} R_c(A(S))$$

$$= \max_{c \in [C]} \frac{1}{Q} \sum_{q=1}^Q R_c(A(S_q))$$

$$\geq \frac{1}{C} \sum_{c=1}^C \frac{1}{Q} \sum_{q=1}^Q R_c(A(S_q))$$

$$= \frac{1}{Q} \sum_{q=1}^Q \frac{1}{C} \sum_{c=1}^C R_c(A(S_q))$$

$$\stackrel{P_x \text{ Unif}}{=} \frac{1}{Q} \sum_{q=1}^Q \frac{1}{C} \sum_{c=1}^C \frac{1}{2n} \sum_{i=1}^{2n} \mathbb{1}_{[A(S_q)(x_i) \neq y_c(x_i)]}$$

$$t := |\{x_1, \dots, x_{2n}\} \setminus \{S_q\}| \geq n$$

Let $\{v_1, \dots, v_t\}$ be the set

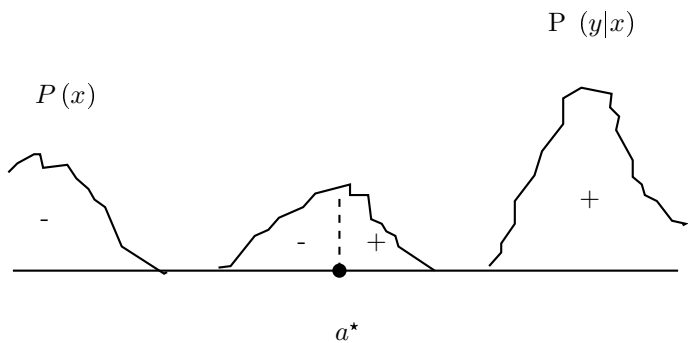
$$\begin{aligned}
 & \text{only consider } \{v\} \geq \frac{1}{Q} \sum_{q=1}^Q \frac{1}{c} \sum_{c=1}^c \frac{1}{2n} \sum_{i=1}^t \mathbb{1}_{[A(S_q)(v_i) \neq y_c(v_i)]} \\
 &= \frac{1}{Q} \sum_{q=1}^Q \frac{1}{2n} \sum_{i=1}^t \underbrace{\frac{1}{c} \sum_{c=1}^c \mathbb{1}_{[A(S_q)(v_i) \neq y_c(v_i)]}}_{\frac{1}{2}} \\
 &= \frac{1}{Q} \sum_{q=1}^Q \frac{1}{2n} \sum_{i=1}^t \frac{1}{2} \\
 &= \frac{1}{Q} \sum_{q=1}^Q \frac{1}{2n} \frac{1}{2} t \\
 &\geq \frac{1}{2} \sum_{q=1}^Q \frac{1}{2n} \frac{1}{2} n \\
 &= \frac{1}{4}
 \end{aligned}$$

$$h_c(x) = \begin{cases} 0 & \text{if } x \notin \mathbb{Z} \\ y_c(x) & \text{otherwise} \end{cases}$$

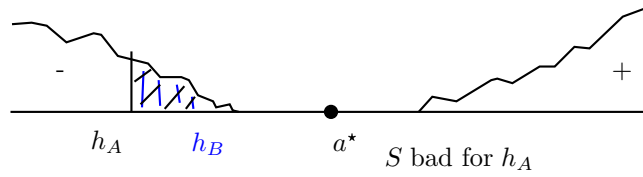
$$|\mathcal{H}| = \infty$$

ex 6.1

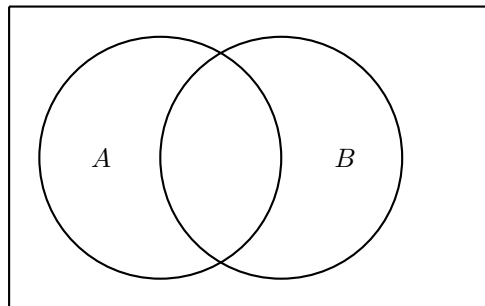
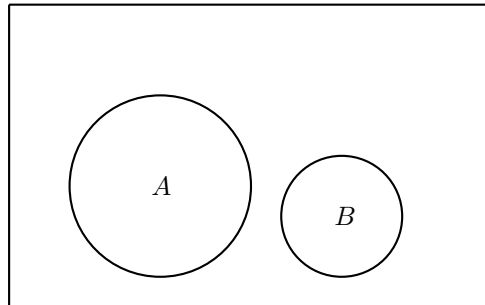
$$\mathcal{H} = \{h_A : a \in \mathbb{R}\}$$



$$h_a(x) = \begin{cases} 1 & \text{if } x \geq a \\ 0 & \text{otherwise} \end{cases}$$



$$\mathbb{P} \left(\bigcup_{h \in \mathcal{H}_\varepsilon} \{S \text{ makes } h, 0 \text{ training risk} \} \right) \leq \sum_{h \in \mathcal{H}_\varepsilon} \mathbb{P}(\{S_{\text{bad}}\})$$



$$P(A \text{ or } B) \leq P(A) + P(B)$$

7 Lecture 7

Recall: finite \mathcal{H}

$$\mathbb{P}_{S \sim P^n} \left(\max_{h \in \mathcal{H}} R(h) - \hat{R}_s(h) \leq \varepsilon \right) \geq 1 - \delta$$

$$\varepsilon = \sqrt{\frac{\log |\mathcal{H}| + \log \frac{1}{\delta}}{2n}}$$

Today: Any $\mathcal{H} = \{h\}, h : X \rightarrow Y$

$$\text{ex. } \mathcal{H} = \left\{ \underbrace{h_a}_{a \in \mathbb{R}}(x) = \text{sign}[\sin(ax)] \right\}$$

$$VC = \infty$$

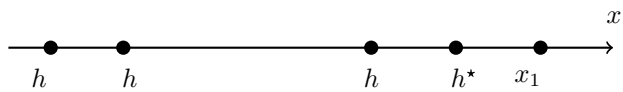
Growth number

$$G(n) := \sup_{x_1, \dots, x_n \in X} \left| \left\{ \mathbb{1}_{[h(x_1) \neq y_1]}, \dots, \mathbb{1}_{[h(x_n) \neq y_n]} : h \in \mathcal{H} \right\} \right|$$

$$(y_i = h^*(x_i))$$

$$d := VC(\mathcal{H}) = \arg \max_n G_{\mathcal{H}}(n) = 2^n$$

$$\text{ex: } \mathcal{H} = \left\{ \underbrace{h_a}_{X=\mathbb{R}}(x) = \{1, x \geq a, 0 \text{ ow}\}, a \in \mathbb{R} \right\}$$

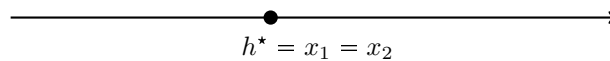
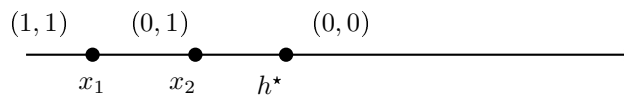
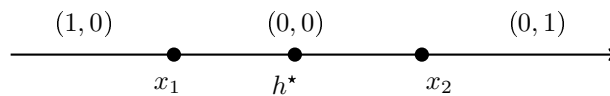


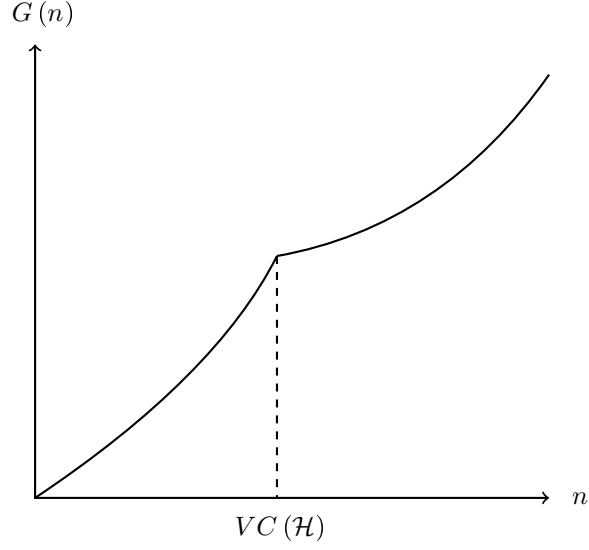
$$\begin{aligned} &(\mathbb{1}_{[h(x_1) \neq 0]}) \\ &|\{(0), (1)\}| = 2 \end{aligned}$$

$$G(1) = 2, VC(\mathcal{H}) = 1$$

$$G(2) = 3$$

$$G(3)$$





Proof outline

- Introduce "ghost sample" $S' \sim P^n$

$$\hat{R}'_{s'}(h) = \frac{1}{n} \sum_{x', y' \in S'} \ell(h(x'), y')$$

- Symmetrization Lemma

$$\boxed{\forall \varepsilon \geq \sqrt{\frac{2 \log 2}{n}}, \mathbb{P}_{s \sim P^n} \left(\sup_{h \in \mathcal{H}} R(h) - \hat{R}_s(h) > \varepsilon \right) \stackrel{\text{Sym lemma}}{\leq} 2 \mathbb{P}_s \left(\sup_{h \in \mathcal{H}} \hat{R}'_{s'}(h) - \hat{R}_s(h) > \frac{\varepsilon}{2} \right)}$$

define

$$\boxed{\left\{ \left(\underbrace{\ell(h(x_1), y_1), \dots, \ell(h(x_n), y_n)}_S, \underbrace{\ell(h(x'_1), y'_1), \dots, \ell(h(x'_n), y'_n)}_{S'} \right) : h \in \mathcal{H} \right\} := \text{Vec}(2n)}$$

$$\stackrel{\text{def Vec}(2n)}{=} 2 \mathbb{P}_s \left(\max_{h \in \text{Vec}(2n)} \hat{R}'_{s'}(h) - \hat{R}_s(h) > \frac{\varepsilon}{2} \right)$$

$$\stackrel{\text{Union b}}{\leq} 2 |\text{Vec}(2n)| \mathbb{P} \left(\hat{R}'_{s'}(h) - \hat{R}_s(h) > \frac{\varepsilon}{2} \right)$$

$$\stackrel{\text{Growth num}}{\leq} 2 G_{\mathcal{H}}(2n) \mathbb{P} \left(\hat{R}'_{s'}(h) - \hat{R}_s(h) > \frac{\varepsilon}{2} \right)$$

$$\stackrel{\text{Hoeffding's (2 samples)}}{\leq} 2 G_{\mathcal{H}}(2n) e^{-\frac{n \left(\frac{\varepsilon}{2} \right)^2}{2}}$$

$$\stackrel{\varepsilon \geq \sqrt{\frac{2 \log 2}{n}}}{\Rightarrow} \text{wp} \geq 1 - \delta, \sup_{h \in \mathcal{H}} R(h) - \hat{R}_s(h) \leq 2 \sqrt{2 \frac{\log G_{\mathcal{H}}(2n) + \log \frac{2}{\delta}}{n}}$$

"Tighter than VC-bound"

8 Lecture 8

Growth number $G_{\mathcal{H}}(n) := \sup_{x_1 \dots x_n \in X} |\{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\}|$

symmetrization lemma

$$\underbrace{\forall \varepsilon \geq \sqrt{\frac{2 \log 2}{n}}}_{\text{Assumption A1}},$$

$$P\left(\sup_{h \in \mathcal{H}} R(h) - \hat{R}(h) > \varepsilon\right) \leq 2P\left(\sup_{h \in \mathcal{H}} \underbrace{\hat{R}'(h)}_{\text{ghost sample}} - \hat{R}(h) > \frac{\varepsilon}{2}\right)$$

$$\Rightarrow P\left(\sup_{h \in \mathcal{H}} R(h) - \hat{R}(h) > \varepsilon\right) \leq 2G_{\mathcal{H}}(2n) e^{-\frac{n\left(\frac{\varepsilon}{2}\right)^2}{2}} := \delta$$

$$\varepsilon := \sqrt{\frac{8(\log 2G_{\mathcal{H}}(2n) - \log \delta)}{n}}$$

VC-dim of $\mathcal{H} := d$

$$d = \max_n n : G_{\mathcal{H}}(n) = 2^n$$

Sauer's Lemma:

Assuming $d < \infty$. Then

$$G_n(n) \leq \sum_{i=0}^d \binom{n}{i}$$

$$\Rightarrow \text{If } n \geq d, G_{\mathcal{H}}(n) \leq \left(\frac{en}{d}\right)^d$$

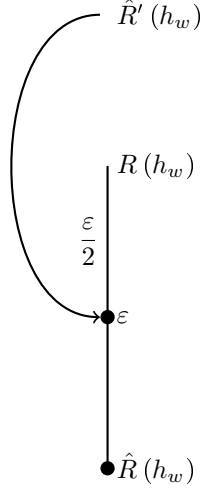
$$P\left(\sup_{h \in \mathcal{H}} R(h) - \hat{R}(h) > \varepsilon\right) \leq 2\left(\frac{2en}{d}\right)^d e^{-\frac{n\varepsilon^2}{8}} := \delta$$

$$\Rightarrow \varepsilon = \sqrt{\frac{d \log n + d \log \frac{2e}{d} + \log \frac{2}{\delta}}{n}} \Rightarrow O\left(\sqrt{\frac{d}{n}}\right)$$

Proof of Sym Lemma. Let a "worst" hypo be $h_w \in \arg \sup_{h \in \mathcal{H}} R(h) - \hat{R}(h)$

$$\mathbb{1}_{[R(h_w) - \hat{R}(h_w) > \varepsilon]} \cdot \mathbb{1}_{\left[R(h_w) - \underbrace{\hat{R}'(h_w)}_{\text{ghost}} < \frac{\varepsilon}{2}\right]}$$

$$= \mathbb{1}_{\left[R(h_w) - \hat{R}(h_w) > \varepsilon \wedge (\hat{R}'(h_w) - R(h_w)) > -\frac{\varepsilon}{2}\right]}$$



$$\begin{aligned}
& \stackrel{\text{implication}}{\leq} \mathbb{1}_{\left[\hat{R}'(h_w) - \hat{R}(h_w) > \frac{\varepsilon}{2}\right]} \\
& \stackrel{\text{expectation over ghost sample } S'}{\Rightarrow} \mathbb{1}_{\left[R(h_w) - \hat{R}(h_w) > \varepsilon\right]} \underbrace{P'}_{\text{wrt } S' \sim P_{XY}^n} \left(R(h_w) - \hat{R}'(h_w) < \frac{\varepsilon}{2}\right) \\
& \leq P' \left(\hat{R}'(h_w) - \hat{R}(h_w) > \frac{\varepsilon}{2}\right) \rightarrow \boxed{1}
\end{aligned}$$

By Hoeffding's Ineq

$$\begin{aligned}
& P' \left(R(h_w) - \hat{R}'(h_w) < \frac{\varepsilon}{2}\right) \geq 1 - e^{-2n \left(\frac{\varepsilon}{2}\right)^2} \stackrel{A1}{\geq} \frac{1}{2} \\
& \boxed{1} \Rightarrow \mathbb{1}_{\left[R(h_w) - \hat{R}(h_w) > \varepsilon\right]} \leq 2P' \left(\hat{R}'(h_w) - \hat{R}(h_w) > \frac{\varepsilon}{2}\right) \\
& \stackrel{h_w \text{ def}}{\Rightarrow} \mathbb{1}_{\left[\sup_{h \in \mathcal{H}} R(h) - \hat{R}(h) > \varepsilon\right]} \leq \text{RHS above} \stackrel{\text{implication}}{\leq} 2P' \left(\sup_{h \in \mathcal{H}} \hat{R}'(h) - \hat{R}(h) > \frac{\varepsilon}{2}\right) \\
& \stackrel{\text{sym } S \text{ and } S'}{=} 2P \left(\sup_{h \in \mathcal{H}} \hat{R}'(h) - \hat{R}(h) > \frac{\varepsilon}{2}\right) \\
& \stackrel{\text{expectation over } S}{\Rightarrow} P \left(\sup_{h \in \mathcal{H}} R(h) - \hat{R}(h) > \varepsilon\right) \leq 2P \left(\sup_{h \in \mathcal{H}} \hat{R}'(h) - \hat{R}(h) > \frac{\varepsilon}{2}\right)
\end{aligned}$$

