

Thory HW2 due May 21st.

Lecture 15.

Recap. [Fundamental theorem of statistical learning]

\mathcal{H} has a finite VC-dim $\implies \mathcal{H}$ has the Uniform Convergence Property.

\iff If $VCdim(\mathcal{H}) = d$, then $m_{\mathcal{H}}^{u.c.}(\epsilon, \delta) \leq C_2 \cdot \frac{d + \log \frac{1}{\delta}}{\epsilon^2}$

we prove a weaker version

$m_{\mathcal{H}}^{u.c.}(\epsilon, \delta) \leq C \cdot \frac{d \cdot \log \frac{1}{\delta}}{\delta^2 \epsilon^2}$

① Lemma [Sauer-Shelah-Perles] *proof in class*

$T_{\mathcal{H}}(m) \triangleq \max_{C \subseteq \mathcal{X}, |C|=m} |\mathcal{H}_C| \leq \sum_{i=1}^{d=VCdim(\mathcal{H})} \binom{m}{i} \leq d m^d$

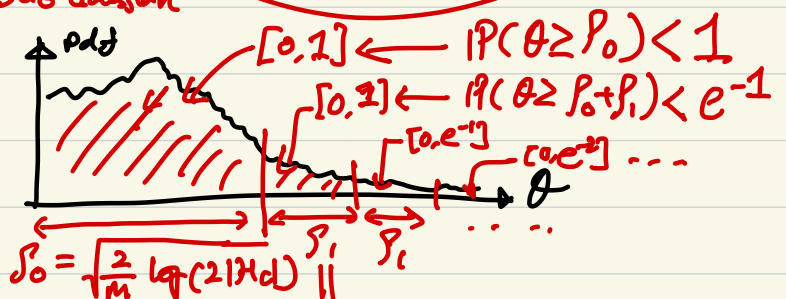
② Theorem [Uniform Convergence] *proof in book*

$|\mathcal{L}_D(h) - \mathcal{L}_S(h)| \leq \frac{4 + \sqrt{\log(T_{\mathcal{H}}(2m))}}{\delta \cdot \sqrt{m/2}}$

* One remaining step. [Lemma A.4 in book]

$P(\overset{R.V.}{\frac{1}{\theta}} > \beta) \leq 2 \cdot |\mathcal{H}_C| \cdot e^{-\frac{\beta^2}{2}}$ $\implies E[\theta] \leq \frac{4 + \sqrt{\log(2|\mathcal{H}_C|)}}{\sqrt{m/2}}$

Proof sketch \rightarrow sub-Gaussian pdf



$E[\theta] = \int_0^{\infty} \theta \cdot f(\theta) d\theta = p_0 + \sum_{i=1}^{\infty} e^{-(i-1)} \cdot (p_0 + i \cdot p_1) \leq \dots$

$4 \cdot \sqrt{\frac{2}{m}} + \underbrace{\sqrt{\frac{2}{m} \log(2|\mathcal{H}_d)}}_{p_0}$

Chapter 9. Linear Predictors.

• Binary classification with linear classifier. $y \in \{\pm 1\}$

• Class of affine functions $\mathcal{H}_d \triangleq \{ h_{w,b}(x) = \langle w, x \rangle + b : w \in \mathbb{R}^d, b \in \mathbb{R} \}$

↑ *not linear*

↑ $w^T x = \text{inner product}$

↓ *bias*

• Prediction at x : $\hat{y} = \text{sign}(\langle w, x \rangle + b)$

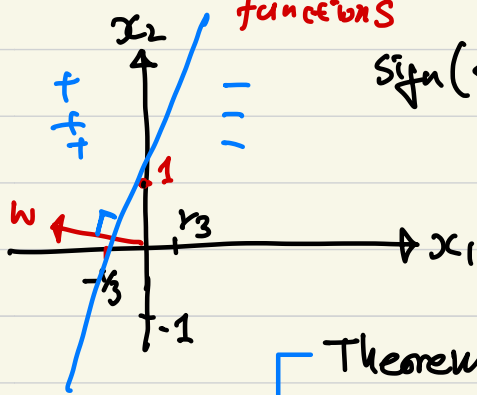
• 0-1 loss: $L_{w,b}(x, y) = \mathbb{I}(\text{sign}(\langle w, x \rangle + b) \neq y)$

• Learning half spaces:

$HS_d = \text{sign} \circ \mathcal{H}_d \triangleq \{ \text{sign}(h_{w,b}(x)) = \text{sign}(\langle w, x \rangle + b) : w \in \mathbb{R}^d, b \in \mathbb{R} \}$

↑ *class of functions*

↑ *composition of a function $\text{sign}(\cdot)$ and a class.*



$$\text{sign}(\langle w, x \rangle + b)$$

$$w = (-3, 1), b = 1$$

$$\text{decision boundary is } [-3, 1] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + 1 = 0$$

Theorem. For Homogeneous class, i.e. $b=0$,

$$V(\dim(HS_d)) = d$$

Proof \Rightarrow ① $\exists C$ with $|C|=d$ s.t. HS_d shatters C .

