

So, let's move on to some hyp. classes.

Linear predictors

Define $\mathcal{L}_d = \{h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R}\}$ } Class of affine functions.

$$h_{w,b}(x) = \langle w, x \rangle + b$$

From \mathcal{L}_d , we can get different predictors via composition with some $\phi: \mathbb{R} \rightarrow \mathbb{Y}$. Further, we can use

$w' = (b, w_1, \dots, w_d)$ and $x' = (1, x_1, \dots, x_d)$ to write

$$h_{w,b}(x) = \langle w', x' \rangle$$

\Rightarrow Affine functions in \mathbb{R}^d can be written as homogeneous linear functions in \mathbb{R}^{d+1} .

$$f(sx_1, \dots, sx_d) = s^k \cdot f(x_1, \dots, x_d)$$

Here: $k=2$ as $\langle kx, k\omega \rangle = \sum_{i=1}^d kx_i \cdot k\omega_i$
 $= k^2 \cdot \langle x, \omega \rangle$

if $\phi: \mathbb{R} \rightarrow \mathbb{Y}$ is $\phi = \text{sign}$

we get HALFSPACE classifiers

$$HS_d = \phi \circ L_d = \{x \mapsto \text{sign}(h_{w,b}(x)), h_{w,b} \in \mathcal{L}_d\}$$

VC-dim. of HALFSpace hypothesis?

$$\mathcal{F} = \{x \mapsto \text{sign}(\langle w, x \rangle) : w \in \mathbb{R}^d\}$$

$$y = \{\pm 1\}$$

Theorem: The VC-dim of \mathcal{F} is $VC(\mathcal{F}) = d!$

Proof: we first show $VC(\mathcal{F}) \geq d$ (lower bound)

let $C = \{c_1, \dots, c_d\}$ where each $c_i \in \mathbb{R}^d$

$$c_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, c_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \dots, c_d = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$$

at 1st coord at 2nd at dth coord

Let $w = \begin{pmatrix} y_1 \\ \vdots \\ y_d \end{pmatrix}$ for any given labeling of $\{c_1, \dots, c_d\}$ ^{$\sum_{i=1}^d y_i = 1$}

c_{ij} denotes the j -th coord of c_i (and w_j the j -th coord of w)

$$\Rightarrow \langle w, c_i \rangle = \sum_{j=1}^d w_j \underbrace{c_{ij}}_{\substack{\text{only } 1 \text{ at } i=j}} = w_i = y_i \quad \checkmark$$

We found a set of size d that is shattered by \overline{T} .

$$\Rightarrow VC(\mathcal{F}) \geq d$$

we will show the upp. bound $V(F) < d+1$ is a contradiction!

Assume $C = \{c_1, \dots, c_d, c_{d+1}\}$ is shatt. by F .

That would mean we have $w_1, \dots, w_{2^{d+1}}$ weights that realize the 2^{d+1} possible labelings.

If we write all these labelings in a matrix like

$$\begin{pmatrix} w_1^T c_1 & w_2^T c_1 & \dots & w_{2^{d+1}}^T c_1 \\ w_1^T c_2 & w_2^T c_2 & \dots & w_{2^{d+1}}^T c_2 \\ \vdots & \vdots & \ddots & \vdots \\ w_1^T c_{d+1} & w_2^T c_{d+1} & \dots & w_{2^{d+1}}^T c_{d+1} \end{pmatrix}$$

we have $d+1$ rows and 2^{d+1} columns!

We can also write this matrix as a product XW with

$$X = \begin{pmatrix} -c_1 & - & \\ -c_2 & - & \\ & & \\ -c_{d+1} & - & \end{pmatrix} \quad \text{and} \quad W = \begin{pmatrix} | & | & & | \\ w_1 & w_2 & \dots & w_{2^{d+1}} \\ | & | & & | \end{pmatrix}$$

$$XW = \begin{pmatrix} \langle w_1, c_1 \rangle & \langle w_2, c_1 \rangle & \dots \\ \vdots & \ddots & \end{pmatrix}$$

$$X \in \mathbb{R}^{(d+1) \times d} \quad \text{and} \quad W \in \mathbb{R}^{d \times 2^{d+1}}$$

Taking $\text{sign}(XW)$ gives all possible labelings! (see step of Φ)
(in the columns of XW)

let $M = XW$. we know that

$$\text{rank}(M) \leq \min(\text{rank}(X), \text{rank}(W)) \leq d$$

as X has only d columns.

Claim: The rows of M are lin. independent (under stoff. ops.)

Proof: Rem. that v_1, \dots, v_k in a vec. space are lin. ind. if

$$\alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_k v_k = 0$$

is only satisfied iff $\forall i: \alpha_i = 0$!