9 Lecture 9

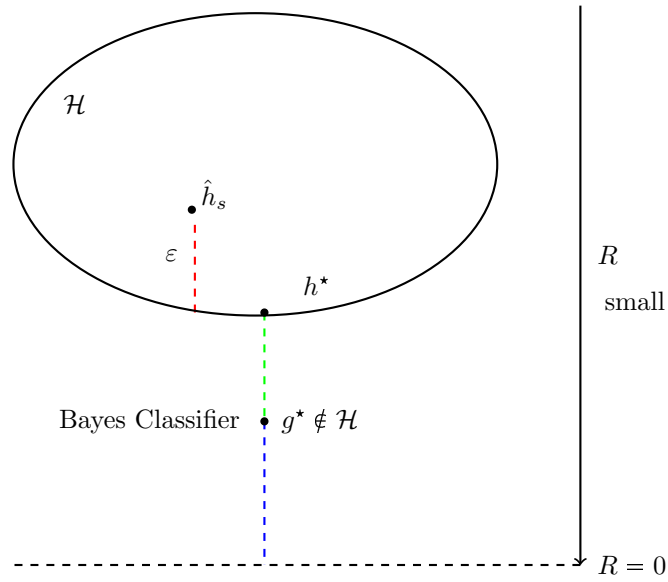
$$\mathbb{P}_{S \sim P^n} \left(\sup_{h \in \mathcal{H}} R(h) - \hat{R}_s(h) \leq \varepsilon\right) \geq 1 - \delta$$

where $\varepsilon \sim O\left(\sqrt{\frac{VC(\mathcal{H}) + \log \frac{1}{\delta}}{n}}\right), \boxed{\varepsilon(n, \delta)}$

$$\mathbb{P}_S \left(R \left(\hat{h}_s^{\text{ERM}} \right) - R(h^*) \leq \varepsilon \right) \geq 1 - \delta$$

where $h^* \in \arg \inf_{h \in \mathcal{H}} R(h)$

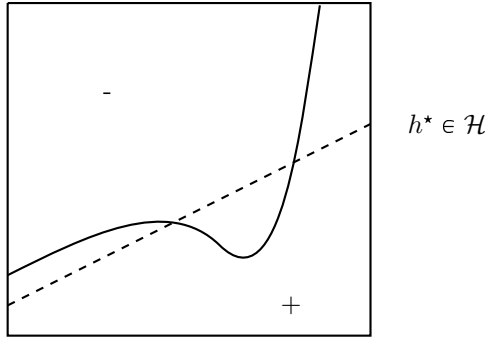
$$R(h^s) = \mathbb{E}_{(x,y) \sim P} \ell(\hat{h}_s(x), y)$$



- (red) Estimation error: $R(\hat{h}_s) - R(h^*)$
- (green) Approximation error: $R(h^*) - R(g^{\text{Bayes}})$
- (blue) Bayes error: $R(g^{\text{Bayes}})$

$$P_X = \text{Unif } [0, 1]^2$$

$$P_{XY}$$

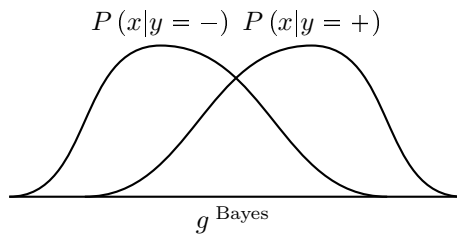


\mathcal{H} linear

$$P_{Y|X}$$

$$1 > P(y = + | x \in + \text{ region}) > \frac{1}{2}$$

$$\frac{1}{2} > P(y = + | x \in - \text{ region}) > 0$$



$$\mathcal{H} = \{ \text{sign} [\sin (\alpha x)] : \alpha \in [1, 2] \}$$

$$g^{\text{Bayes}} \in \arg \max_{y \in Y} P(y|x)$$

Consider $\mathcal{H}_1, \mathcal{H}_2, \dots$

$$VC\left(\mathcal{H}_i\right)<\infty,\,\forall\,i\in N$$

$$\mathbb{P}\left(\sup_{h\in\mathcal{H}}R\left(h\right)-\hat{R}_s\left(h\right)>\varepsilon\left(n,\delta\right)\right)\leqslant\delta$$

$$\Rightarrow\mathbb{P}\left(\sup_{h\in\mathcal{H}_i}R\left(h\right)-\hat{R}_s\left(h\right)>\varepsilon\left(n,w_i\delta\right)\right)\leqslant w_i\delta,\,\forall\,i$$

$$\mathbb{P}\left(\,\forall\,i,\sup_{h\in\mathcal{H}_i}R\left(h\right)-\hat{R}_s\left(h\right)>\varepsilon\left(n,w_i\delta\right)\right)\leqslant\left(\sum_{i=1}^{\infty}w_i\right)\delta,\,\forall\,i$$

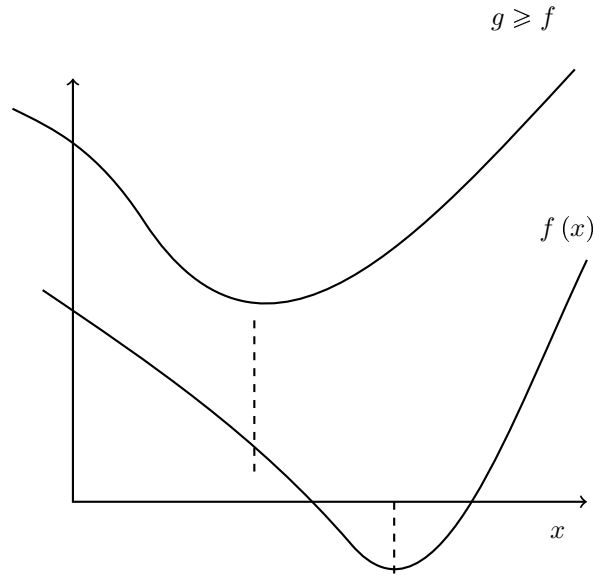
where $\sum_{i=1}^{\infty}w_i\leqslant 1, w_i>0$

$$R(h) - \hat{R}_s(h) \leq \varepsilon_i(n, w_i \delta)$$

$$R(h) \leq \hat{R}_s(h) + \varepsilon_i(n, w_i \delta)$$

"Alg 1"

$$\hat{h}_s \in \arg \inf_{h \in \bigcup_{i=1}^{\infty} \mathcal{H}_i} \hat{R}_s(h) + \min_{i: h \in \mathcal{H}_i} \varepsilon_i(n, w_i \delta)$$



"Alg 2" \leftarrow Stuctural Risk Minimization

$$i^*(h) = \left[\arg \min_{i=1,2,\dots} h \in \mathcal{H}_i \right], \hat{h}_s \in \arg \inf_h \hat{R}_s(h) + \varepsilon_{i^*(h)}(n, w_{i^*(h)} \delta)$$

$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i y_i) + \frac{\lambda}{2} \|\theta\|^2$$

$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i y_i) + \frac{\lambda}{2} \|\theta - 86\|^2$$

10 Lecture 10

Structural Risk Minimization

Input: $\mathcal{H}_1, \mathcal{H}_2, \dots$ each with finite VC

$$w_1, w_2, \dots, \text{ such that } \sum_{i=1}^{\infty} w_i \leq 1, w_i \geq 0$$

$$\varepsilon_i(N, w_i \delta) \approx \sqrt{\frac{VC(\mathcal{H}_i) + \log \frac{1}{w_i \delta}}{N}}$$

$$\mathbb{P}_s \left(\sup_{h \in \mathcal{H}_i} R(h) - \hat{R}(h) > \varepsilon_i(N, w_i \delta) \right) < w_i \delta$$

Union bound over $i = 1, 2, \dots$

$$\mathbb{P}_s \left(\exists i, \sup_{h \in \mathcal{H}_i} R(h) - \hat{R}(h) > \varepsilon_i(N, w_i \delta) \right) < \left(\sum_{i=1}^{\infty} w_i \right) \delta$$

$$\mathbb{P}_s \left(\forall i, \sup_{h \in \mathcal{H}_i} R(h) - \hat{R}(h) \leq \varepsilon_i(N, w_i \delta) \right) \geq 1 - \left(\sum_{i=1}^{\infty} w_i \right) \delta \begin{matrix} \left(\sum_i w_i \delta \leq 1 \right) \\ \geq \end{matrix} 1 - \delta$$

For each \mathcal{H}_i

$$\forall \varepsilon, \delta, P, \exists N, \forall n > N$$

$$\mathbb{P}_{s \sim P^n} \left(\sup_{h \in \mathcal{H}_i} R(h) - \hat{R}_s(h) > \varepsilon \right) < \delta$$

$$N \approx \frac{VC(\mathcal{H}_i) + \log \frac{1}{\delta}}{\varepsilon^2}$$

$$\varepsilon_i(w, \delta) \approx \sqrt{\frac{VC(\mathcal{H}_i) + \log \frac{1}{\delta}}{N}}$$

SRM:

$$\hat{h}^{\text{SRM}} \in \arg \min_{i, h \in \mathcal{H}_i} \hat{R}(h) + \varepsilon_i(N, w_i \delta)$$

Special Case of SRM

Assumption: \mathcal{H} countable

$$\mathcal{H} = \{h_1, h_2, \dots\}$$

$$= \{h_1\} \cup \{h_2\} \cup \dots$$

$$\mathcal{H}_1 = \{h_1\}, \mathcal{H}_2 = \{h_2\} \dots$$

Need w_i for \mathcal{H}_i , equiv denote as $w_h, h \in \mathcal{H}$

$$\sum_{h \in \mathcal{H}} w_h \leq 1, w_h \geq 0 \forall h \in \mathcal{H}$$

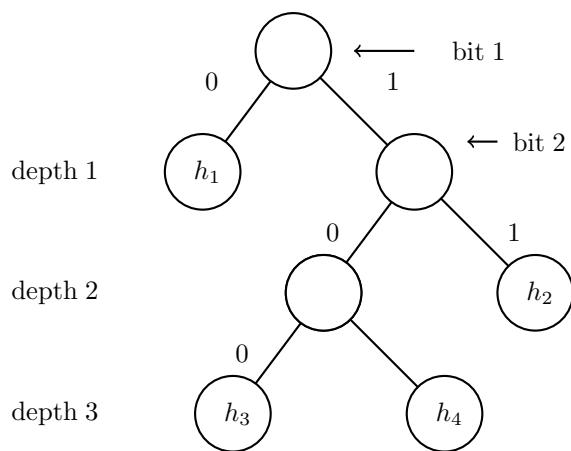
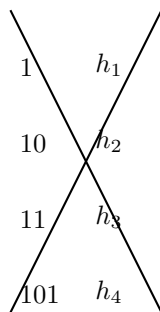
SRM

$$\hat{h}^{\text{SRM}} \in \arg \min_{h \in \mathcal{H}} \hat{R}(h) + \varepsilon_i(N, w_i \delta) = \arg \min_{h \in \mathcal{H}} \hat{R}(h) + \sqrt{\frac{\log \frac{1}{w_i} + \log \frac{2}{\delta}}{2N}}$$

\mathcal{H}_i singleton \Rightarrow Hoeffding's

$$\varepsilon_i = \sqrt{\frac{\log \frac{2}{w_i \delta}}{2N}}$$

Special ² case prefix binary code



$$\begin{array}{lll}
0 & h_1 & w_1 = \frac{1}{2} \\
11 & h_2 & w_2 = \frac{1}{4} \\
100 & h_3 & w_3 = \frac{1}{8} \\
101 & h_4 & w_4 = \frac{1}{8}
\end{array} \Rightarrow$$

$$w_h := \frac{1}{2^{\text{depth}(h)}}$$

Kreft's Thm

$$\sum_{h \in \mathcal{H}} w_h \leq 1$$

$$\arg \min \hat{R}(h) + \sqrt{\frac{\text{depth}(h) \log 2 + \log \frac{1}{\delta}}{2N}}$$

minimum description length (MDL)

\Rightarrow Occam's Razor

PAC-Bayes bound

Set prior P over \mathcal{H}

For any $Q = A(S)$ over \mathcal{H} that learner produces (i.e. equiv of \hat{h}^{ERM})

$$R(Q) := \mathbb{E}_{(x,y) \sim P_{\text{dist}}} \mathbb{E}_{h \sim Q} \ell(h(x), y)$$

11 Lecture 11

PAC-Bayes bound

Prior P over \mathcal{H} , loss $\ell \in [0, 1]$

$$\text{Risk } R(h) = \mathbb{E}_{P_{XY}} \ell(h(x), y), \hat{R}_s(h) = \frac{1}{|S|} \sum_{i=1}^n \ell(h(x_i), y_i)$$

"Gibbs classifier" $h \sim P$

$$R(P) = \mathbb{E}_{h \sim P} R(h), \hat{R}_s(P) = \mathbb{E}_{h \sim P} \hat{R}_s(h)$$

Sample $S \sim P_{XY}^n$, $A(S) := Q$ distribution over \mathcal{H}

Theorem 2. P given, $\forall P_{XY}, \ell \in [0, 1], \varepsilon, \delta,$

$$\mathbb{P}_{S \sim P_{XY}^n} \left(\forall Q, R(Q) \leq \hat{R}_s(Q) + \sqrt{\frac{KL(Q\|P) + \log \frac{n}{\delta}}{2(n-1)}} \right) \geq 1 - \delta$$

$$KL(Q\|P) := \sum_{h \in \mathcal{H}} Q(h) \log \frac{Q(h)}{P(h)}$$

Let Q_1 concentrated on \hat{h}_s^{ERM} , Q_1 minimizes $\hat{R}_s(Q)$

$$KL(Q_1\|P) = Q(\hat{h}_s^{ERM}) \log \frac{Q(\hat{h}_s^{ERM})}{P(\hat{h}_s^{ERM})} = \log \frac{1}{P(\hat{h}_s^{ERM})}$$

$$\lim_{x \rightarrow 0} x \log x = \lim_{x \rightarrow 0} \frac{(\log x)'}{\left(\frac{1}{x}\right)'} = \frac{\frac{1}{x}}{-\frac{1}{x^2}} = -x = 0$$

(read textbook for proof)

11.1 Rademacher Complexity

$$\mathcal{F} := \ell \circ \mathcal{H} = \{f(\cdot, \cdot) := \ell(h(\cdot), \cdot), h \in \mathcal{H}\}$$

$$R(f) = \mathbb{E}_{(x,y) \sim P_{XY}} f(x, y), \hat{R}_s(f) = \frac{1}{n} \sum_{(x,y) \in S} f(x, y)$$

Definition 1. Rademacher Complexity of \mathcal{F}

$$\mathbb{R}(\mathcal{F} \circ S) := \frac{1}{n} \mathbb{E}_{\vec{\sigma}} \sup_{f \in \mathcal{H}} \vec{\sigma}^T \vec{f}_s$$

Depends on S , (also n)

$$[\sigma_1, \dots, \sigma_n] := \vec{\sigma}, \sigma_i \in \{-1, 1\}, P(\sigma_i = 1) = \frac{1}{2} \forall i$$

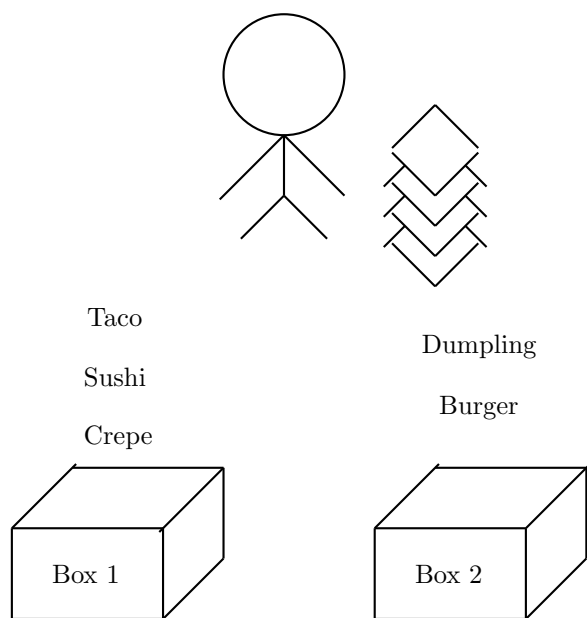
$$\begin{bmatrix} f(x_1, y_1) \\ \dots \\ f(x_n, y_n) \end{bmatrix} (x_i, y_i) \in S, f \in \mathcal{F}$$