$C = \{e_1, e_2, \dots, e_d\}$ is shattered by HS_d .

$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_d \end{bmatrix}$ For any (y_1, y_2, \dots, y_d) , we need to find w s.t. $\text{sign}(\langle e_i, w \rangle) = y_i, \forall i \in [d]$.

Let, $w = (y_1, y_2, \dots, y_d)$, then \uparrow

* more generally there is d -degrees of freedom...

(2) $\forall C$ with $|C| = d+1$, cannot be shattered by HS_d

Consider a set $C = \{x_1, x_2, \dots, x_d, x_{d+1}\}$ and labels $y_1, y_2, \dots, y_d, y_{d+1}$ to be assigned.

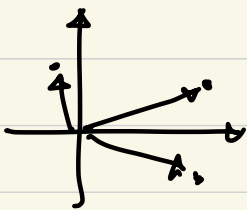
proof by contradiction: suppose C can be shattered by HS_d (any y 's can be assigned).

$$\exists w \in \mathbb{R}^d \text{ s.t. } \begin{aligned} \text{sign}(\langle w, x_1 \rangle) &= y_1 \\ \text{sign}(\langle w, x_2 \rangle) &= y_2 \\ &\vdots \end{aligned} \quad \text{for any } (y_1, \dots, y_{d+1})$$

* Construction of counter example.

for given $C = \{x_1, \dots, x_{d+1}\}$, exists $a = (a_1, a_2, \dots, a_{d+1})$ s.t.

$$a_1 x_1 + a_2 x_2 + \dots + a_d x_d + a_{d+1} x_{d+1} = 0$$



because $d+1$ vectors in \mathbb{R}^d , are linearly dependent.
from Linear Algebra

we let $(y_1 = \text{sign}(a_1), y_2 = \text{sign}(a_2), \dots) = y$

let $I \triangleq \{i \in [d+1] : a_i > 0\}$ at least one of them is non-empty set
 $J \triangleq \{j \in [d+1] : a_j < 0\}$

Note that $\sum_{i \in I} a_i x_i + \sum_{j \in J} a_j x_j = 0$ by construction

$$\rightarrow \sum_{i \in I} a_i x_i = \sum_{j \in J} |a_j| x_j$$

$$\text{sign}(a_i) = \text{sign}(\langle x_i, w \rangle) \quad \rightarrow \quad \sum_{i \in I} a_i \langle x_i, w \rangle = \langle \sum_{i \in I} a_i x_i, w \rangle = \langle \sum_{j \in J} |a_j| x_j, w \rangle = \sum_{j \in J} |a_j| \langle x_j, w \rangle$$

Theorem 9.3 $\forall \text{dim}(HS_d)$ with parameter $w \in \mathbb{R}^d$, \exists ERM is $d+1$

Proof \Rightarrow (1) $\exists C$ with $|C|=d+1$ s.t. C is shattered by HS_d
 $C = \{0, e_1, e_2, \dots, e_d\}$
 label y_0, y_1, \dots, y_d

(2) $\forall C$ with $|C|=d+2$, C is not shattered by HS_d .
 same proof with $(w, b) \in \mathbb{R}^{d+1}$ & $d+2$ vectors $C = \{x_1, \dots, x_{d+2}\}$

by Fundamental Theorem of Statistical Learning, we know that

ERM achieves

$$\underset{\uparrow \text{ERM}}{L_D(h_S)} \leq \epsilon$$

[Thm 6.8]

w.p $1-\delta$.

$$\text{if } m \geq C \cdot \frac{d+1}{\epsilon}$$

, when \mathcal{D} is realizable.

*hard to find ERM

when not realizable.

