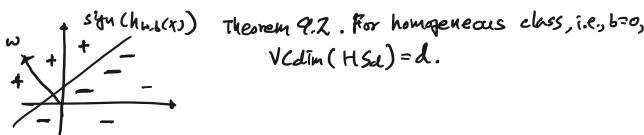


Chapter 9. Linear Predictors.

- class of affine functions $H_d = \{ h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R} \}$
- Prediction at x : $\hat{y} = \text{sgn}(\langle w, x \rangle + b)$
- D1 loss: $l_{\text{D1}}(x, y) = \mathbb{I}(\text{sgn}(\langle w, x \rangle + b) \neq y)$
- Learning half spaces = linear classifier for binary classification $\mathcal{Y} = \{-1, +1\}$

$$HS_d = \text{sgn} \circ H_d = \{ \text{sgn}(h_{w,b}(x)) = \text{sgn}(\langle w, x \rangle + b) : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$

↓
composition of a function and a class.



Proof 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 , because
any (y_1, y_2, \dots, y_d) , $w = (y_1, y_2, \dots, y_d)$ satisfies $\text{sgn}(\langle w, e_i \rangle) = y_i$

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

Label: $C = \{x_1, x_2, \dots, x_d, x_{d+1}\}$, suppose C can be shattered by HS_d
s.t. $\text{sgn}(\langle w, x_1 \rangle) = y_1, \dots, \text{sgn}(\langle w, x_{d+1} \rangle) = y_{d+1}$

But $a_1 \cdot x_1 + a_2 \cdot x_2 + \dots + a_{d+1} \cdot x_{d+1} = 0$. \leftarrow this contradicts because d+1 vectors in \mathbb{R}^d .

$I \triangleq \{i : a_i > 0\}, J \triangleq \{j : a_j < 0\}$, I & J both cannot be empty.
Suppose both nonempty.

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

$$0 < \sum_{i \in I} a_i \cdot \langle x_i, w \rangle = \left\langle \sum_{i \in I} a_i x_i, w \right\rangle = \left\langle \sum_{j \in J} |a_j| x_j, w \right\rangle = \sum_{j \in J} |a_j| \langle x_j, w \rangle < 0$$

$$\uparrow$$

• Hence, contradiction.

• If J empty some contradicts with $\sum_{j \in J} |a_j| \langle x_j, w \rangle \geq 0$.

Theorem 9.3. $\text{VCdim}(HS_d)$ with parameters w, b in \mathbb{R}^d is $d+1$.

① $\exists C$ with $|C|=d+1$ shattered by HS_d .

$C = \{e_1, e_2, \dots, e_d\}$
Label $y_1, y_2, \dots, y_d \rightarrow$ we let $b=y_0$, $w=y_1, \dots, w_d=y_d$, then
 $\forall x \in C, \langle w, x \rangle + b = y_0$

② Same proof in \mathbb{R}^{d+1} and \mathbb{R}^2 vectors.

\Rightarrow ERM achieves $L_D(HS) \leq \epsilon$ comp. \leftarrow
GRM

if $m \geq C \cdot \frac{d+1+\epsilon}{\epsilon}$, when D is realizable * Hard if
Not
Realizable

* From Geometric view of Fundamental Theorem of Statistical Learning [Thm 8-8]

Q. How do you find ERM solution half-spaces? (in the nonseparable case) $\rightarrow \boxed{\min_w L_S(w)}$

① Linear Program: optimization with linear objective and linear constraints.

↑
Convex
optimization

$$\begin{aligned} \text{Max}_{w \in \mathbb{R}^d} \quad & \langle w, u \rangle \\ \text{Subject to} \quad & \frac{1}{d} \sum_{i=1}^d w_i y_i \geq 1 \leftarrow \text{entrywise inequality} \\ & \|w\|_1 \leq 1 \end{aligned}$$

$$\min_w L_S(w)$$

• generally hard
• easy when realizable.

ERM \leftarrow L.P. solvers.

find w s.t. $\text{sign}(w \cdot x_i) = y_i$

$\Leftrightarrow y_i \cdot \langle w, x_i \rangle > 0$. (such w exists under realisability)

claim $\Leftrightarrow y_i \cdot \langle \tilde{w}, x_i \rangle \geq 1$, because we can use $\tilde{w} = C \cdot w$ instead with arbitrary C .

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

$$y_i \cdot x_i^T \xrightarrow{m} \boxed{\geq 1}$$

② Another algorithm that finds an ERM for realizable linear classification.

Iterative algorithm:

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

initialise: $w^{(0)} = (0, \dots, 0)$

for $t=1, 2, \dots$

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

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

else

output $w^{(t)}$.