Theorem 3. 26.5 (3) *loss bound*, $|\ell(\cdot)| \leq c$

$$\forall h^* \in \mathcal{H}, \text{ wp } \geq 1 - \delta, R(\hat{h}_s^{ERM}) \leq R(h^*) + 2\mathbb{R}(\mathcal{F} \circ S) + 5c \sqrt{\frac{2 \log \frac{8}{\delta}}{n}}$$

data dependent bound

12 Lecture 12

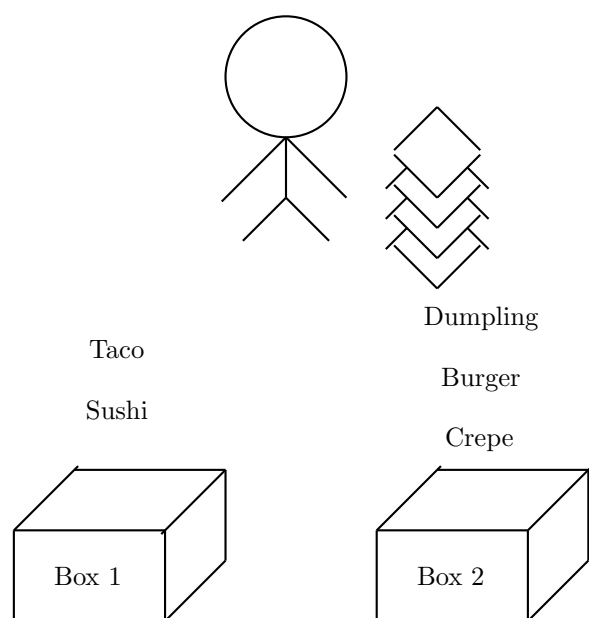
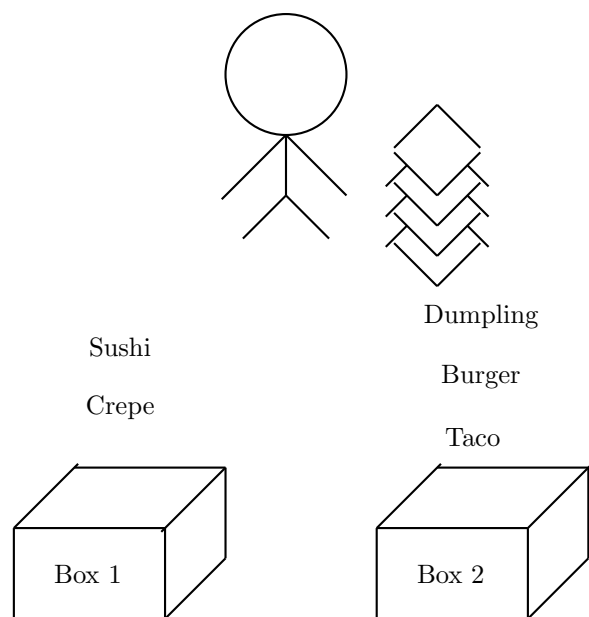


$X = \text{vocabulary}$

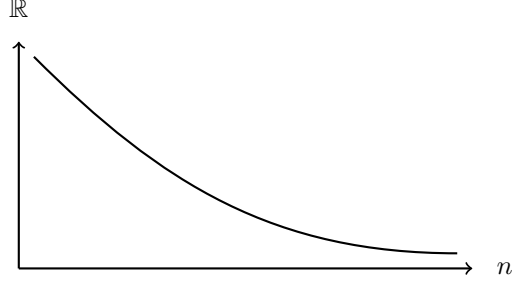
$$h \in \underbrace{\mathcal{H}}_{\text{all rules in your mind}}, X \rightarrow \{ \text{Box1}, \text{Box2} \}$$

$$\frac{1}{m} = \sum_{j=1}^m \sup_{h \in \mathcal{H}} \sum_{i=1}^5 \sigma_i^{(j)} h(x_i)$$

$$\approx \mathbb{E}_{\sigma_1 \dots \sigma_5} \sup_{h \in \mathcal{H}} \vec{\sigma}^T \vec{h}$$



$\mathbb{R}(\mathcal{H})$



$$A \subset \mathbb{R}^n$$

$$\mathbb{R}(A) = \frac{1}{n} \mathbb{E}_{\sigma_1 \dots \sigma_n} \sup_{a \in A} \vec{\sigma}^T a$$

Lemma 26.2, $F = \ell \circ \mathcal{H}$

$$\mathbb{E}_s \left[\sup_{f \in \mathcal{F}} R(f) - \hat{R}_s(f) \right] \leq 2 \mathbb{E}_s \mathbb{R}(F \circ S)'$$

Proof:

$$\begin{aligned} \sup_{f \in F} R(f) - \hat{R}_s(f) &\stackrel{\text{ghost } s'}{=} \sup_{f \in F} \mathbb{E}_{s'} \left[\hat{R}_{s'}(f) - \hat{R}_s(f) \right] \\ &\stackrel{\text{Jensen's}}{\leq} \mathbb{E}_{s'} \sup_{f \in F} \left[\hat{R}_{s'}(f) - \hat{R}_s(f) \right] \end{aligned}$$

$$\mathbb{E}_s \left[\sup Rf - \hat{R}_s f \right] \leq \mathbb{E}_{s,s'} \sup_f \left[\hat{R}_{s'} f - \hat{R}_s f \right] \quad (1)$$

$$S = \{z_1, \dots, z_n\}$$

$$S' = \{z'_1, \dots, z'_n\}$$

$$\begin{aligned} (1) \text{ RHS} &= \frac{1}{n} \mathbb{E}_{s',s} \sup_f \sum_{i=1}^n (f(z'_i) - f(z_i)) \stackrel{\text{introduce } \sigma = \{\sigma_1, \dots, \sigma_n\}}{=} \frac{1}{n} \mathbb{E}_{s',s,\sigma} \sup_f \sum_{i=1}^n \sigma_i (f(z'_i) - f(z_i)) \\ &\leq \frac{1}{n} \mathbb{E}_{s',s,\sigma} \left\{ \left[\sup_{f \in F} \sum_i \sigma_i f(z'_i) \right] + \left[\sup_{g \in F} \sum_i \sigma_i (-g(z_i)) \right] \right\} \\ &= \mathbb{E}_s \mathbb{R}(F \circ S) + \mathbb{E}_s \mathbb{R}(F \circ S) = 2 \mathbb{E}_s \mathbb{R}(F \circ S) \\ &= \mathbb{E}_{s',s} \sup_f f(z'_1) - f(z_1) + \sum_{i=2}^n (f(z'_i) - f(z_i)) \\ &\stackrel{s,s' \text{ iid}}{=} \mathbb{E}_{s',s} \sup_f f(z_1) - f(z'_1) + \sum_{i=2}^n (f(z'_i) - f(z_i)) \\ &= \mathbb{E}_{s',s,\sigma_1} \sup_f \sigma_1 (f(z'_1) - f(z_1)) + \sum_{i=2}^n (f(z'_i) - f(z_i)) \end{aligned}$$

$$\mathbb{P}_s \left(\overbrace{\sup_{f \in F} R(f) - \hat{R}_s(f)}^X \geq \varepsilon \right) \leq \frac{\mathbb{E}[\sup]}{\varepsilon} \leq \frac{2\mathbb{E}_s \mathbb{R}(F \circ S)}{\varepsilon} := \delta$$

$$\Rightarrow \varepsilon = \frac{2\mathbb{E}_s \mathbb{R}(F \circ S)}{\delta}$$

Markov ineq r.v. $\boxed{X \geq 0} \forall a > 0$

$$P(X > a) \leq \frac{\mathbb{E}[X]}{a}$$

$$\boxed{\text{VC style } \varepsilon, \varepsilon = \sqrt{\frac{VC(\mathcal{H}) + \log \frac{1}{\delta}}{n}}}$$

McDiarmid's Ineq, $f : V^n \rightarrow \mathbb{R}, V \subseteq \mathbb{R}$

foral $i \in [n], \forall x_i, x'_i \in V, |f(x_1, x_2, \dots, x_i, \dots, x_n) - f(x_1, x_2, \dots, x'_i, \dots, x_n)| \leq c$

Lecture X_1, \dots, X_n be r.v. in V

$$\mathbb{P} \left(|f(X_1, \dots, X_n) - \mathbb{E}f(X_1, \dots, X_n)| \leq c \sqrt{\left(\log \frac{2}{\delta} \right) \frac{n}{2}} \right) \geq 1 - \delta$$

13 Lecture 13

$$\mathbb{P}_s \left(\sup_{f \in F} R(f) - \hat{R}_s(f) > \varepsilon \right)$$

Lemma 26.2

$$\mathbb{E}_s \left[\sup_{f \in F} R(f) - \hat{R}_s(f) \right] \leq 2\mathbb{E}_s \mathbb{R}(F \circ S)$$

Assumption: $|\ell| \leq c$

McDiarmid's Ineq: For f which satisfies $|f(x_1 \dots x_i \dots x_n) - f(x_1 \dots x'_i \dots x_n)| \leq c, \forall x, i, x'_i$

$$\text{wp} \geq 1 - \delta, |f(X_1 \dots X_n) - \mathbb{E}f(X_1 \dots X_n)| \leq C_0 \sqrt{\frac{n \log \frac{2}{\delta}}{2}}$$

$$\sup_{f \in F} R(f) - \hat{R}_s(f) = \sup_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim P_{XY}} \ell(h(x), y) - \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

$$s = (x_1, y_1), \dots, (x_n, y_n)$$

$$\Rightarrow C_0 = \frac{2C}{n}$$

$$\stackrel{\text{McDiarmid}}{\Rightarrow} \text{wp} \geq 1 - \delta, \left| \left[\sup_{f \in F} R(f) - \hat{R}_s(f) \right] - \mathbb{E}_{s'} \left[\sup_{f \in F} R(f) - \hat{R}_{s'}(f) \right] \right| \leq \frac{2c}{n} \sqrt{\frac{n \log \frac{2}{\delta}}{2}} = c \sqrt{\frac{2 \log \frac{2}{\delta}}{n}}$$

$$\forall x, i, x'_i \rightarrow \text{wp} \geq 1 - \delta, \sup_{f \in F} R(f) - \hat{R}_s(f) \leq \mathbb{E}_{s'} \left[\sup_{f \in F} R(f) - \hat{R}_{s'}(f) \right] + c \sqrt{\frac{2 \log \frac{2}{\delta}}{n}}$$

$$\stackrel{\text{Lemma 26.2}}{\leq} 2\mathbb{E}_{s'} \mathbb{R}(F \circ S') + c \sqrt{\frac{2 \log \frac{2}{\delta}}{n}}$$

Thm 26.5 (1)

$\text{ERM } \arg \min_{h \in \mathcal{H}} \hat{R}_s(h)$

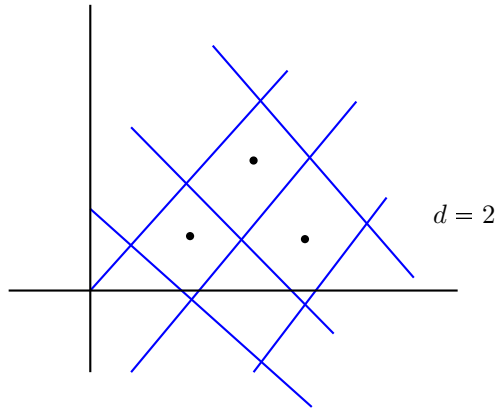
$$\mathcal{H} = \{h_{w,b} : X \rightarrow Y, w \in \mathbb{R}^d, b \in \mathbb{R}, h_{w,b}(x) = \text{sign}(w^T x + b)\}$$

$$X \subseteq \mathbb{R}^d, Y = \{-1, 1\}$$

$$\text{sign}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ -1, & \text{if } z < 0 \end{cases}$$

$$\{x : w^T x + b = 0\}$$

$$VC(\mathcal{H}) = d + 1$$



$$\ell = 0 - 1 \text{ loss}$$

$$\text{Assume: } \exists h^* \in \mathcal{H}, \forall (x, y) \in S, y = h^*(x)$$

(Batch) Perceptron Alg

$$\text{Init, } t = 0, w_0 = \vec{0} \in \mathbb{R}^d$$

$$\text{pick } (x, y) \in S$$

if $\text{sign}(w_t^T x) \neq y$

$$w_t = w_t + yx$$

Terminates?!

How many mistakes were made? $\neq \hat{R}_s(h)$

$$\begin{aligned} C &:= \min_{i \in |S|} y_i (w^*)^T x_i \\ C &> 0 \\ w^{**} &:= \frac{w^*}{c} \\ \forall i, y_i (w^{**})^T x_i &\geq 1, \boxed{1} \\ B &:= \min \|w\|, w \in \{w : \forall i \in S : y_i w^T x_i \geq 1\} \\ w^* &\in \arg \min \end{aligned}$$

$$\begin{aligned} \frac{w_t^T w^*}{\|w_t\| \|w^*\|} &\stackrel{\cos}{\leq} 1 \\ w_{t+1}^T w^* - w_t^T w^* &= (w_{t+1} - w_t)^T w^* = (y_t x_t)^T w^* \stackrel{\boxed{1}}{\geq} 1 \\ \sum_{t=0}^T (w_{t+1}^T w^* - w_t^T w^*) &= w_{T+1}^T w^* - \underbrace{w_0^T w^*}_0 \geq T, \boxed{2} \end{aligned}$$

14 Lecture 14

Batch Perceptron (Ch. 9)

Give $(x_i, y_i)_{i=1:n}, Y = \{-1, 1\}$

Assume realizability $\exists w \in \mathbb{R}^d \forall i \in [n], y_i w^T x_i > 0$

$$\begin{aligned} &\Rightarrow \exists w \forall i, y_i w^T x_i \geq 1 \\ \bar{W} &:= \{w : \forall i, y_i w^T x_i \geq 1\} \\ w^* &\in \arg \min_{w \in \bar{W}} \|w\| \\ B &:= \|w^*\| \end{aligned}$$

Alg: $w_0 = \emptyset$

if $y_i w_t^T x_t \leq 0$ (misclassification)

$w_{t+1} = w_t + y_t x_t$ (repeat)

Claim: Alg makes bounded number of mistakes on any sequence x_t, y_t (in training set)

Proof:

$$\cos(w_{t+1}, w^\star) = \frac{w_{t+1}^T w^\star}{\|w_{t+1}\| \cdot \|w^\star\|} \leq 1$$

step 1 : $w_{t+1}^T w^\star$ grows $O(t)$

$$\begin{aligned} w_{t+1}^T w^\star &= (w_t + y_t x_t)^T w^\star = w_t^T w^\star + y_t \underbrace{x_t^T w^\star}_{\geq 1 \text{ (def } w^\star \in \bar{W})} \\ &\Rightarrow w_{t+1}^T w^\star - w_t^T w^\star \geq 1 \\ &\Rightarrow \sum_{t=1}^{T-1} (w_{t+1}^T w^\star - w_t^T w^\star) \geq T \\ &\Rightarrow w_T^T w^\star - w_0^T w^\star \geq T \\ &\underbrace{\Rightarrow}_{w_0 = \emptyset} w_T^T w^\star \geq T \end{aligned}$$

