

Ok, we've analyzed that for a given  $\mathcal{H}$   
 Uniform convergence  
 (Agnostic) PAC learnability  
 ERM being a (n Agnostic) PAC learner  
 $VCdim(\mathcal{H})$  is bounded  
 are equivalent!

ERM often isn't efficient, certainly in the agnostic setting...

But it is efficient for realizable linear separators:  $\mathcal{H}_{lin-d}$

Pick  $w \in \mathbb{R}^d$   
 $b \in \mathbb{R}$  s.t.  $sgn(\langle w, x_i \rangle + b) = y_i \quad \forall (x_i, y_i) \in S$

$$VCdim(\mathcal{H}_{lin-d}) = d+1$$

LB

$$\begin{matrix} x_0 = 0 \\ x_1 = e_1 \\ \vdots \\ x_d = e_d \end{matrix}$$

shatterable set:

$$\begin{matrix} \text{For } y_0 \\ \vdots \\ y_d \end{matrix}$$

$$\begin{aligned} \text{Pick } b = y_0 &\Rightarrow \langle w, x_0 \rangle + b = b = y_0 \checkmark \\ w_i = y_i - b &\Rightarrow \langle w, x_i \rangle + b \\ &= \langle w, e_i \rangle + b \\ &= w_i + b = y_i \checkmark \end{aligned}$$

UB

Take any set  $|C| > d$  in  $\mathbb{R}^d$ . [Let's ignore  $b$  and just argue about the bias-less ones]  
 WTS it cannot be shattered by  $\mathcal{H}_{lin-d}$   
 [find a set of labels we can't induce]  
 Since  $|C| > d+1$ ,  $C$  is not linearly  $\perp$   
 $\Rightarrow \exists \vec{a} \neq 0$  s.t.

$$\vec{a} \cdot \vec{x} = 0$$

$$\sum_{i \in P} a_i x_i = \sum_{j \in N} |a_j| x_j$$

$$\begin{aligned} P &\triangleq \{i : a_i > 0\} \\ N &\triangleq \{j : a_j < 0\} \end{aligned}$$

Neither can be empty. Why?

$\vec{a} \neq \vec{0}$  so at least one is nonempty

[If one is  $\neq \emptyset$ , the other can't be...]

Suppose  $w$  induces labels  $\text{sgn}(a_i)$  for  $x_i$   
 This means  $a_i \cdot \langle w, x_i \rangle > 0 \quad \forall i$ .

$$0 \leq \sum_{i \in P} a_i \langle w, x_i \rangle = \langle \sum_{i \in P} a_i x_i, w \rangle = \langle \sum_{j \in N} |a_j| x_j, w \rangle = \sum_{j \in N} |a_j| \langle x_j, w \rangle < 0$$

A contradiction!

So  $C$  can't be shattered  
 by  $\mathbb{R}^{n-d}$

Since  
 $\text{sgn} \langle x_j, w \rangle = -1$

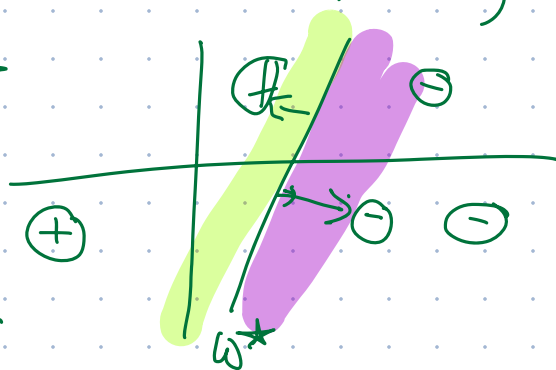
This means ERM on  $m \geq C \cdot \frac{d + \log \frac{1}{\epsilon}}{\epsilon}$  samples  
 PAC-learns linear <sup>realizable</sup> separators in  $\mathbb{R}^d$ .


[When not realizable, it's computationally hard!]

[Simple Alg: use Linear Programming! Other things also work..]

