# 18. Linear Predictors (part 1: binary classification)

Monday, February 19, 2024     4:28 PM

Starting Part Two of SS, "From Theory to Applications"

## Usefulness of linear predictors

- used in practice (or their generalizations like SVM)
- can be supercharged via kernel methods
- easy to understand, interpret and train
- used for classification and regression

↙ We really should say "affine"

**Def** Linear classifier space $L_D$

$$L_d = \{ h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R} \}, \quad h_{w,b}(x) = \langle w, x \rangle + b \quad \text{is affine}$$

(often work in $d+1$ dim and take $b=0$ wlog, i.e. $x_{(d+1)} := 1$, $w_{(d+1)} = b$

Later, we'll take $H = \phi \circ L_d$ where $\phi$ is identity for regression
$$\phi = \text{sign} \quad \text{for classification}$$

## §9.1 Binary Classification and Halfspaces

$$HS_d := \phi \circ L_d, \quad \phi = \text{sign}.$$

We've seen $VCdim(HS_d) = d+1 < \infty$ so it's PAC learnable

$$\left( \text{need } m = \Omega \left( \frac{d \cdot \log(1/\epsilon) + \log(1/\delta)}{\epsilon} \right) \text{ samples} \right)$$

$m = \Omega(g(m))$ means
$g(m) = O(m)$

...and ERM is a PAC learner!

## How to solve ERM for $HS_d$ ?

If data isn't separable (ie. not realizable) it can be very hard. In practice, use a surrogate loss function that's at least continuous (unlike 0-1 loss), eg. hinge (SVM) or logistic

So for now assume data is separable
and wlog look at homogeneous $(b=0)$ case

Mostly of theoretical and historical interest

So assume $\exists \; w_{oracle}$ st. $y_i = \text{sign}(\langle w_{oracle}, x_i \rangle)$

Find $w$ st.
$$y_i = \text{sign}(\langle w, x_i \rangle) \quad \Longleftrightarrow \quad y_i \cdot \langle w, x_i \rangle > 0$$

# 18a. Linear Predictors (part 1: binary classification)

Idea: solve this via <u>linear programming</u> $\left( \begin{array}{l} \min\limits_{w} \quad c^T w \\ \text{st. } A \cdot w = b \\ \quad\quad C \cdot w \geq d \end{array} \right)$

which are efficiently solvable in both theory <u>and</u> practice

$\quad\quad\quad$ sort of... depends on your model of computation.

But... our constraint was ... $> 0$ not ... $\geq 0$. In optimization, that's a big deal!

Fix: positive scaling of $w$ won't affect final output

$\quad$ so $\quad \forall w$ that correctly classifies the data, i.e. $y_i < w, x_i > \, > 0$

$\quad \exists$ a scaled version $\tilde{w} = c \cdot w$ $(c > 0)$ st. $\quad$ i.e. $c = \dfrac{1}{\min\limits_i <w, x_i>}$

$$y_i \cdot < \tilde{w}, x_i > \, \geq 1$$

**Alternative (historical) method:** Rosenblatt's "Perceptron" for ERM, 1958

"Batch perceptron": Data input $\left( (x_i, y_i) \right)_{i=1}^{m}$, $x_i \in \mathbb{R}^d$, $y_i = \pm 1$

$w^{(0)} = 0 \in \mathbb{R}^d$

For $t = 1, 2, \ldots$

$\quad\quad$ If $\exists \, i \in [m]$ s.t. $y_i < w^{(t)}, x_i > \, \leq 0$ $\quad$ i.e. misclassified

$\quad\quad\quad\quad w^{(t+1)} = w^{(t)} + y_i \cdot x_i$ $\quad\quad$ pushes in right direction on $x_i$

$\quad\quad$ Else

$\quad\quad\quad\quad$ Done! 0 training loss

Note $\quad y_i < w^{(t+1)}, x_i > \, = \underbrace{y_i < w^{(t)}, x_i >}_{\text{OLD}} + \underbrace{\|x_i\|^2}_{\text{pushes it positive.}}$

$\quad\quad\quad\quad \underbrace{\phantom{y_i}}_{\text{NEW}}$

**Thm** Perceptron converges (if data separable)

One way to prove: exercise 14.3, this is subgradient descent on

$$f(w) = \max_{i \in [m]} \left( 1 - y_i < w, x_i > \right)$$

$\quad\quad\quad\quad$ max preserves convexity

$\quad\quad\quad\quad\quad\quad$ no gradient if argmax not unique.

**Def** subgradient If $f$ is convex, $f : \mathbb{R}^d \to \mathbb{R}$, then $g \in \mathbb{R}^d$ is

$\quad$ a subgradient of $f$ at $w$ if $\forall w'$, $f(w') \geq f(w) + <g, w' - w>$

Sometimes we go in depth on convexity here...