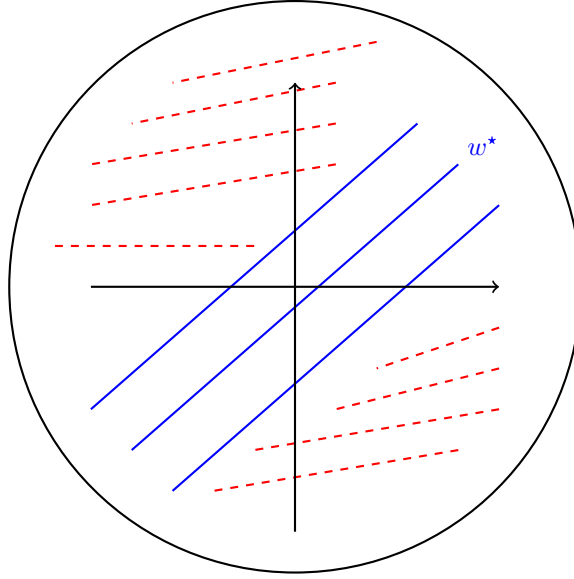
step 2 : $\|w_{t+1}\| \sim o(t)$

$$\begin{aligned} \|w_{t+1}\|^2 &= \|w_t + y_t x_t\|^2 = (w_t + y_t x_t)^T (w_t + y_t x_t) \\ &= w_t^T w_t + 2y_t w_t^T x_t + y_t^2 x_t^T x_t \\ &\underbrace{=}_{Y \in \{-1, 1\}} \|w_t\|^2 + 2 \underbrace{y_t w_t^T x_t}_{\text{misclassification} \leq 0} + \|x_t\|^2 \\ &\leq \|w_t\|^2 + \|x_t\|^2 \\ &\Rightarrow \|w_{t+1}\|^2 - \|w_t\|^2 \leq \|x_t\|^2 \leq R^2 \text{ (Assumption 2, } \forall i \in [n], R \geq \|x_t\|) \\ &\Rightarrow \|w_T\|^2 \leq TR^2 \\ &\Rightarrow \|w_T\| \leq \sqrt{T}R \end{aligned}$$

$$1 \geq \frac{w_{t+1}^T w^\star}{\|w_{t+1}\| \cdot \|w^\star\|} \underbrace{\geq}_{\text{steps 1,2}} \frac{T}{\sqrt{T} \cdot R \cdot B} = \frac{\sqrt{T}}{RB}$$

$$\sqrt{T} \leq RB$$

$$T \leq R^2 B^2$$



$$S = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$s \in \bigcup_{T=1}^{\infty} S^T$$

$$X, Y \in \{-1, 1\}$$

$$s \in \bigcup_{T=1}^{\infty} (X \times Y)^T$$

Interaction protocol:

0 world picks s

At time t :

1 world shows x_t in S

2 learner predict $\hat{y}_t \in Y$

3 world reveals y_t in S (learner "learns, updates")

ADD Realizability assumption: \mathcal{H}

$$\exists h^* \in \mathcal{H} : \forall s_x \in \bigcup_{T=1}^{\infty} X^T, \forall x_t \in s_x, y_t = h^*(x_t)$$

Special case: $|\mathcal{H}| < \infty$

Version space

$$VS := \{h \in \mathcal{H} : h \text{ is consistent with all data so far} \}$$

15 Lecture 15

Online learning:

→ realizable assumption: mistake bound

→ randomization: regret

Given $\mathcal{H}, \exists h^* \in \mathcal{H}$ in each iteration t : world chooses $x_t \in X$ (not necessarily iid), learner predicts $\hat{y}_t \in Y$, world rewards $y_t := h^*(x_t)$, learner incurs error $\mathbb{1}_{(\hat{y}_t \neq y_t)}$, "learns"

Given sequence $S = x_1, \dots, x_n$, alg A , define the number of mistakes A makes on S by $M_A(S)$

If $\exists A : \max_{S \in \bigcup_{n=0}^{\infty} X^n} M_A(S)$ is finite (mistake bound), then \mathcal{H} is online-learnable.

Also, what's $\min_A \max_{S \in \bigcup_{n=0}^{\infty} X^n} M_A(S)$?

"Baby steps" Assume $|\mathcal{H}| < \infty$

$\exists A$ consistent or A version space, maintains a version space

init: $VS_0 = \mathcal{H}$

At time t , pick any $h \in VS_t$, predict $\hat{y}_t = h(x_t)$

Update $VS_{t+1} = \{h' \in VS_t : h'(x_t) = y_t\}$

$$\begin{aligned} X &= \{x^1, \dots, x^m\} \\ h_1 &= 1, 0, \dots, 0 \\ h_2 &= 0, 1, 0, \dots, 0 \\ &\dots \\ h_m &= 0, \dots, 0, 1 \\ h_{m+1} &= h^* = 0, \dots, 0 \\ M &\leq |\mathcal{H}| - 1 \end{aligned}$$

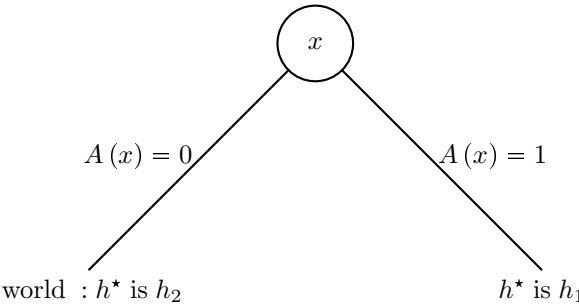
$$\begin{aligned} &A_{\text{halving}} \\ \hat{y}_t &= \arg \max_{y \in \{0,1\}} |\{h \in VS : h(x_t) = y\}| \\ 1 &\stackrel{h^* \in VS}{\leq} |VS| \leq \frac{|\mathcal{H}|}{2^M} \Rightarrow 2^M \leq |\mathcal{H}|, M \leq \log_2 |\mathcal{H}| \end{aligned}$$

$$d := \max_{S \in \bigcup_{n=0}^{\infty} X^n} M_A(S) \text{ (no longer assume } |\mathcal{H}| < \infty \text{) (still } \exists h^* \in \mathcal{H} \text{)}$$

Fix any A . Assume $\exists x \in X$ such that $\exists h_1, h_2 \in \mathcal{H}$ such that $h_1(x) \neq h_2(x)$

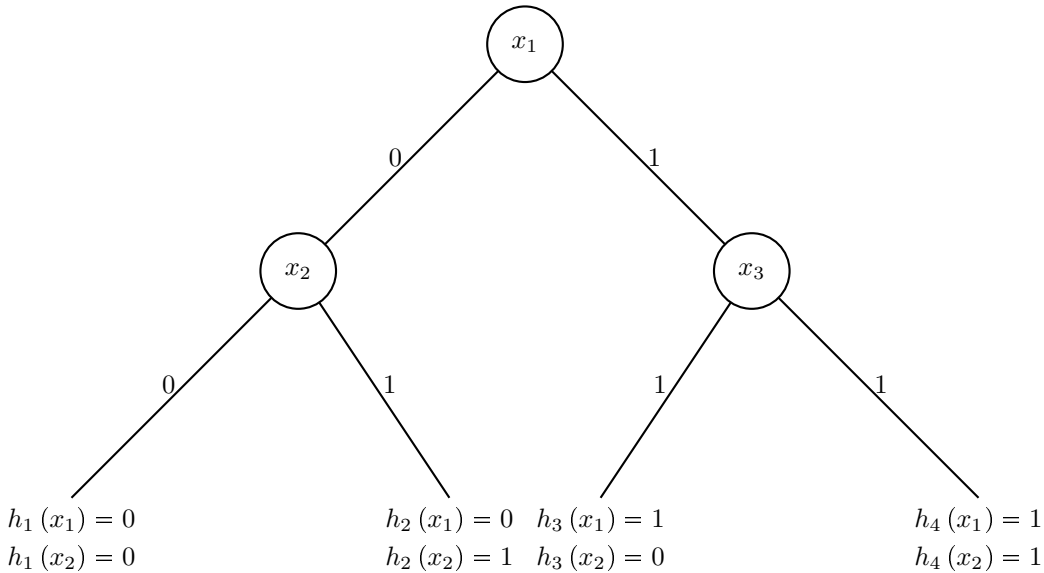
$$h_1(x) = 0$$

$$h_2(x) = 1$$



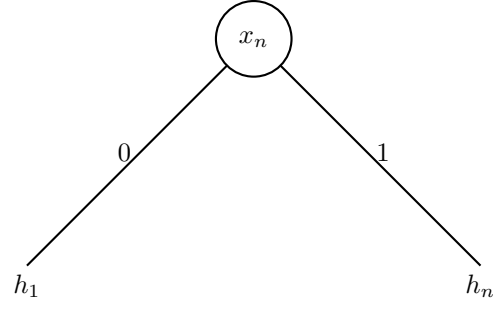
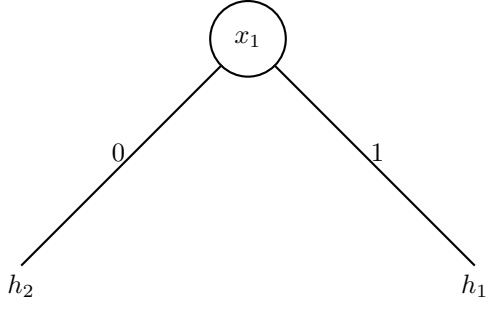
World wants

\mathcal{H} such that \exists



World asks \mathcal{H} "is there a tree with h at the leaves?"

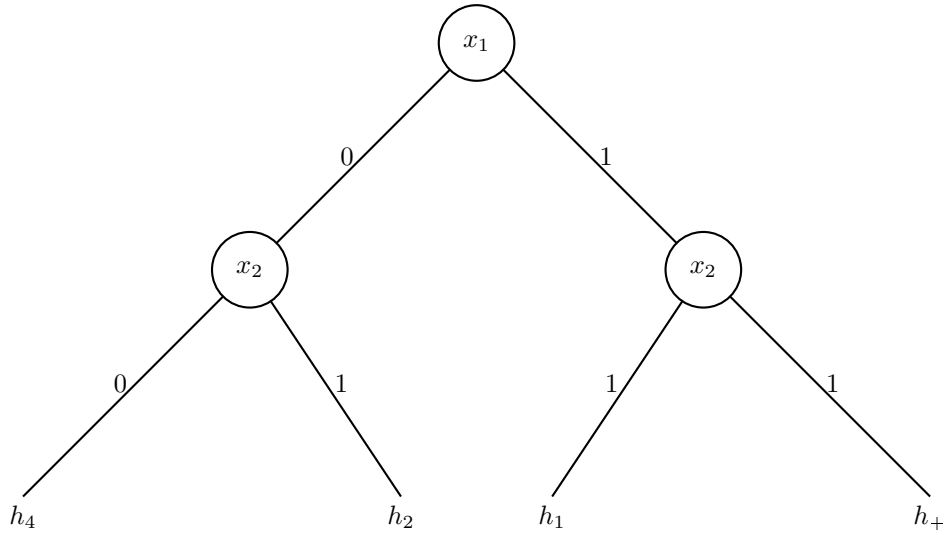
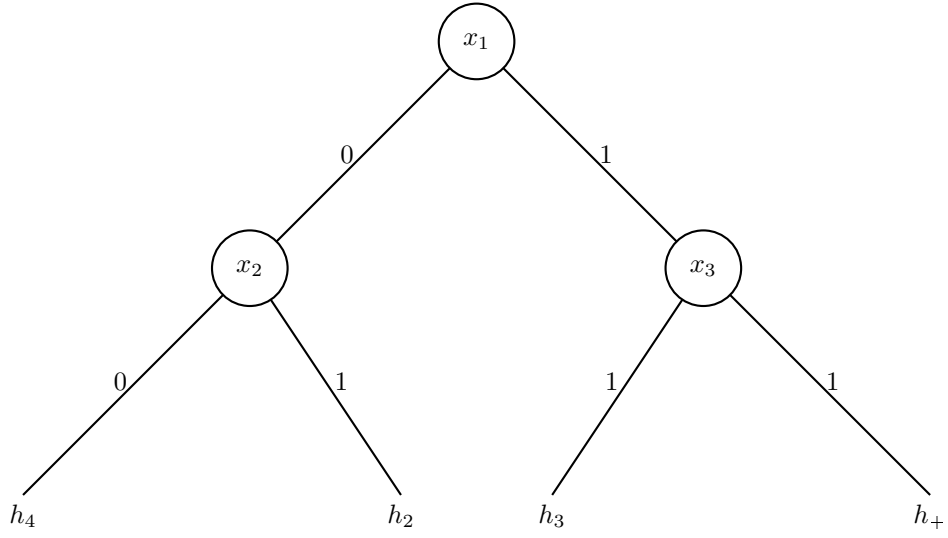
.	x_1	x_2	x_3	...
h_1	0	0	*	
h_2	0	1	*	
h_3	1	*	0	
h_4	1	*	1	
...				



$$\mathcal{H} = \{h_i, i \in [n]\}$$

$$h_i^* = \mathbb{1}_{(x=x_i)} \quad \forall i \in [n]$$

$$\mathcal{H}^2 = \mathcal{H} \cup \{h_+ : h_+(x) = 1, \forall x\}$$



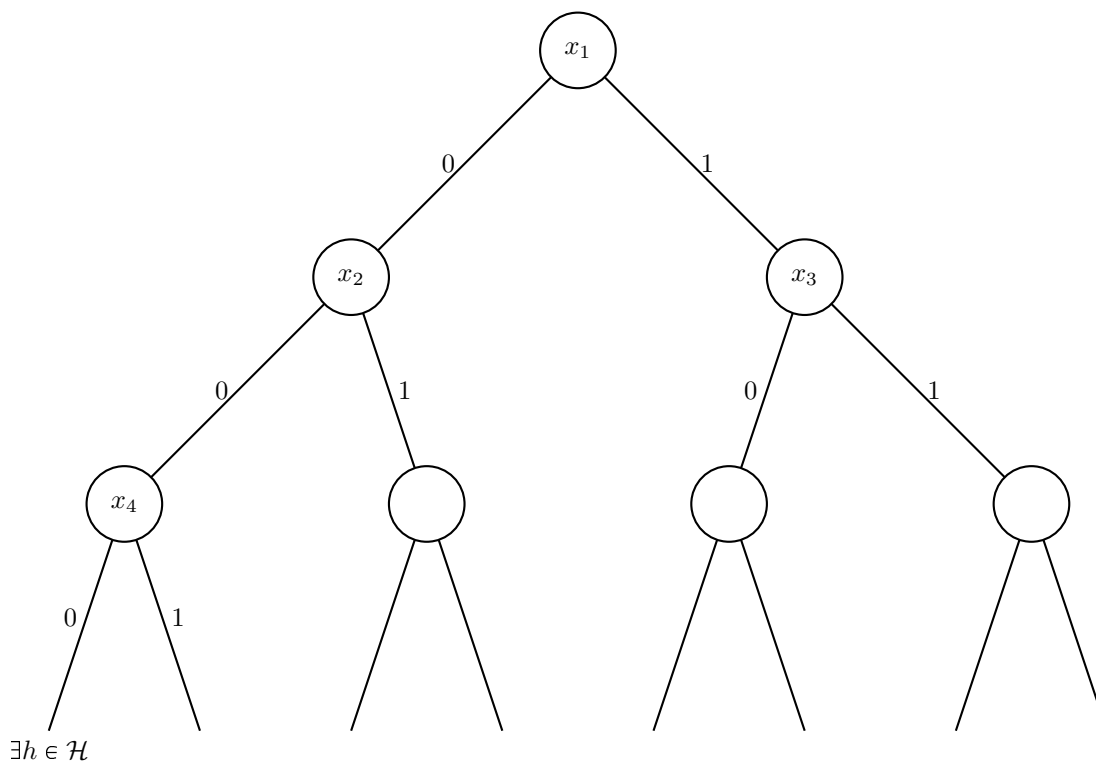
What is the deepest complete tree with h at the leaves?
 depth = Littlestone dimension of \mathcal{H} .