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

$$\Leftrightarrow y_i \cdot \langle w^{(t)}, x_i \rangle > 0$$

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

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

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

Theorem 9.1

Assume S is separable, $B = \max_{i \in [m]} \|x_i\| : \forall i \in [m] y_i \cdot \langle w, x_i \rangle \geq 1$

$$R = \max_i \|x_i\|. \quad \boxed{\text{depends on } B \text{ is suboptimal}}$$

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

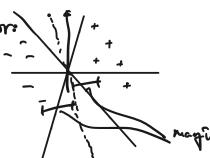
• When it stops it finds an ERM with $y_i \cdot \langle w^{(t)}, x_i \rangle > 0, \forall i \in [m]$

• Perceptron is PAC learner with $M_H(\epsilon, \delta) \geq \frac{\log \frac{1}{\delta}}{\epsilon^2}$

Proof:

Let $w^* \in \arg \min_w \|w\|$ s.t. $y_i \cdot \langle w, x_i \rangle \geq 1, \forall i \in [m]$

\hookrightarrow max margin separator



$$\text{Claim: } \frac{\langle w^*, w^{(T+1)} \rangle}{\|w^*\| \cdot \|w^{(T+1)}\|} \geq \frac{\sqrt{T}}{RB} \Rightarrow T \leq (RB)^2$$

$$\uparrow \|w^*\| = B, \langle w^*, w^{(T+1)} \rangle \geq T, \|w^{(T+1)}\| \leq R\sqrt{T}$$

$$\begin{aligned} (*) \quad \langle w^*, w^{(T+1)} \rangle - \langle w^*, w^{(T)} \rangle &= \langle w^* - w^{(T+1)} - w^{(T)}, w^{(T)} \rangle \\ &= \langle w^*, y_i x_i \rangle \\ &= y_i \cdot \underbrace{\langle w^*, x_i \rangle}_{\geq 1} \geq 1. \end{aligned}$$

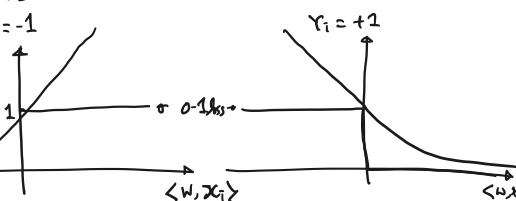
$$\Rightarrow \langle w^*, w^{(T+1)} \rangle = \sum_{i=1}^T \langle w^*, w^{(T+1)} - w^{(T)} \rangle \leq T$$

$$\begin{aligned} (***) \quad \|w^{(T+1)}\|^2 &= \|w^{(T)} + y_i x_i\|^2 \\ &= \|w^{(T)}\|^2 + y_i^2 \|x_i\|^2 + 2y_i \langle w^{(T)}, x_i \rangle \\ &\stackrel{1 \leq R}{\leq} \|w^{(T)}\|^2 + R^2 \\ &\leq \|w^{(T)}\|^2 + R^2 \\ \|w^{(T+1)}\|^2 &\leq T \cdot R^2 \end{aligned}$$

③ Logistic Regression: surrogate loss

true loss: 0-1 loss; $l_w(x) = \mathbb{I}(\text{sign}(\langle w, x \rangle) \neq y)$

• non-convex
• non-smooth
• no gradient



Example: $f = \lambda_1 \cdot \mathbb{I}(\langle w, x \rangle \geq 0) + \lambda_2 \cdot \text{logistic}(\langle w, x \rangle)$

• convex ERM: minimize $\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \langle w, x_i \rangle))$

- smooth

* Online Perceptron Algorithm [much more on online learning at

CSE 541. Interactive Machine Learning]

↓
data come in an online fashion.

1, 2, ..., t, t+1, ...

receive x_t

choose y_t

predict $p_t = \text{sign}(\langle w^{(t)}, x_t \rangle)$

receive y_t

pay loss: $l_h(x_t, y_t)$

initialize $w^{(0)} = (0, 0, \dots, 0)$

for $t=1, \dots, T$

 receive x_t

 predict $p_t = \text{sign}(\langle w^{(t)}, x_t \rangle)$

 if $y_t \langle w^{(t)}, x_t \rangle \leq 0$

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

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

we care about, in an online learning problem, how many mistakes you make.

Theorem 21.16. $R = \max_t \|x_t\|$, if realizable, i.e. $\exists w^* \text{ s.t. } \langle w^*, x_t \rangle \geq 1 \forall t$,

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

proof [HW1 Prob 5.]