Q. How do you find ERM solution halfspaces? (in the realizable case)
 train data is "separable"

$$\min_{h \in HS_d}$$

$$L_S(h)$$

\uparrow
0-1 loss

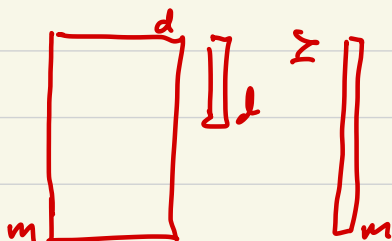
generally hard

easy when realizable.

Algorithm 1. Linear Program. $\hat{=}$ optimization with linear objective
 in Convex Program & linear constraints

Generic LP: Maximize $\langle \underline{u}, w \rangle$
 $w \in \mathbb{R}^d$

Subject to $\underline{A} \cdot w \geq \underline{v}$: entrywise inequality



ERM.

find $w \in \mathbb{R}^d$

Subject to $\text{sign}(\langle w, x_i \rangle) = y_i$ $(x_i) \in S$

\iff

$$y_i \cdot \langle w, x_i \rangle > 0$$

: under realizability solution exists.

\iff

$$y_i \cdot \langle w, x_i \rangle \geq 1.$$

: we can always rescale w .

minimize 0

s.t. $y_i \cdot \langle w, x_i \rangle \geq 1, \forall i \in [m]$

\implies LP-solver.

\downarrow
ERM solution.

Algorithm 2. finds ERM for realizable linear classification.

Iterative Algorithm:

input: $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

initialize: $w^{(1)} = (0, \dots, 0) \in \mathbb{R}^d$

for $t=1, 2, \dots$

if $\exists i \in [m]$ s.t. $y_i \cdot \langle w^{(t)}, x_i \rangle \leq 0$ then

$$w^{(t+1)} \leftarrow w^{(t)} + y_i \cdot x_i$$

else

output $w^{(t)}$

[Perceptron Algorithm]

we want $y_i \cdot \langle x_i, w \rangle > 0$

$$y_i \cdot \langle w^{(t+1)}, x_i \rangle$$

$$= y_i \cdot \langle w^{(t)} + y_i \cdot x_i, x_i \rangle$$

$$= y_i \cdot \langle w^{(t)}, x_i \rangle + y_i^2 \cdot \langle x_i, x_i \rangle$$

$$= y_i \cdot \langle w^{(t)}, x_i \rangle + 1 \cdot \|x_i\|^2$$

$$> y_i \cdot \langle w^{(t)}, x_i \rangle$$

Theorem 9.1.

Assume S is separable, $B = \min \{ \|w\| : \forall i \in [n], y_i \langle w, x_i \rangle \geq 1 \}$
 $R = \max_i \|x_i\|.$

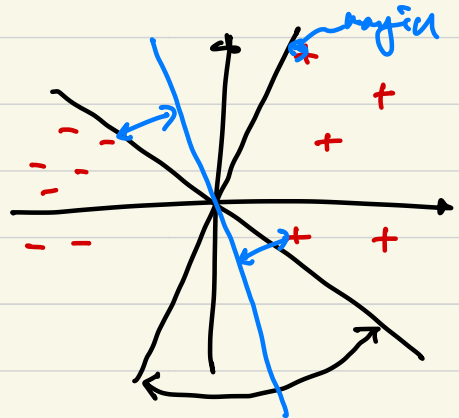
then Perceptron Algorithm stops after at most $(RB)^2$ iterations.

↑ depends on B
is suboptimal!

When it stops, it finds an ERM with $y_i \langle w^, x_i \rangle > 0, \forall i \in [n]$

Proof > let $w^* \in \arg \min_w \|w\|$ s.t. $y_i \langle w, x_i \rangle \geq 1 \forall i \in [n].$

↳ Max-margin separator



claim: $\frac{\langle w^*, w^{(T+1)} \rangle}{\|w^*\| \cdot \|w^{(T+1)}\|} \geq \frac{\sqrt{T}}{R \cdot B}$

because L.H.S = $\cos(\angle(w^*, w^{(T+1)})) \leq 1$

↳ $T \leq R^2 B^2$

proof of claim > $\|w^*\| \leq B, \langle w^*, w^{(T+1)} \rangle \geq T, \|w^{(T+1)}\| \leq R\sqrt{T}$
(*) (**)

(*) $\langle w^*, w^{(T+1)} \rangle - \langle w^*, w^{(1)} \rangle = \langle w^*, w^{(T+1)} - w^{(1)} \rangle$