16 Lecture 16

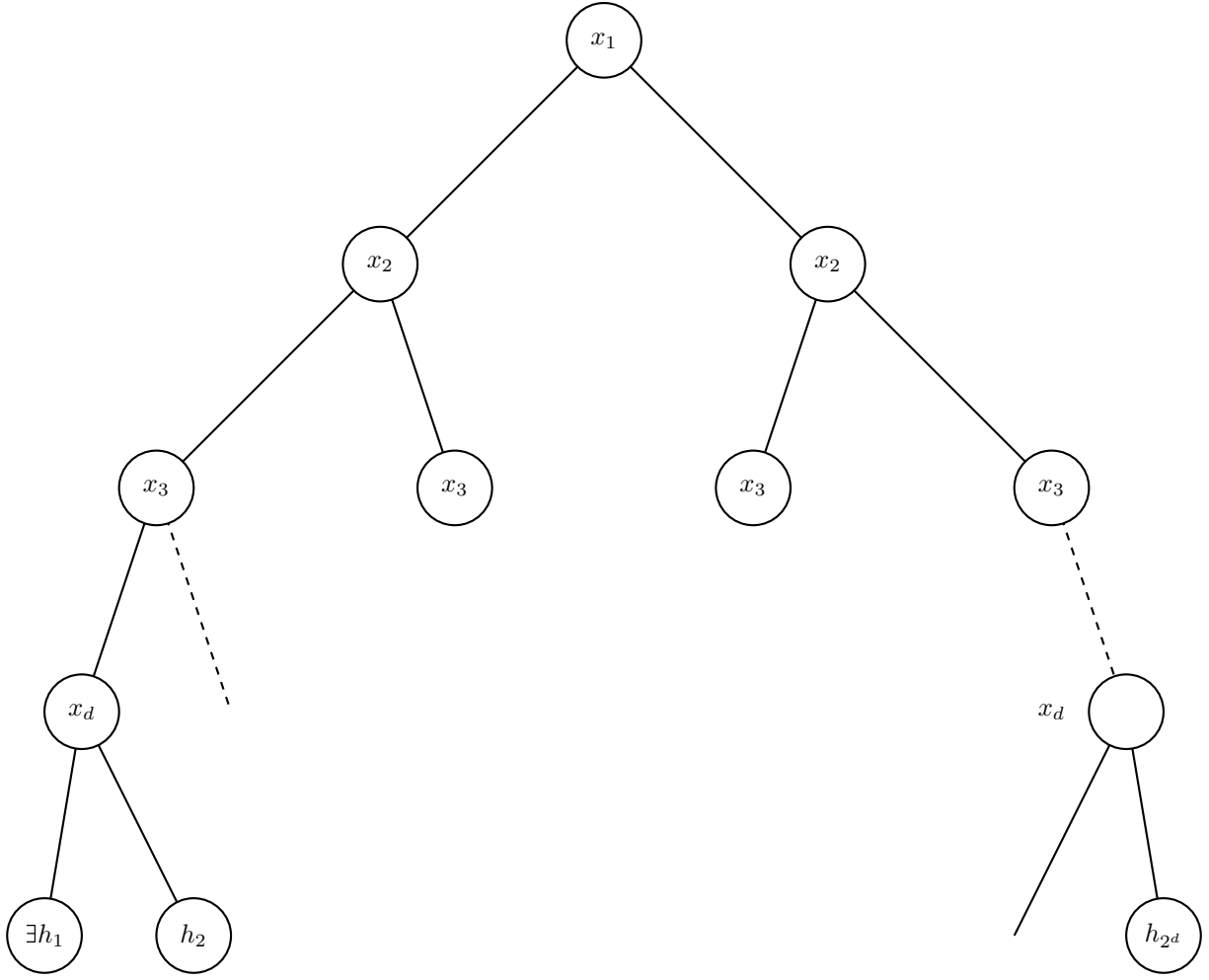
Realizable online learning

$$\max_A \max_{h^* \in \mathcal{H}} \max_{s \in \bigcup_{n=0}^{\infty} X^n} M_A(S) \geq \text{depth of "that tree"} := Ldim(\mathcal{H})$$

Given X, \mathcal{H} , find the deepest complete tree of the form



Ex. Suppose $VC(\mathcal{H}) = d$.
 What about $Ldim(HsC)$?
 $\exists x_1 \dots x_d$ shattered by \mathcal{H} .

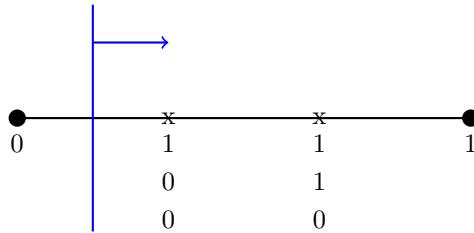
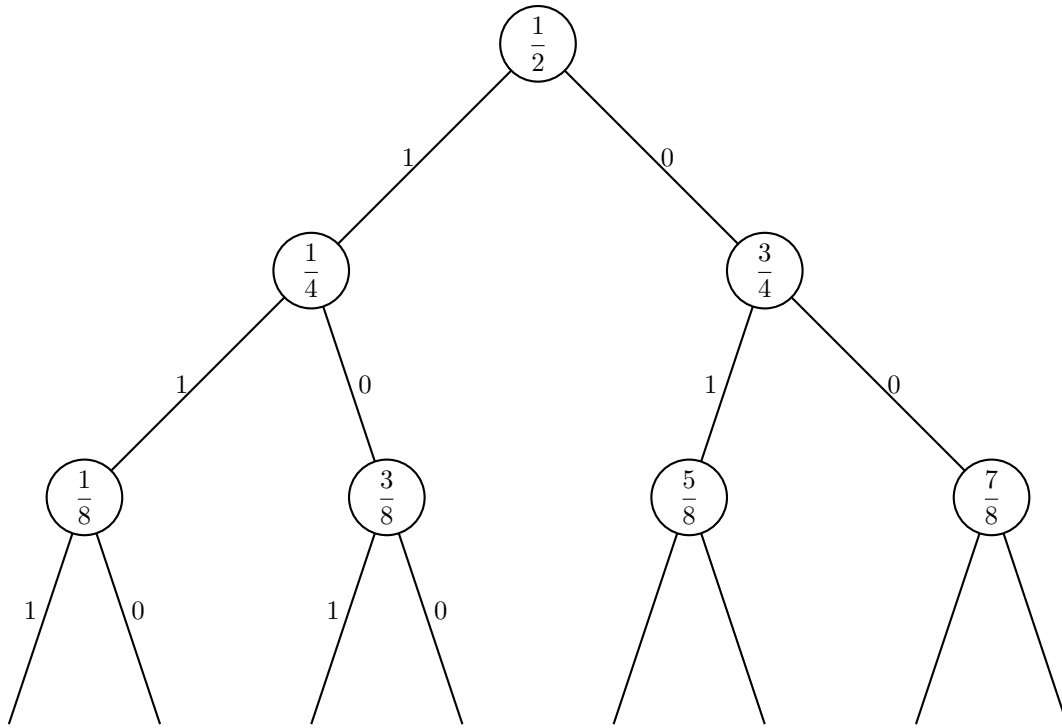


Theorem 4. $Ldim(\mathcal{H}) \geq VC(\mathcal{H})$

Ex: "One-hot" $\{h_x : x \in C\} := \mathcal{H}$
 $Ldim(\mathcal{H}) = 1$ possibly $|\mathcal{H}| = \infty$
 Ex: $\mathcal{H} = \{h_a, a \in [0, 1] : h_a(x) = \mathbb{1}_{(x \geq a)}\},$

$$VC(\mathcal{H}) = 1$$

$$Ldim(\mathcal{H}) = \infty$$



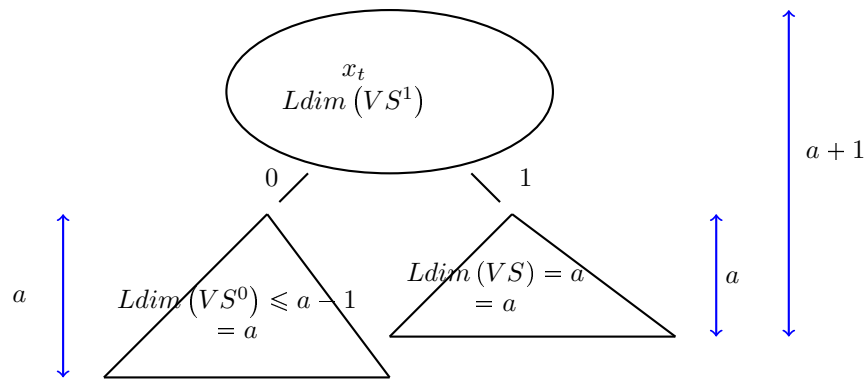
Standard Optimal Alg: A_{SOA}
 receive x_t

$$VS^0 \cup VS^1, VS^p = \{h \in VS : h(x_t) = p\}$$

$$\hat{y}_t = \arg \max_{p \in Y} Ldim(VS^p)$$

receive y_t

$$VS \leftarrow VS^{y_t}$$



World: x_t
 Alg: predicts: \hat{y}_t
 World: give $y_t \in Y$

$$\text{regret} : \sum_{t=1}^T \ell(\hat{y}_t(x_t), y_t) - \min_{h \in \mathcal{H}} \sum_{t=1}^T \ell(h(x_t), y_t) \geq \text{Show this! } \frac{T}{2}$$

17 Lecture 17

No longer assume $\exists h^* \in \mathcal{H}$ such that $y_t = h^*(x_t)$
 (expected) regret wrt \mathcal{H} (want: $o(T)$, in fact \sqrt{T})

$$\sup_{(x,y)_{1:T}} \mathbb{E} \sum_{t=1}^T \mathbb{1} \left(\underbrace{\hat{y}_t(x_t)}_{\text{made by } A} \neq y_t \right) - \min_{h \in \mathcal{H}} \sum_{t=1}^T \mathbb{1}_{(h(x_t) \neq y_t)}$$

Ex: $\mathcal{H} = \{h_0, h_1\}$, A deterministic, regret $\geq \frac{T}{2}$
 A randomized

17.1 Weighted Majority

(learn from expert advice)

Given d experts, horizon T , stepsize $\eta > 0$

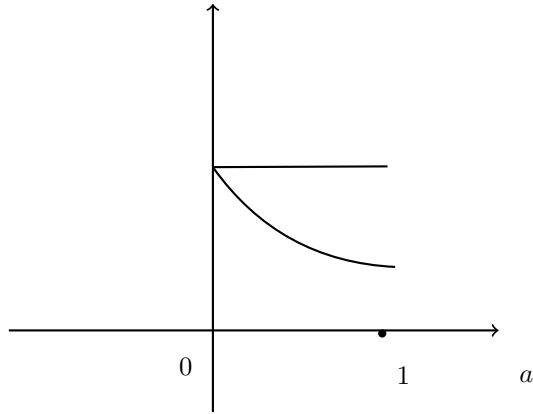
1. init $\tilde{w}^{(1)} = \underbrace{(1, \dots, 1)}_d$ (unnormalized weights)
2. For $t = 1, 2, \dots, T$
3. $w^{(t)} = \frac{\tilde{w}^{(t)}}{Z_t}$ where $Z_t = \sum_{i=1}^d \tilde{w}_i^{(t)}$
4. Choose expert $i \sim w^{(t)}$

5. Observe loss vector $V_t = \left(v_{t_1}, \dots, \overbrace{\boxed{v_{t_i}}}^{MAB}, \dots, v_{t_d} \right) \in [0, 1]^d$ "bounded", pay expected loss $w^{(t)T} V_t$
6. update $\tilde{w}_j^{(t+1)} = \tilde{w}_j^{(t)} e^{-\eta v_{t_j}} \quad \forall j \in [d]$

$$\left(\sum_{t=1}^T w^{(t)T} V_t \right) - \min_{j \in [d]} \sum_{t=1}^T V_{t_j} \leq \sqrt{2 (\log d) T}$$

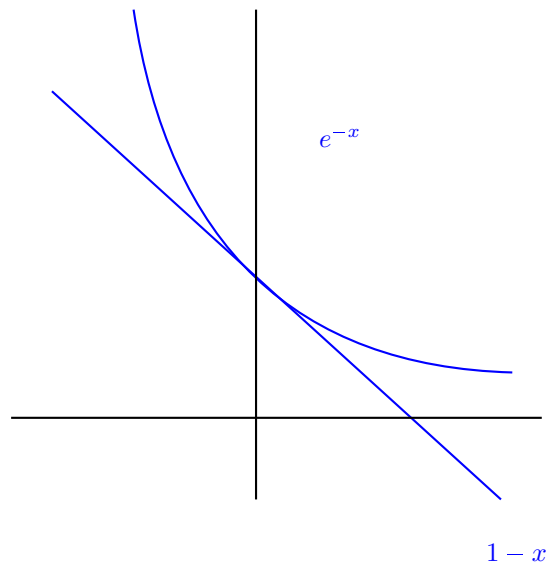
Proof:

$$\begin{aligned} \log \frac{Z_{t+1}}{Z_t} &= \log \frac{\sum_{j=1}^d \tilde{w}_j^{(t)} e^{-\eta V_{t_j}}}{Z_t} \\ &= \log \sum_j w_j^{(t)} e^{-\eta V_{t_j}} \\ &\leq \log \sum_j w_j^{(t)} \left[1 - \eta V_{t_j} + \frac{\eta^2 V_{t_j}^2}{2} \right] \\ &= \log \left[1 - \underbrace{\sum_j w_j^{(t)} \left(\eta V_{t_j} - \frac{\eta^2 V_{t_j}^2}{2} \right)}_{:= b \in (0,1)} \right] \\ &\leq \log e^{-b} \\ &= -\eta \left(w^{(t)T} V_t \right) + \eta^2 \sum_j w_j^{(t)} \frac{V_{t_j}^2}{2} \\ &\leq -\eta w^{(t)T} V_t + \frac{\eta^2}{2} \end{aligned}$$



$$a \in (0, 1)$$

$$e^{-a} \leq 1 - a + \frac{a^2}{2}$$



$$1 - x \leq e^{-x}$$

18 Lecture 18

18.1 Online learning

(no assumption on $\exists h^* \in \mathcal{H}, y_t = h^*(x_t)$)

- subroutine: wighted Majority

Input: $d = \#$ experts, T rounds

init $\tilde{w}^{(1)} = (1, \dots, 1)_d$

for $t = 1, 2, \dots$

$$Z_t = \tilde{w}^{(t)^T} \mathbf{1}, w^{(t)} = \frac{\tilde{w}^{(t)}}{Z_t}$$

pick $i \sim w^{(t)}$ to predict

suffer expected loss $w^{(t)^T} V^{(t)}$

$$\tilde{w}^{(t+1)} = \tilde{w}^{(t)} e^{-\eta V_i^{(t)}} \quad \forall i = 1 \dots d$$

$$V_i^{(t)} := \ell \left(\overbrace{\text{expert}_i(x_t)}^{h_i(x_t)}, y_t \right)$$

$$V^{(t)} := \begin{bmatrix} V_1^{(t)} \\ \dots \\ V_d^{(t)} \end{bmatrix} \in [0, 1]^d$$

Theorem 5. $\left(\sum_{t=1}^T w^{(t)T} V^{(t)} \right) - \left(\min_{\substack{i \in [d] \\ \text{"best expert"}}} \sum_{t=1}^T V_i^{(t)} \right)$

$$\begin{aligned} \log \frac{Z_{t+1}}{Z_t} &= \log \sum_i^d w_i^{(t)} e^{-\eta V_i^{(t)}} \\ &\begin{bmatrix} e^{-x} \leq 1-x + \frac{x^2}{2} \\ \leq \end{bmatrix} \log \left[\sum_i^d w_i^{(t)} \left[1 - \eta V_i^{(t)} + \frac{\eta^2 V_i^{(t)^2}}{2} \right] \right] \\ &= \log \left[1 - \sum_i w_i^{(t)} \left(\eta V_i^{(t)} - \frac{\eta^2 V_i^{(t)^2}}{2} \right) \right] \\ &\stackrel{1-x \leq e^{-x}}{\leq} \log e^{-\sum_i w_i^{(t)} \left(\eta V_i^{(t)} - \frac{\eta^2 V_i^{(t)^2}}{2} \right)} \\ &= -\eta w^{(t)T} V^{(t)} + \sum_i w_i^{(t)} \frac{\eta^2 V_i^{(t)^2}}{2} \\ &\stackrel{V \in [0,1]}{\leq} -\eta w^{(t)T} V^{(t)} + \frac{\eta^2}{2} \leftarrow \boxed{1} \end{aligned}$$

Telescope $\boxed{1}$ over t :

$$\begin{aligned} \log Z_{T+1} - \log d &= \log Z_{T+1} - \log Z_1 \\ &= \sum_{t=1}^T \log \frac{Z_{t+1}}{Z_t} \\ &\stackrel{\boxed{1}}{\leq} -\eta \sum_t w^{(t)T} V^{(t)} + \frac{\eta^2 T}{2} \leftarrow \boxed{2} \end{aligned}$$

$$\begin{aligned} Z_{T+1} &= \sum_i \tilde{w}_i^{(T+1)} \\ &= \sum_i 1 \cdot e^{-\eta \sum_t V_i^{(t)}}, \\ \log Z_{T+1} &= \log \sum_i e^{-\eta \sum_t V_i^{(t)}} \end{aligned}$$

$$\begin{aligned}
&\geq \max_i \log e^{-\eta \sum_t^T V_i^{(t)}} \\
&= \max_i -\eta \sum_t^T V_i^{(t)} \\
&= -\eta \left(\min_i \sum_t^T V_i^{(t)} \right) \leftarrow \boxed{3}
\end{aligned}$$

”best expert”

$$\begin{aligned}
&\boxed{2}, \boxed{3} \Rightarrow \\
&-\eta \left(\min_i \sum_t^T V_i^{(t)} \right) - \log d \stackrel{\boxed{3}}{\leq} \log Z_{T+1} - \log d \\
&\leq -\eta \sum_t^T w^{(t)^T} V^{(t)} + \frac{\eta^2 T}{2} \\
&\sum_t^T w^{(t)^T} V_i^{(t)} - \min_{i \in [d]} \sum_t^T V_i^{(t)} \leq \frac{\log d}{\eta} + \frac{\eta T}{2} \rightarrow \text{sublinear}
\end{aligned}$$

$$\begin{aligned}
&\frac{\partial \text{RHS}}{\partial \eta} = 0 \\
&\Rightarrow -\frac{\log d}{\eta^2} + \frac{T}{2} = 0 \\
&\Rightarrow \frac{1}{\eta} = \sqrt{\frac{T}{2 \log d}} \\
&\Rightarrow \text{RHS} = \sqrt{2T \log d}
\end{aligned}$$

If $|\mathcal{H}| < \infty$

$$\text{regret} \leq \sqrt{2T \log |\mathcal{H}|}$$

when $|\mathcal{H}| = \infty$?