If our data is not only realizable but also  
 (separable)

Has a margin



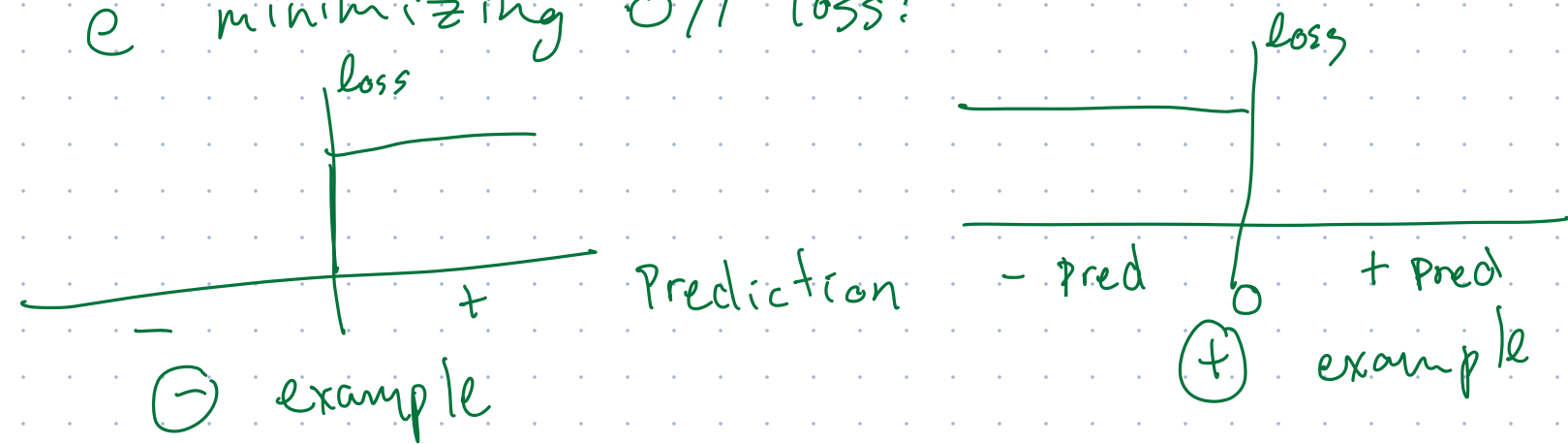
min width  
 of   
 is the margin  
 of  $w^*$

then we can learn  
 with many fewer  
 samples

....  $w^* = \underset{w}{\operatorname{argmin}} \|w\|$  st.  
 $\Rightarrow \quad \|x\| \leq R$

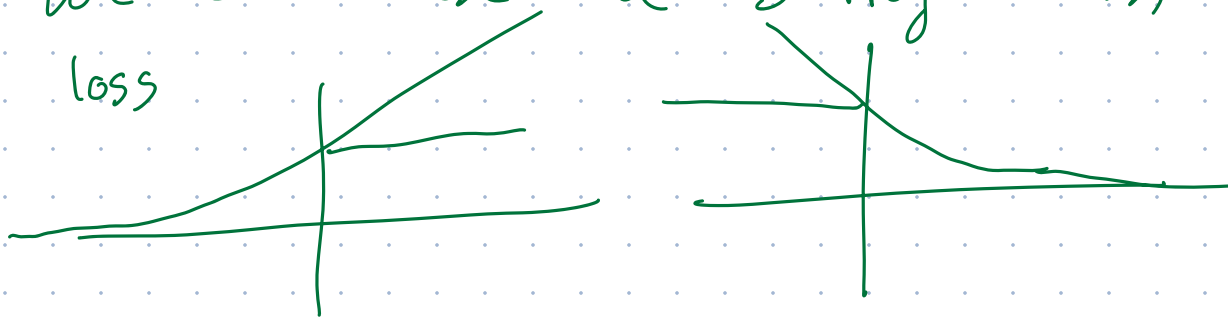
$$y_i \langle w, x_i \rangle \geq 1 \quad \forall i \in [n]$$

Finally, what is hard about non-realizable setting?  
 Our 2nd-fave alg / Gradient Descent, isn't good @ minimizing 0/1 loss:



It's non convex,  
 non-smooth  
 doesn't have a gradient!

We can use a surrogate loss, like logistic loss



convex, smooth, has gradients  $\ell_w(x, y) = \log(1 + \exp(-y \langle w, x \rangle))$

So,  $|L_D(\text{logistic}) - L_S(\text{logistic})|$  will be small (uniform convergence!) but  
 no guar that  $0-1-L_S(\text{logistic})$  will be small  
 even if  $\exists w$  w. small  $0-1 L_S$

What's another set of methods for arguing about good  $L_D(h)$ ?

$$L_D(h) = |L_D(h) - L_S(h)| + L_S(h)?$$