$= \langle w^*, y_i x_i \rangle$

← by Algorithm

$= y_i \langle w^*, x_i \rangle \geq 1$

$= 1 \geq 1$ by def

$\Rightarrow \langle w^*, w^{(T+1)} \rangle = \sum_{i=1}^T \{ \langle w^*, w^{(i+1)} \rangle - \langle w^*, w^{(i)} \rangle \} \geq T$

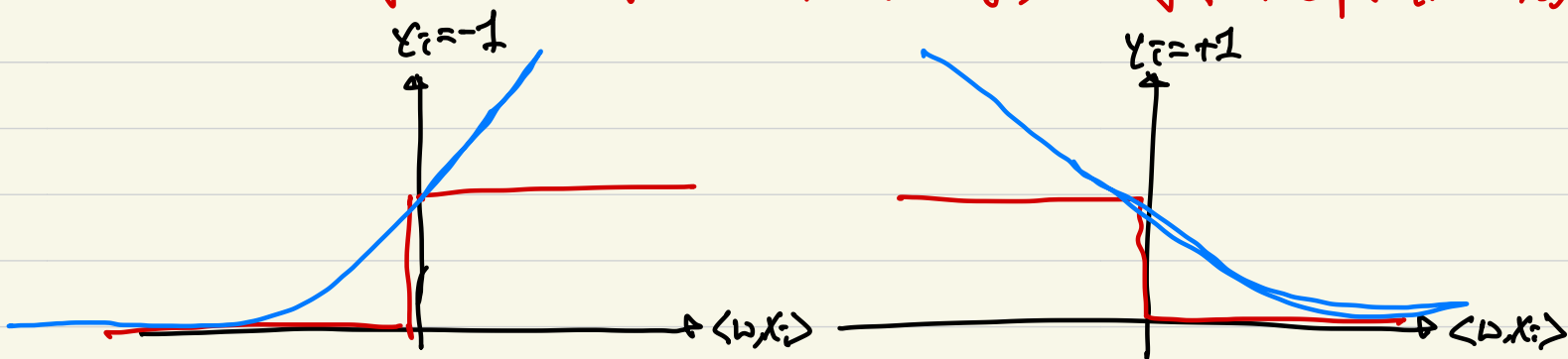
$$\begin{aligned}
 (**) \quad \|w^{(t+1)}\|^2 &= \|w^{(t)} + \gamma_i x_i\|^2 \\
 &= \|w^{(t)}\|^2 + \underbrace{\gamma_i^2}_{=1} \cdot \underbrace{\|x_i\|^2}_{\leq R^2} + \underbrace{2\gamma_i \langle w^{(t)}, x_i \rangle}_{\leq 0 \text{ by algorithm}} \\
 &= \|w^{(t)}\|^2 + R^2
 \end{aligned}$$

$$\|w^{(T+1)}\|^2 \leq T \cdot R^2$$

Algorithm 3. Logistic Regression : surrogate loss

True loss: 0-1 loss, $l_{h_{w, \sigma}}(x_i, y_i) = \mathbb{I}(y_i \neq \text{sign}(\langle w, x_i \rangle))$

Surrogate loss: Logistic loss, $Q_w(x, y) = \log(1 + \exp(-\gamma_i \langle w, x_i \rangle))$



* we choose Convex, Smooth surrogate loss

$$\text{minimize } \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-\gamma_i \langle w, x_i \rangle))$$

Algorithm 4. online Perceptron Algorithm [much more on online learning at CSE541 Interactive meeting layout.]

Data Come in an online fashion

time: 1, 2, ..., t, t+1

- receive x_t
- choose h_t , Predict $h_t(x_t)$
- receive y_t
- loss: $l_{h_t}(x_t, y_t)$

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

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

- receive x_t
 - Predict $P_t = \text{sign}(\langle w^{(t)}, x_t \rangle)$
 - If $y_t \langle w^{(t)}, x_t \rangle \leq 0$,
- $$w^{(t+1)} \leftarrow w^{(t)} + \gamma_t \cdot x_t$$

else

$$w^{(t+1)} \leftarrow w^{(t)}$$

In an online problem, we care about how many mistakes we make over time.

Theorem 21.16 $R = \max_t \|X_t\|$, if realizable, with $\exists w^*$
 $\forall_i \langle w^*, x_i \rangle \geq 1, \forall t.$

$$\# \text{ of mistakes} \leq R^2 \|w^*\|^2$$

proof \rightarrow [HW3 Problem 5].