- subroutine 2

SOA: $VS = \mathcal{H}$

for $t = 1, 2, \dots$

receive x_t

$$\begin{aligned}
&VS^0 \cup VS^1 \\
&\hat{y}_t = \arg \max_y Ldim(VS^y)
\end{aligned}$$

$$VS \leftarrow VS^{y_t}$$

assuming $\exists h^* \in \mathcal{H}$

s.t. $y_t = h^*(x_t) \ \forall t$

then SOA mistakes $\leq Ldim(\mathcal{H})$

expert (i_1, i_2, \dots, i_L)

$$1 \leq i_1 < i_2 < \dots i_L \leq T, L \leq Ldim(\mathcal{H}) < \infty$$

init $VS = \mathcal{H}$

for $t = 1, 2, \dots, T$

receive x_t

$VS^0 = \{h \in VS, h(x_t) = 0\}$, same for VS^1

if $t \in \{i_1, \dots, i_L\}$

$$\hat{y}_t = \arg \min_y Ldim(VS^y) \leftarrow \text{anti-SOA}$$

else

$$\hat{y}_t = \arg \max_y Ldim(VS^y) \leftarrow \text{SOA}$$

$$VS \leftarrow VS^{\hat{y}_t}$$

Given $x_1 \dots x_T$

$$\forall h \in \mathcal{H}$$

$$\Rightarrow h(x_1) \dots h(x_T)$$

want: $\exists L, i_1 \dots i_L$ s.t. expert $(i_1 \dots i_L)$ produces the same predictions

Run SOA on input

$$x_1, h(x_1)$$

...

$$x_T, h(x_T)$$

Lemma 21.13

19 Lecture 19

19.1 Online Convex Optimization

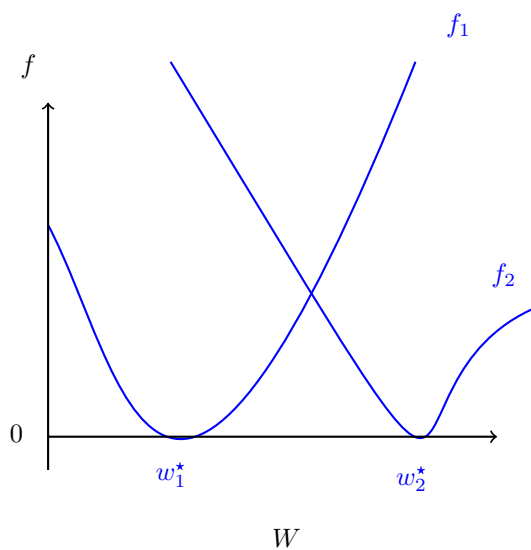
for $t = 1, 2, \dots$

learner chooses $w^{(t)} \in \text{Convex set } W$

environment chooses loss function $f_t : W \rightarrow \mathbb{R}$ convex (subgradient)

learner suffers $f_t(w^{(t)})$

$$\text{regret} := \sum_{t=1}^T f_t(w^{(t)}) - \inf_{w \in W} \sum_{t=1}^T f_t(w)$$



Online gradient descent

$$w^{(1)} = 0$$

for $t = 1, 2, \dots$

predict $w^{(t)}$

receive $f_t(\cdot)$, suffer $f_t(w^{(t)})$

$$w^{(t+1)} = \text{Proj}_W \left[\underbrace{w^{(t)} - \underbrace{\eta}_{\text{stepsize}} \cdot \underbrace{\partial}_{\text{subgradient}} f_t(w^{(t)})}_{W^{(t+1/2)}} \right]$$

Theorem 6. (Thm 21.15)

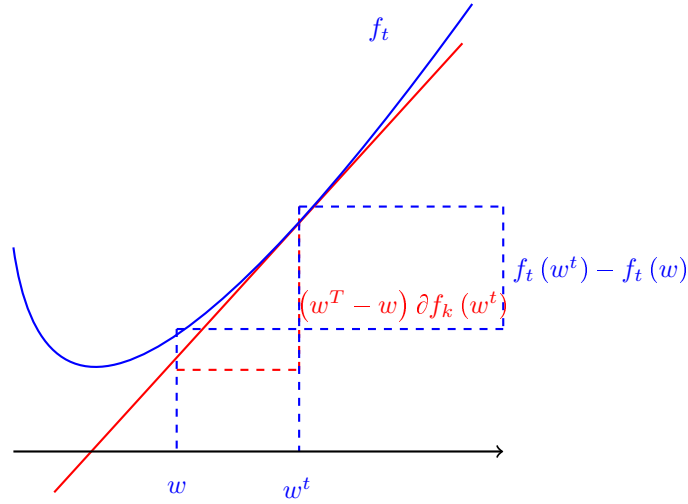
$$\text{Regret} \leq \inf_{w \in W} \frac{\|w\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\partial f_t(w^{(t)})\|^2$$

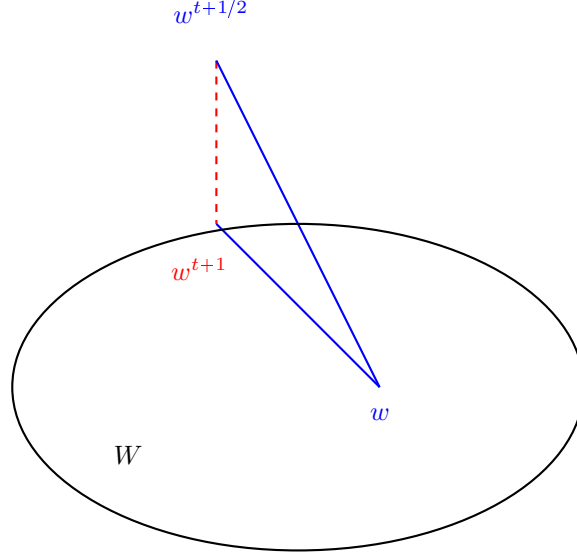
$$\text{OR } \forall w \in W, \text{Regret}(w) \leq \frac{\|w\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\partial f_t(w^{(t)})\|^2$$

Proof. Fix $w \in W$

$$\begin{aligned} \|w^{(t+1)} - w\|^2 - \|w^{(t)} - w\|^2 &= \left(\|w^{t+1} - w\|^2 - \|w^{t+1/2} - w\|^2 \right) + \left(\|w^{t+1/2} - w\|^2 - \|w^t - w\|^2 \right) \\ &\leq 0 + \left(\|w^t - \eta \partial f_t(w^t) - w\|^2 - \|w^t - w\|^2 \right) \\ &= \left(-2(w^t - w)^T \eta \partial f_t(w^t) + \eta^2 \|\partial f_t(w^t)\|^2 \right) \\ &\leq 2\eta(f_t(w^t) - f_t(w)) + \eta^2 \|\partial f_t(w^t)\|^2 \leftarrow \boxed{1} \end{aligned}$$

□





$$\begin{aligned}
\sum_{t=1}^T \text{LHS } \boxed{1} &= \left\| w^{(T+1)} - w \right\|^2 - \left\| w^{(1)} - w \right\|^2 \\
&\leq \sum_{t=1}^T \text{RHS } \boxed{1} = 2\eta \sum_{t=1}^T (f_t(w^t) - f_t(w)) + \eta^2 \sum_{t=1}^T \left\| \partial f_t(w^t) \right\|^2 \\
\sum_{t=1}^T f_t(w^t) - \sum_{t=1}^T f_t(w) &\leq \frac{-\left\| w^{T+1} - w \right\|^2 + \left\| w \right\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \left\| \partial f_t(w^t) \right\|^2 \\
&\leq \frac{\left\| w \right\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \left\| \partial f_t(w^t) \right\|^2 \leftarrow \boxed{2}
\end{aligned}$$

Further assumptions:

1. W is norm bounded: $\forall w \in W, \|w\| \leq B$
2. $f_t, \forall t$ is ρ Lipschitz $\|\partial f_t(w)\| \leq \rho \forall w \in W$

$$\begin{aligned}
\boxed{2} &\stackrel{\text{Assump 1.2}}{\Rightarrow} \sum_{t=1}^T f_t(w^t) - \sum_{t=1}^T f_t(w) \leq \frac{B^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \rho^2 \\
&= \underbrace{\frac{B^2}{2\eta} + \frac{\eta T \rho^2}{2}}_{\boxed{B\rho\sqrt{T}}} \\
&\quad - \frac{B^2}{2\eta^2} + \frac{T\rho^2}{2} = 0 \\
\eta &= \sqrt{\frac{B^2}{T\rho^2}} = \frac{B}{\sqrt{T}\rho}
\end{aligned}$$

doubling trick

run OGD on $t = 1$
 run OGD on $t = 2, 3$
 run OGD on $t = 4, 5, 6, 7$

20 Lecture 20

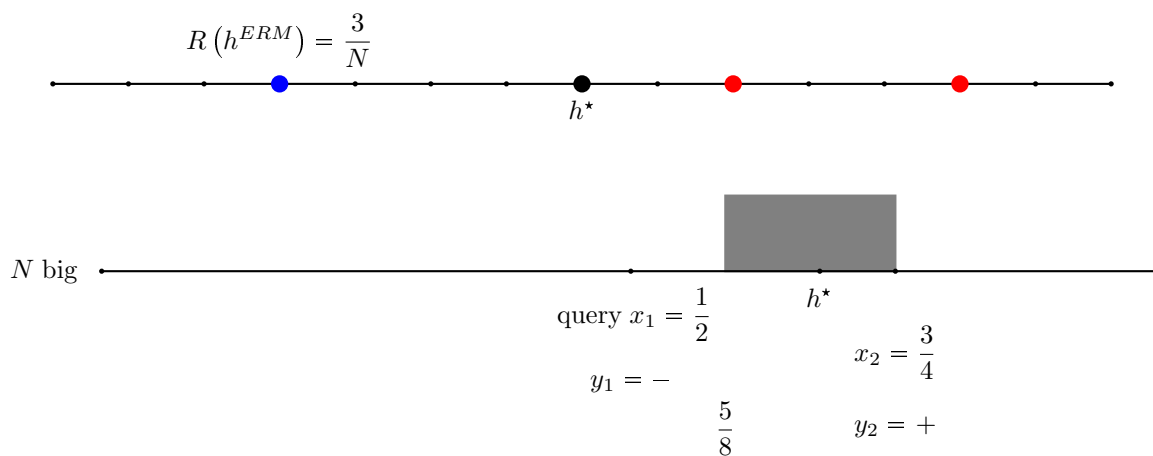
passive learning: $(x, y) \sim P_{XY}$ "unit cost" PAC w.p. $\geq 1 - \delta, n = O\left(\frac{|\mathcal{H}|}{\varepsilon}\right)$

active learning: $x \sim P_X$ free, query x , "oracle" gets $y \sim P_{Y|X=x}$, realizable ($y = h^*(x)$), unit cost, alg can choose x ! adaptively

e.x. $X = 0 : \frac{1}{N} : 1, \mathcal{H} = \{\mathbb{1}_{(x \geq a)}, a \in X\}, P_X \text{ unif}(X) \stackrel{iid}{\sim} n \text{ training items (passive)}$

ERM, $h^{ERM} \in \text{Version Space}, V = \{h \in \mathcal{H} : h(x_i) = y_i, i = 1 : n\}$

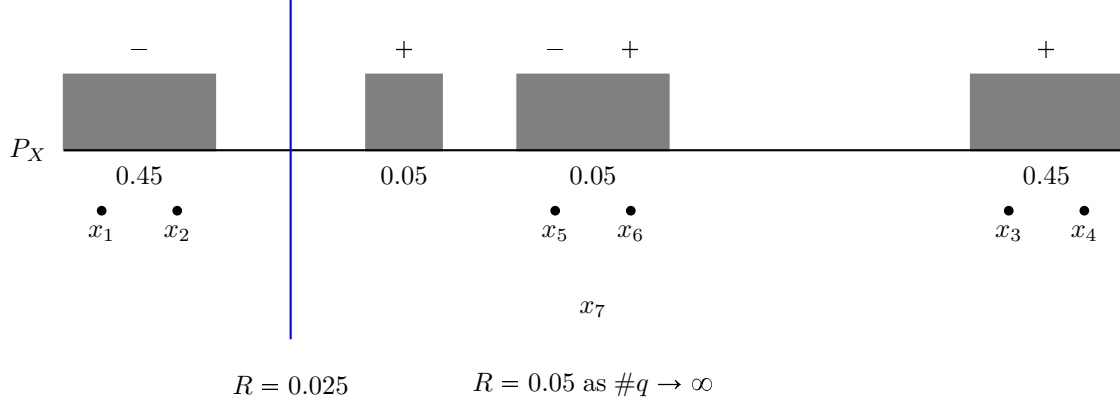
$$R(h^{ERM}) = O\left(\frac{1}{n}\right)$$



n queries

$$n = O\left(\log\left(\frac{1}{\varepsilon}\right)\right) \leq R(h \in VS) = O\left(\frac{1}{2^n}\right)$$

Unertainty-based Active Learning



CAL

init $V_1 = \mathcal{H}$ version space

for epoch $r = 1, 2, \dots$

$$x \sim P_X$$

if V_r disagrees on x (ie $\exists h, h' \in V_r, h(x) \neq h'(x)$)

query x 's label, oracle gives y

$$V_{r+1} = \{h \in V_r, h(x) = y\}$$

init $V_1 = \mathcal{H}$ version space

for epoch $r = 1, 2, \dots, R$

(make sure this happens k times)

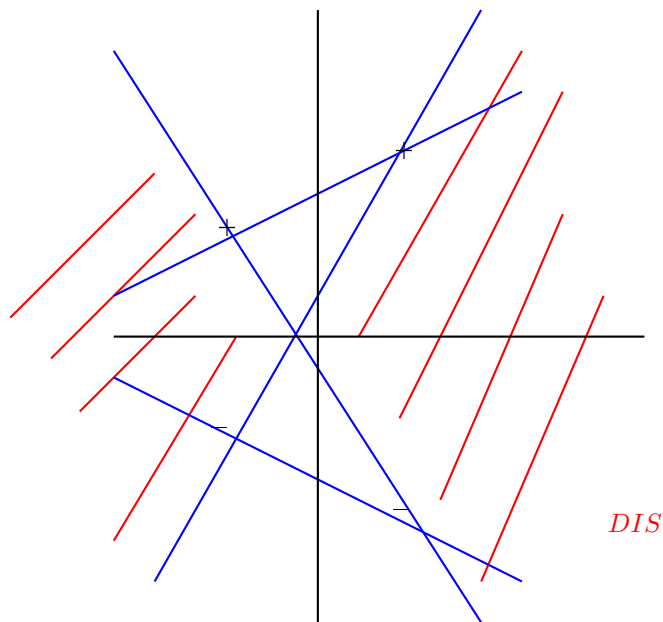
- $x \sim P_X$
- if V_r disagrees on x (ie $\exists h, h' \in V_r, h(x) \neq h'(x)$)
- query x 's label, oracle gives y

$$V_{r+1} = \{h \in V_r, h(x_i) = y_i, i = 1 \dots k\}$$

Output any $h \in V_{R+1}$

version space V_r

Disagreement region $DIS(V_r) := \{x \in X : \exists h, h' \in V_r : h(x) \neq h'(x)\}$