What is M of size?

$$M = \begin{pmatrix} w_1^T c_1 & w_2^T c_1 & & w_{2^{d+1}}^T c_1 \\ w_1^T c_2 & w_2^T c_2 & & \\ \vdots & & & \\ w_1^T c_{d+1} & w_2^T c_{d+1} & \dots & w_{2^{d+1}}^T c_{d+1} \end{pmatrix} \begin{matrix} \rightarrow \text{row } 1 \\ \\ \\ \rightarrow \text{row } d+1 \end{matrix}$$

So,

$$Q_1 \cdot \text{row } 1 + Q_2 \cdot \text{row } 2 \dots + Q_{d+1} \text{row } d+1 =$$

$$Q_1 \cdot \begin{pmatrix} w_1^T c_1 & w_2^T c_1 & \dots & w_{2^{d+1}}^T c_1 \end{pmatrix} +$$

$$Q_2 \cdot \begin{pmatrix} w_1^T c_2 & w_2^T c_2 \end{pmatrix}$$

\vdots

$$Q_{d+1} \cdot \begin{pmatrix} w_1^T c_{d+1} & w_2^T c_{d+1} \end{pmatrix}$$

$$\left(\underbrace{Q^T X w_1 \quad Q^T X w_2}_{\text{and so on}} \right) \stackrel{?}{=} 0 \text{ when}$$

Fact is that there is a k s.t

$$\text{sign}(a) = \text{sign}(Xw_k) \text{ due to shattering assumption}$$

So, whatever a is (unless $\vec{0}$), there is always a k s.t Xw_k matches its sign $\Rightarrow a^T Xw_k = \text{sum of positive numbers and}$

$$\left(\begin{array}{ccccccc} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{array} \right) \in \mathbb{R}^{2^{d+1}}$$

at k is non-zero

\Rightarrow but ind. $d+1$ rows \square

This gives $\boxed{\text{rank}(M) = d+1}$

\Rightarrow CONTRADICTION (implying $VC(\mathcal{F}) < d+1$)

With $VC(\mathcal{F}) \geq d$ and $VC(\mathcal{F}) \leq d+1$ we have $VC(\mathcal{F}) = d!$

Learning halfspace classifiers

$$HS_d = \{x \mapsto h_{w,b}(x), h_{w,b} \in \mathcal{H}_d\}$$

we are going to look at the realizable case!

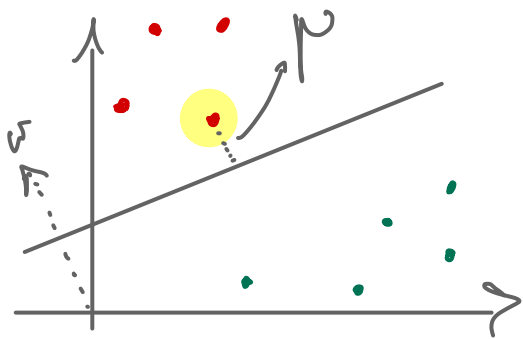
$$S = ((x_1, y_1), \dots, (x_m, y_m))$$

Training data $(y = \{\pm 1\})$ ^{with}

ERM on HS_d means to find $w \in \mathbb{R}^d$ (or \mathbb{R}^{d+1}) _s with bias, s.t.

$$\forall i \in \{1, \dots, m\} : \text{sign}(\langle w, x_i \rangle) = y_i(x)$$

or, equivalently $y_i \langle w, x_i \rangle > 0$ for all $i \in \{1, \dots, m\}$



$$\text{Let } \rho = \min_{i \in \{1, \dots, m\}} y_i \langle w^*, x_i \rangle$$

where w^* is a weight vector that satisfies (x).

Letting $\bar{w} = \frac{w^*}{\rho}$, we have

$$\begin{aligned} \forall i \in \{1, \dots, m\}: \quad y_i \langle \bar{w}, x_i \rangle &= y_i \langle \frac{w^*}{\rho}, x_i \rangle \\ &= \frac{1}{\rho} \cdot y_i \langle w^*, x_i \rangle \geq 1 \end{aligned}$$

This shows that a vector $w \in \mathbb{R}^d$ with

$\forall i: \boxed{y_i \langle w, x_i \rangle \geq 1}$ has to exist! (xx)

We can write (xx) as $Aw \geq v$ with

$$A = \begin{pmatrix} x_{11}y_1 & x_{12}y_1 & \dots & x_{1d}y_1 \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1}y_m & x_{m2}y_m & \dots & x_{md}y_m \end{pmatrix}, \quad w = \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix}, \quad v = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

all ones!

A linear program (LP) is

$$\underbrace{\max_{w \in \mathbb{R}^d} \langle u, w \rangle}_{\text{linear objective}} \quad \text{subject to} \quad \underbrace{Aw \geq v}_{\text{linear inequalities}}$$

In our case, we just set $u \in \mathbb{R}^d$ to some "dummy" vector, eg.

$$u = (1, \dots, 1)^T$$