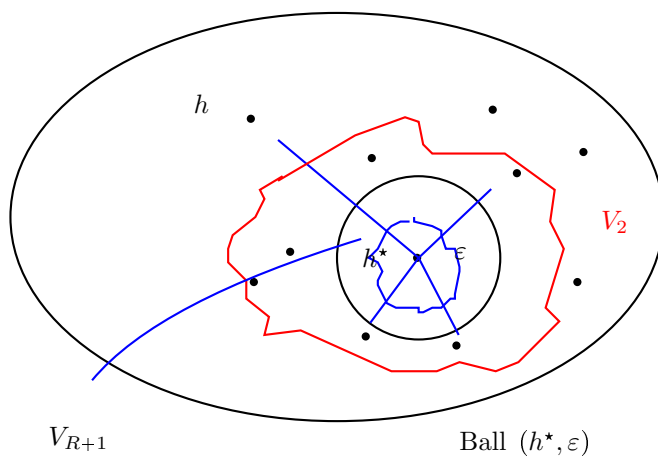
$$\Delta(V_r) := P_X(DIS(V_r))$$

$$R(h^{CAL})$$

pseudometric $d(h, h'), h, h' \in \mathcal{H} = \mathbb{E}_{x \sim P_X} \mathbb{1}_{(h(x) \neq h'(x))}$

$$\Rightarrow R(h) = d(h, h^*)$$

$$\mathcal{H} = V_1$$



21 Lecture 21

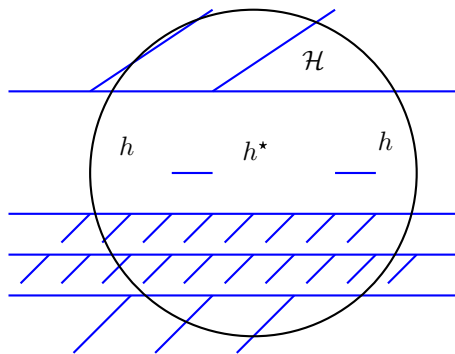
21.1 CAL (mini-batch version)

1. init $V_1 = \mathcal{H}$ (assume $|\mathcal{H}| < \infty$)
2. FOR epochs $i = 1 \dots n$
3. Collect k items $x_1, \dots, x_k \stackrel{iid}{\sim} P_X$ such that they $\in DIS(V_i) = \{x \in X : \exists h, h' \in V_i, h(x) \neq h'(x)\}$
4. query them. Oracle gives their labels $y_1 \dots y_k$
5. $V_{i+1} = \{h \in V_i : h(x_i) = y_i, \forall i \in [k]\}$
6. return any $h \in V_{n+1}$

Want: query complexity

$$O\left(\log \frac{1}{\varepsilon}\right)$$

w.p. $\leq 1 - \delta$



$$P_X(DIS(V_i)) \geq R(h), \forall h \in V_i$$

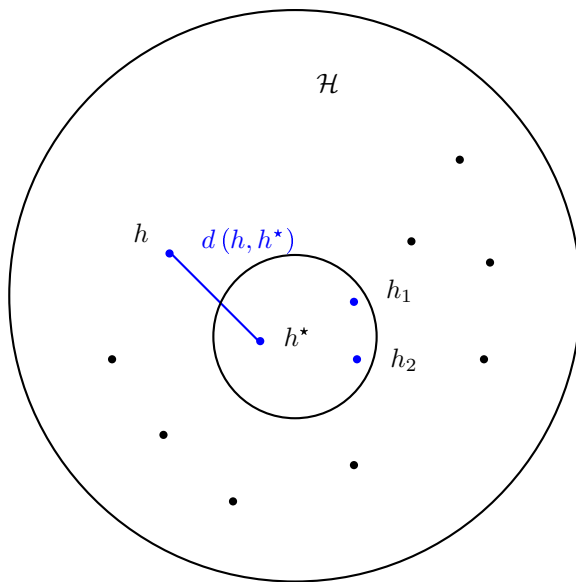
Define (pseudo) metric

$$\begin{aligned} d(h, h') &= P_x(DIS(\{h, h'\})) \\ \Rightarrow d(h, h^*) &= \mathbb{E}_{x \sim P_X} \mathbb{1}_{(h(x) \neq h^*(x))} = R(h) \end{aligned}$$

Want:

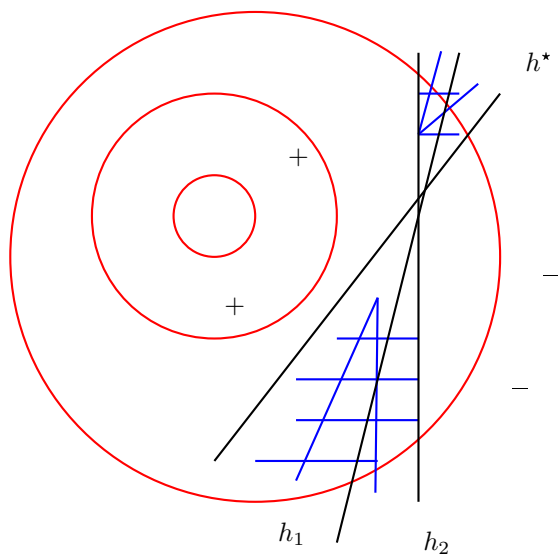
$$\boxed{P_X(DIS(V_{i+1})) \leq \frac{1}{2} P_X(DIS(V_i)) \quad \forall i \in [n]}$$

$$\Rightarrow R(h) \in V_{n+1} \leq P_X(DIS(V_{n+1})) \leq \frac{1}{2^n} := \varepsilon \Rightarrow R(h) < \varepsilon \text{ if } n = \log\left(\frac{1}{\varepsilon}\right)$$



$$B(h^*, r) := \{h \in \mathcal{H}, d(h, h^*) \leq r\}$$

$$P_X(DIS(B(h^*, r)))$$



$$V_i = \{h^*, h_1, h_2\}$$

Step 3 = draw iid $P_{X|V_i}$

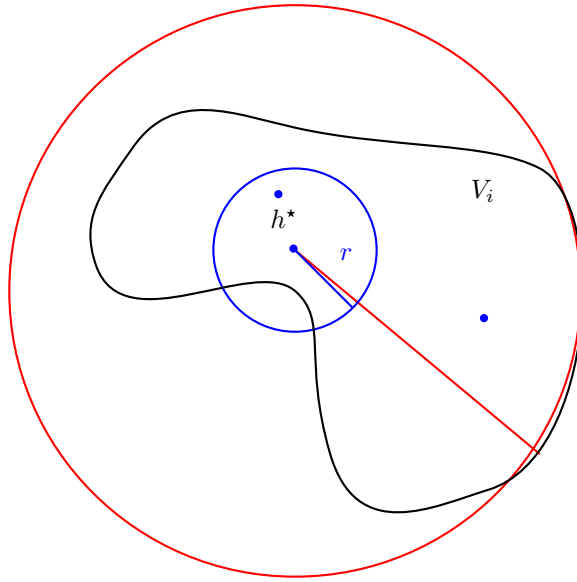
$$P_{X|V_i} = \frac{P(x)}{\int_{x' \in DIS(V_i)} P(x') dx'} = \frac{P(x)}{P_X(DIS(V_i))}$$

Define $Q_{XY} := P_{X|V_i} \cdot P(Y|X)$

$$R_Q(h) = \frac{R_P(h)}{P_X(DIS(V_i))}$$

Want:

$$P_X(DIS(V_{i+1})) \leq f(r_{V_{i+1}}), r_V := \max_{h \in V} d(h, h^*)$$



$$d(h, h^*) = R_P(h) = R_X(DIS(V_i)) \cdot R_Q(h), \forall h \in V_i$$

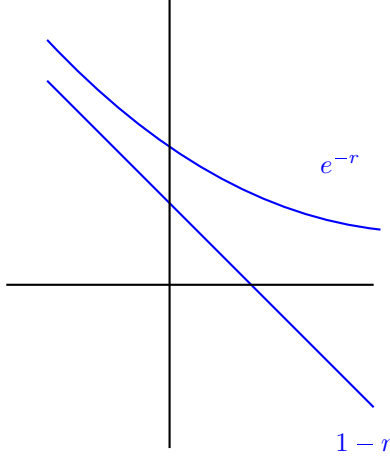
Suppose $h \in V_i$ survives k iid Q

$$\Rightarrow R_Q(h) \text{ small} := r$$

$$|V_i| \cdot \left[(1-r)^k \right] \leq |V_i| e^{-rk} := \frac{\delta}{n}$$

$$-rk = \log \frac{\delta}{n|V_i|}$$

$$\Rightarrow k = \frac{\log \frac{n|V_i|}{\delta}}{r}$$



22 Lecture 22

22.1 Missed Lecture

$$V_i$$

$$Q := P_{X|DIS(V_i)} \cdot P_{Y|X}$$

$$R_Q(h) := \frac{R_p(h)}{P_X(DIS(V_i))}$$

k labeled items $\stackrel{iid}{\sim} Q$

$$wp \geq 1 - \frac{\delta}{n}, \quad \underbrace{\forall h \in V_{i+1}}_{\text{agrees on those } k \text{ items}}, R_Q(h) \leq r \text{ if } k \geq \frac{\log \frac{n|V_i|}{\delta}}{r}$$

$$\Rightarrow \underbrace{R_p(h)}_{=d(h, h^*)} \leq P_X(DIS(V_i)) \cdot r$$

$$\Rightarrow V_{i+1} \subseteq B(h^*, P_X(DIS(V_i)) \cdot r)$$

$$d(h, h^*) = \mathbb{E}_{x \sim P_X} \mathbb{1}_{(h(x) \neq h^*(x))}$$

$$= P_X(DIS(\{h, h^*\}))$$

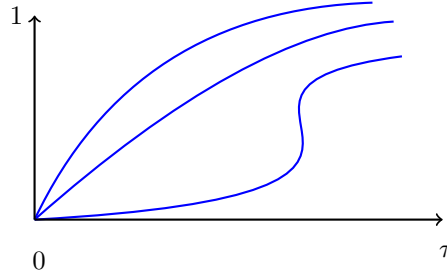
want:

$$P_X(DIS(V_{i+1})) \leq \frac{1}{2} P_X(DIS(V_i))$$

$$\Rightarrow \underbrace{P_X(DIS(V_{n+1}))}_{d(h, h^*) \forall h \in V_{n+1}} \leq \frac{1}{2^n} P_X(DIS(V_1 = \mathcal{H})) \leq \frac{1}{2^n} = \varepsilon$$

set $n = \log \frac{1}{\varepsilon}$

$$\begin{aligned}
& P_X \left(DIS \left(\overbrace{B(h^*, \tau)}^{\{h \in \mathcal{H}: d(h, h^*) \leq \tau\}} \right) \right) \\
& \theta := \sup_{\tau \in (0,1)} \frac{P_X(DIS(B(h^*, \tau)))}{\tau} \text{ (Assume } \theta < \infty \text{)} \leftarrow \text{ (disagreement coefficient)} \\
& \Rightarrow \forall \tau \in (0,1), P_X(DIS(B(h^*, \tau))) \leq \theta \tau \\
& P_X(DIS(V_{i+1})) \leq P_X(DIS(B(h^*, P_X(DIS(V_i))r))) \leq \theta P_X(DIS(V_i))r
\end{aligned}$$



$$\text{Choose } r = \frac{1}{2\theta}, k \geq \frac{\log \frac{n|V_i|}{\delta}}{\frac{1}{2\theta}}$$

$$\text{Set } k = \frac{\log \frac{n|V_i|}{\delta}}{\frac{1}{2\theta}} = 2\theta \left\lceil \log \log \frac{1}{\varepsilon} + \log \frac{|\mathcal{H}|}{\delta} \right\rceil$$

Total number queries by CAL

$$k \cdot n = 2\theta \left\lceil \log \log \frac{1}{\varepsilon} + \log \frac{|\mathcal{H}|}{\delta} \right\rceil \log \frac{1}{\varepsilon}$$

22.2 End Missed Lecture

22.3 Stochastic Bandits

Arms $1 \dots K$

Unknown but fixed reward distributions $V_1 \dots V_k$ with means $U_1 \dots U_k \in \mathbb{R}$

n total rounds

(pseudo) regret

$$U^* = \max_{i \in [k]} U_i$$

Let $I_t \in [k]$ be the arm you pull at round $t, t \in [n]$

Let X_t be the reward you see at round t

$$nU^* - \sum_{t=1}^n \mathbb{E} U_{I_t}$$

exploration then exploitation

1. pull each arm m times, estimate $\hat{U}_i, i \in [k]$
2. for $n - mk$ rounds, pull $I_t := \arg \max \hat{U}_i$

23 Lecture 23

23.1 Subgaussian tail bounds

Let $X_{i=1\dots n} - u$ be independent σ subgaussian random variables. Then,

$$\mathbb{P}(\hat{u} \geq u + \varepsilon) \leq e^{-\frac{n\varepsilon^2}{2\sigma^2}}$$

$$\text{wp} \geq 1 - \delta, u \leq \hat{u} + \sqrt{\frac{2\sigma^2 \log \frac{1}{\delta}}{n}}$$

23.2 K-arm bandit (stochastic)

"environment" k arms with reward

distributions $V_1 \dots V_k$, 1-subgaussian

with mean $U_1 \dots U_k$

def: $U^* = \max_{i \in [k]} U_i$

(assume $U_1 \geq U_2 \geq \dots \geq U_k$)

def: $\Delta_i = U_1 - U_i, i \in [k]$

"learner, agent"

knows: k , 1-subgaussian, n time horizon, $I_t = \text{"policy"}$ $\text{alg}(I_1, X_1, \dots, I_{t-1}, X_{t-1}), X_t \sim V_{I_t}, T_i(t)$, number of arm i pulls up to time t

Goal: minimize (psuedo) regret

$$\text{Reg} := nU^* - \sum_{t=1}^n \mathbb{E}U_{I_t}$$

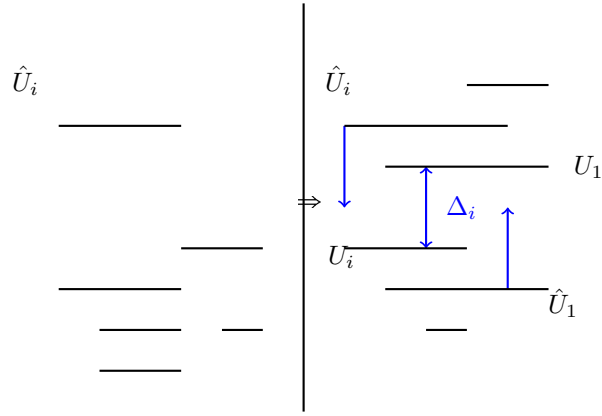
Alg: exploration-then-exploitation (m)

1. pull each arm m times $\Rightarrow \hat{U}_t, i \in [k]$
2. For remaining $n - mk$ pulls, $\arg \max_{i \in [k]} \hat{U}_i$

$$m\Delta_1 + m\Delta_2 + \dots + m\Delta_k = m \sum_{i=1}^k \Delta_i = \frac{n}{k} \sum \Delta_i$$

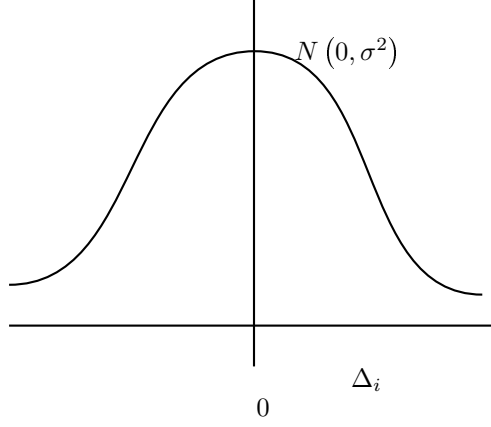
$$\text{Reg} = \sum_{i=1}^k \mathbb{E}(T_i(n)) \Delta_i$$

$$\begin{aligned}
\mathbb{E}T_i(n) &= m + (n - mk) \mathbb{P}(I = i) \\
&= m + (n - mk) \mathbb{P}\left(\hat{U}_i \geq \max_{j \in [k], j \neq i} \hat{U}_j\right) \\
&\leq m + (n - mk) \mathbb{P}\left(\hat{U}_i \geq \hat{U}_1\right) \\
&= m + (n - mk) + \mathbb{P}\left(\hat{U}_i \leq \hat{U}_1 + U_i - U_i + \Delta_i\right) \\
&= m + (n - mk) + \mathbb{P}\left((\hat{U}_i - U_i) - (\hat{U}_1 - U_1) \geq \Delta_i\right) \\
&\leq m + (n - mk) e^{-\frac{\Delta_i^2}{2\left(\frac{2}{m}\right)}}
\end{aligned}$$



$$\begin{aligned}
\hat{U}_i &= \frac{1}{m} \sum_{t=1, I_t=i}^{mk} X_t, X_t \stackrel{iid}{\sim} V_i, X_t - U_i, 1 - \text{subgaussian} \\
&\Rightarrow \hat{U}_i - U_i, \frac{1}{\sqrt{m}} - \text{subgaussian} \forall i \\
&\Rightarrow (\hat{U}_i - U_i) - (\hat{U}_1 - U_1), \sqrt{\frac{2}{m}} - \text{subgaussian}
\end{aligned}$$

$$\mathbb{P}(X \geq \Delta_i) \leq e^{-\frac{\Delta_i^2}{2\sigma^2}}$$



$$\begin{aligned}
Reg &= \sum_{i=1}^k \mathbb{E} T_i(n) \\
&\leq \sum_{i=1}^k \left(m + (n - mk) e^{-\frac{\Delta_i^2 m}{4}} \right) \Delta_i \\
&\leq \sum_i^k \left(m + n e^{-\frac{\Delta_i^2 m}{4}} \right) \Delta_i \\
&\leq \sum_i^k \left(m + n e^{-\frac{\Delta_2^2 m}{4}} \right) \Delta_2
\end{aligned}$$

$$\begin{aligned}
&\frac{\partial \text{RHS}(m)}{\partial m} = 0 \\
&\Rightarrow \sum_{i=1}^k \Delta_i \left(1 - k e^{-\frac{\Delta_i^2 m}{4}} + (n - mk) e^{-\frac{\Delta_i^2 m}{4}} \left(-\frac{\Delta_i^2}{4} \right) \right) \\
&= \sum_i \Delta_i \left(1 - \left[k + (n - mk) \frac{\Delta_i^2}{4} \right] e^{-\frac{\Delta_i^2 m}{4}} \right) \\
&\approx \sum_i \Delta_i \left(1 - \left[k + n \frac{\Delta_i^2}{4} \right] e^{-\frac{\Delta_i^2 m}{4}} \right) \\
&\dots
\end{aligned}$$

$$\frac{\partial \text{RHS}(m)}{\partial m} = 0$$

$$\begin{aligned}
&\Rightarrow \sum_i \Delta_i \left(1 - \frac{n\Delta_2^2}{4} e^{-\frac{\Delta_2^2 m}{4}} \right) = 0 \\
&\Rightarrow \sum \Delta_i = \left(\sum \frac{n\Delta_2^2 \Delta_i}{4} \right) e^{-\frac{\Delta_2^2 m}{4}} \\
&\Rightarrow \frac{4}{\Delta_2^2} \log \frac{\sum \frac{n\Delta_2^2 \Delta_i}{4}}{\sum \Delta_i} = m = \frac{4}{\Delta_2^2} \log \frac{n\Delta_2^2}{4}
\end{aligned}$$

24 Lecture 24

$$U_1 \geq \dots \geq U_k$$

$$\Delta_i = U_1 - U_i, i \in [k]$$

$$Reg := nU_1 - \sum_{t=1}^n \mathbb{E}[U_{I_t}]$$

$$\begin{aligned}
&= \sum_{i=1}^k \Delta_i \mathbb{E}[T_i(n)] \\
&\leq \left(m + n e^{-\frac{\Delta_2^2 m}{4}} \right) \left(\sum \Delta_i \right) \leftarrow \boxed{1}
\end{aligned}$$

$$\text{optimal } m = \frac{4}{\Delta_2^2} \log \frac{n\Delta_2^2}{4}$$

$$\begin{aligned}
\boxed{1} &\stackrel{m}{=} \left(\frac{4}{\Delta_2^2} \log \frac{n\Delta_2^2}{4} + n [e^m]^{-\frac{\Delta_2^2}{4}} \right) \sum \Delta_i \\
&= \left(\frac{4}{\Delta_2^2} \log \frac{n\Delta_2^2}{4} + n \left[\frac{n\Delta_2^2}{4} \right]^{\frac{4}{\Delta_2^2} \left(-\frac{\Delta_2^2}{4} \right)} \right) \sum \Delta_i \\
&= \left(\frac{4}{\Delta_2^2} \log \frac{n\Delta_2^2}{4} + \frac{4}{\Delta_2^2} \right) \sum \Delta_i \\
&= \frac{4}{\Delta_2^2} \left(1 + \log \frac{n\Delta_2^2}{4} \right) \left(\sum \Delta_i \right) \\
&\stackrel{k=2, \underline{\Delta} := \Delta_2}{=} \frac{4}{\Delta} \frac{1 + \log(n\Delta^2)}{4} \leftarrow \boxed{2}
\end{aligned}$$

Alg needs Δ, k, n

Worst case Δ

$$\frac{\partial \boxed{2}}{\partial \Delta} = 0$$

$$\begin{aligned}
-\frac{4}{\Delta^2} \left(1 + \log \frac{n\Delta^2}{4} \right) + \frac{4}{\Delta} \frac{4}{n\Delta^2} \frac{2n\Delta}{4} &= 0 \\
-\frac{4}{\Delta^2} - \frac{4}{\Delta^2} \log \frac{n\Delta^2}{4} + \frac{8}{\Delta^2} &= 0 \\
\frac{4}{\Delta^2} &= \frac{4}{\Delta^2} \log \frac{n\Delta^2}{4} \\
\frac{n\Delta^2}{4} &= e \\
\Delta^* &= \sqrt{\frac{4e}{n}} \\
\boxed{2}^{\Delta^*} &\triangleq \frac{4\sqrt{n}}{\sqrt{4e}} \left(1 + \log \frac{n\frac{4e}{4}}{4} \right) \\
&= \frac{8\sqrt{n}}{\sqrt{4e}} = \frac{4}{\sqrt{e}} \sqrt{n}
\end{aligned}$$

24.1 Upper Confidence Bound (UCB)

$$x_1 \dots x_t \sim V(1 - \text{subgaussian})$$

$$U := \frac{1}{t} \sum_{\tau=1}^t x_\tau \leq \hat{U}_t + \sqrt{\frac{2 \log \frac{1}{\delta}}{t}} \text{ wp } 1 - \delta$$

$$\text{UCB}(\text{arm } i, \text{count } t \text{ of arm } i \text{ pulls}) :=$$

$$\begin{cases} \hat{U}_{it} + \sqrt{\frac{2 \log \frac{1}{\delta}}{t}} & \text{if } t > 0 \\ \infty & \text{if } t = 0 \end{cases}$$

Alg:

for $j = 1 \dots n$

pull $I_j \in \arg \max_{i \in [t]} \text{UCB}(i, t_i)$

$$Reg := \sum_{i=1}^k \Delta_i \mathbb{E} T_i(n)$$

idea:

$$[\forall t : UCB(1, t) > U_1] := \text{Event 1}$$

Fix $i \neq 1$, fix a magic number $\tau_i \in [n]$

$$[UCB(i, \tau_i) < U_1] := \text{Event 2}$$

Define event $E_i = \text{Event 1} \wedge \text{Event 2}$

$$\begin{aligned} \mathbb{E} T_i(n) &= \mathbb{E} \mathbb{1}_{(E_i)} T_i(n) + \mathbb{E} \mathbb{1}_{(E_i^c)} T_i(n) \\ &\stackrel{\text{Claim 1}}{\leq} \tau_i + \mathbb{E} \mathbb{1}_{(E_i^c)} T_i(n) \\ &\stackrel{T_i(n) \leq n}{\leq} \tau_i + n \mathbb{P}(E_i^c) \end{aligned}$$

Claim 1, $E_i \Rightarrow T_i(n) \leq \tau_i$

$$\begin{aligned} \mathbb{P}(E_i^c) &\stackrel{\text{Union}}{\leq} \mathbb{P}(\exists t, UCB(1, t) \leq U_1) + \mathbb{P}(UCB(i, \tau_i) > U_1) \\ &\leq n \mathbb{P}\left(\hat{U}_{1t} + \sqrt{\frac{2 \log \frac{1}{\delta}}{t}} \leq U_1\right) + \mathbb{P}(UCB(i, \tau_i) > U_1) \\ &\stackrel{\text{subgaussian}}{\leq} n\delta + \mathbb{P}\left(\hat{U}_{i, \tau_i} + \sqrt{\frac{2 \log \frac{1}{\delta}}{\tau_i}} > U_i + \Delta_i\right) \\ &= n\delta + \mathbb{P}\left(U_{i\tau_i} - U_i > \underbrace{\Delta_i - \sqrt{\frac{2 \log \frac{1}{\delta}}{\tau_i}}}_{\frac{1}{2}\Delta_i}\right) \leftarrow \boxed{3} \end{aligned}$$

Choose τ_i such that

$$\begin{aligned} \Delta_i - \sqrt{\frac{2 \log \frac{1}{\delta}}{\tau_i}} &= \frac{1}{2}\Delta_i \leftarrow \boxed{4} \\ \boxed{3} &\leq n\delta + e^{-\frac{\tau_i \left(\frac{1}{2}\Delta_i\right)^2}{2}} \leftarrow \boxed{5} \end{aligned}$$

$$\begin{aligned}
\boxed{4} &\Rightarrow \frac{\Delta_i}{2} = \sqrt{\frac{2 \log \frac{1}{\delta}}{\tau_i}} \\
\tau_i &= \frac{4 \cdot 2 \log \frac{1}{\delta}}{\Delta_i^2} \\
\boxed{5} &= n\delta + e^{-\frac{8 \log \frac{1}{\delta}}{\Delta_i^2} \frac{\Delta_i^2}{2 \cdot 4}} \\
&= (n+1)\delta \\
\mathbb{E}T_i(n) &\leq \frac{8 \log \frac{1}{\delta}}{\Delta_i^2} + n(n+1)\delta
\end{aligned}$$

25 Lecture 25

One version of UCB(δ)

$$UCB(i, t_i) := \hat{U}_{i, t_i} + \sqrt{\frac{2 \log \frac{1}{\delta}}{t_i}}$$

repeat: pull $\arg \max_{i \in [k]} UCB(i, t_i)$

Last time

$$\begin{aligned}
\mathbb{E}T_i(n) &\leq \tau_i n(n+1)\delta \\
\tau_i &= \left\lceil \frac{8 \log \frac{1}{\delta}}{\Delta_i^2} \right\rceil \\
Reg &= \sum_{i=1}^k \Delta_i \mathbb{E}T_i(n) \leq \sum_{i=1}^k \Delta_i (\tau_i + n(n+1)\delta) \\
&= \sum_i \Delta_i \left(\left\lceil \frac{8 \log \frac{1}{\delta}}{\Delta_i^2} \right\rceil + n(n+1)\delta \right), \delta := \frac{1}{n^2} \\
&= \sum_i \Delta_i \left(\left\lceil \frac{16 \log n}{\Delta_i^2} \right\rceil + 1 + \frac{1}{n} \right) \\
&\leq \sum_i \Delta_i \left(\frac{16 \log n}{\Delta_i^2} + 3 \right) \\
&= \left(\sum_i \frac{16 \log n}{\Delta_i} \right) + 3 \sum_i \Delta_i
\end{aligned}$$

25.1 Contextual bandit

context vector $f_t := \begin{bmatrix} f_{\text{user}}(t) \\ f_{\text{arm}}(i_t) \end{bmatrix}$

$$\exists w^*, \mathbb{E}[r_{\text{reward}}(t)] = f_t^T w^*$$

25.2 Optimal Teaching

Learner is $\text{ERM} = V(S)$

\mathcal{H} realizable, S , ℓ is 0 – 1 loss

$$\arg \min_{h \in \mathcal{H}} \frac{1}{|S|} \sum_{(x,y) \in S} \ell(h, x, y) = \text{Version Space } V(S)$$

learner that only takes S

Learner is a function $\{S\} \rightarrow 2^{\mathcal{H}}$

\mathcal{H} finite

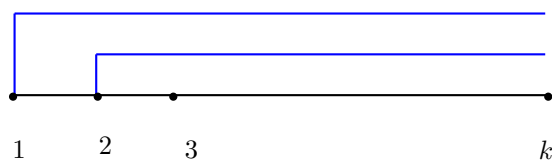
learner $(S = (x^1, y^1) \dots (x^n, y^n))$ creates a lookup table

$$T = \begin{bmatrix} h_1 & x_1 \\ \dots & \dots \\ h_{|S|} & x_{|S|} \end{bmatrix}$$

returns $T(x')$

$$S = (x^*, y^*) \dots$$

$$X = [k]$$



$$L(S) = \{h^*\}$$

$$S = \{(2, -), (3, +)\}$$

$$\min_{S \in \mathbb{S}} |S| \text{ such that } L(S) = \{h^*\}$$

$$\mathbb{S} = \bigcup_{n=0}^{\infty} (X \times Y)^n$$

Teaching dimension (h^*, \mathcal{H})

—	x_1	x_2	...	x_k
h_1	1	0	...	0
...	0	1	...	0
...
h_k	0	...	0	1
h_{k+1}	0	0

$$TD(\mathcal{H}) := \max_{h \in \mathcal{H}} TD(h, \mathcal{H})$$

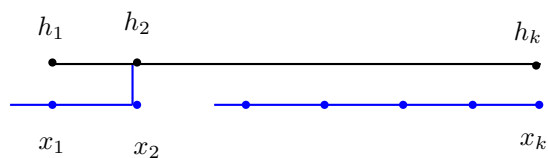
—	x_1	x_2	...	x_k
h_1	1	0	...	0
...	0	1	...	0
...
h_k	0	...	0	1
h_{k+1}	1	1

—	x_1	x_2	...	x_k
h_1	1	0	...	0
...	0	1	...	0
...
h_k	0	...	0	1
h_{k+1}	1	1
h_{k+2}	0	0

26 Lecture 26

learner: $L(S) \in 2^{\mathcal{H}}$

Teacher: $T(h) = S$ such that $L(T(h)) = \{h\}$



$$T(h_2) = \{(x_1, -), (x_2, +), (x_3, +)\}$$

Teaching Dimension (h, \mathcal{H}, L_{VS})

$$= \min_{S \in 2^X} |S| \text{ such that } L(S) = \{h\}$$

$$TD(\mathcal{H}) = \max_{h \in \mathcal{H}} TD(h, \mathcal{H})$$

—	x_1	...	x_d	x_{d+1}	...	x_{d+2^d}
h_1	—	—	—	1	0	0
...	—	Whole truth Table	—	0	1	0
h_{2^d}	—	—	—	0	0	1

$$VC(\mathcal{H}) = d$$

$$TD(\mathcal{H}) = 1$$

—	x_1	...	x_d
h_1	1	0	0
...	0	1	0
h_d	0	0	1
h_{d+1}	0	0	0

$$VC(\mathcal{H}) = 1$$

$$TD(\mathcal{H}) = d$$

”Collusion”

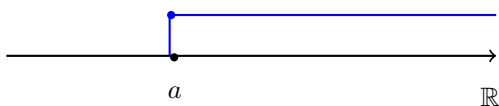
ONE definition of collusion-free teaching

If S teaches h , Then $\forall S' \supset S$ (labeled consistently by h) also teaches h

—	x_1	x_d
h_1	1	1	0	0
h_2	0	1	1	0	...	0
...
h_{d-1}	0	0	1	1
h_d	1	0	0	1

no-clash teaching

$T(h)$, if $\forall h, h', T(h)$ is INconsistent with h' OR $T(h')$ is INconsistent with h



$$\mathcal{H} = \{h_a(x) = \mathbb{1}_{(x \geq a)} : a \in \mathbb{R}\}$$

Teaching Dimension (h, \mathcal{H}, L_{VS})

$$= \min_{S \in 2^X} |S| \text{ such that } L(S) = \{(h - \varepsilon, h + \varepsilon)\}$$

Teaching Dimension $(h, \mathcal{H}, L_{SVM} \text{ (hard margin)})$

$$= \min_{S \in 2^X} |S| \text{ such that } L(S) = \{h\}$$

Teaching Dimension $(h, \mathcal{H}, L_{SVM} \text{ (soft margin, logistic regression)})$

$$= \min_{S \in 2^X} |S| \text{ such that } L(S) = \{h\}$$

26.1 Reinforcement Learning

$$Reg \sim O\left(\sqrt{n \log n}\right)$$