

Last time:

- Window 2 for  $F(\sigma) = \sigma$ -sep. monotone LTFs
- Perceptron alg, PCT:  $\frac{1}{\sigma^2}$  m.b. for unit  $x$ ,  
origin-sep.  $v \cdot x \geq 0$  target  $f_{\eta}$ , with "margin"  $\sigma$

Today:

- example of applying PCT
- Dual Perceptron, "kernelization"
- (maybe) start generic bounds for online learning

Questions?

---

Ex of applying PCT:

King with  $n$  advisors ( $n$  odd)

King feels subset  $S \subseteq [n]$  give bad advice

Binary dec: gets all  $n$  opinions, negates votes in  $S$ ,  
goes w/ maj vote.

You're political analyst. Observe repeatedly.

How many mistakes? (predicting King's decision  
given advisors' opinions)

A: Use Perceptron.

$\text{MAJ}(x_1, \dots, x_n) : \{0, 1\}^n \rightarrow \{0, 1\}$

$$\ll \mathbb{I} \left[ x_1 + \dots + x_n \geq \frac{n}{2} \right]$$

Target: MAJ( $x_1, \bar{x}_2, \bar{x}_3, x_4, x_5, \bar{x}_6, x_7, x_8$ )

$$S = \{2, 3, 6\}$$

$\pm 1$  view on input:

$$\{\pm 1\}^n \rightarrow \{\pm 1\}$$

$$\text{MAJ}(x) = \text{sign}(x_1 + \dots + x_n)$$

"S-regated maj":  $\text{sign}(v_1 x_1 + \dots + v_n x_n)$

$$v_i = \begin{cases} +1 \\ -1 \end{cases} \quad i \in S$$

Perc alg wants  $v$  unit vector:

normalize

$$v_i = \begin{cases} +1/\sqrt{n} \\ -1/\sqrt{n} \end{cases}$$



each  $x_i = \begin{cases} +1/\sqrt{n} \\ -1/\sqrt{n} \end{cases}$

every  $\|x\|=1, \|v\|=1$

PCT:

$$|v \cdot x| \geq ? \quad \frac{1}{n}$$

$$|v \cdot x| = |v_1 x_1 + v_2 x_2 + \dots + v_n x_n| =$$

$\begin{matrix} \swarrow & \downarrow \\ \pm 1/\sqrt{n} & \pm 1/\sqrt{n} \end{matrix}$

sum of  $n$   $\pm 1/n$

$\rightarrow \geq \frac{1}{n}$

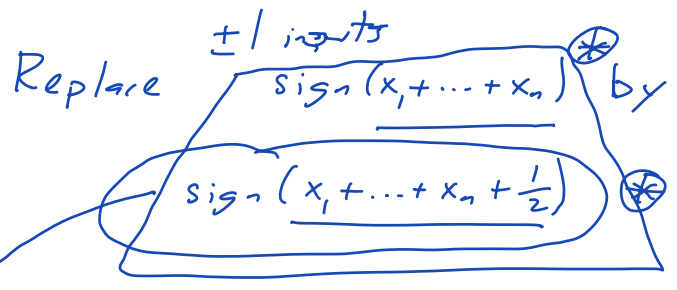
$$\mathcal{J} = \frac{1}{n}$$

So PCT  $\Rightarrow$  make  $\leq \frac{1}{\sqrt{2}} \leq n^2$  mist. (3)

---

What if  $n$  even, + MAJ is  $\mathbb{I}[x_1 + \dots + x_n \geq \frac{n}{2}]$  ?

$\pm 1$   $x_1 + \dots + x_n$  could be 0



$$\text{sign}(t) = \begin{cases} +1 & t \geq 0 \\ -1 & t < 0 \end{cases}$$

over  $\{\pm 1\}^n$ , logically equivalent

---

### Dual Perceptron + "Kernel Functions"

Nice thing about Perc: has a "dual form" (exactly equivalent, looks different - will make it possible to use "kernels").

Key idea of dual Perc: to run Perc, only need to be able to compute inner prod. of 2 examples.

Running Perc: init wt vector  $w = 0^n$ .

After 1st mist (on  $x^i$ ), now  $w$  is  $x^i$  or  $-x^i$ .

So  $w \cdot x$  is  $\underline{x^1 \cdot x}$  or  $\underline{-x^1 \cdot x}$  (inner prod of 2 ex).

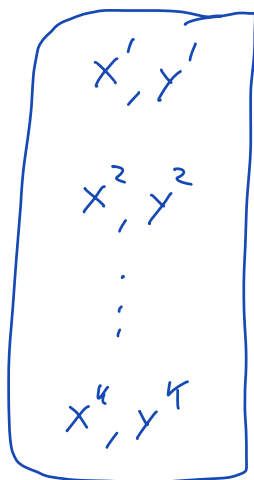
2<sup>nd</sup> mist, on  $x^2$ :  $w$  becomes, say,  $\underline{-x^1 + x^2}$

To compute  $w \cdot x$ , compute  $\underline{-x^1 \cdot x} + \underline{x^2 \cdot x}$

After  $k$  updates, can compute  $w \cdot x$  as  $\pm 1$  sum of  $k$  i.p.'s of examples.

More formal descrip. of dual Perc:

Maintain hyp as list of pairs (lab ex's)



$x^i = i^{\text{th}}$  ex on which  
we made mistake  
 $y^i \in \{\pm 1\} = \overset{\text{correct}}{\text{label of } x^i}$ .

Given new  $x$ ,  $w \cdot x = \sum_{i=1}^k y^i (x^i \cdot x)$ .

Can

Exactly simulate Perceptron this way.

---

Why? Seems ineff - after  $k$   
mist, hyp. is of size  $(n \cdot k)$ .

True!

But... can replace "inner product" with any "kernel function".  
(this lets us run dual Perce. over any space for which we have an eff. evaluable kernel function!)

Kernel functions + feature expansions.

Let  $X =$  original feature space  
 $\{0, 1\}^n$  or  $\mathbb{R}^n$   
( $X$  has an inner product).

Let  $X' =$  aspirational feature space,  
also with an inner product:

$\{0, 1\}^N$  or  $\mathbb{R}^N$

A feature expansion is a mapping  $\Phi: X \rightarrow X'$ .

Ex:

$$X = \{0, 1\}^n, \quad X' = \{0, 1\}^{3^n = N}$$

$$\Phi: \{0, 1\}^n \rightarrow \{0, 1\}^{3^n}$$

$(x_1, \dots, x_n) \rightarrow$  (list of all  $3^n$  conj over  $x_1, \dots, x_n$ )

$$n=2: \overline{\Phi}(x_1, x_2) =$$

empty conj

$$(1, x_1, \overline{x_1}, x_2, \overline{x_2}, x_1 \wedge x_2, \overline{x_1} \wedge \overline{x_2}, x_1 \wedge \overline{x_2}, \overline{x_1} \wedge x_2)$$

$\overline{\Phi}(x)$  is "expanded" version of  $x$ .

The kernel function for feat. exp.  $\Phi: X \rightarrow X'$

is  $K: X \times X \rightarrow \mathbb{R}$ ,

$$K(a, b) = \overbrace{\underbrace{\Phi(a)}_{\text{in } X'} \cdot \underbrace{\Phi(b)}_{\text{in } X'}}^{\text{in } \mathbb{R}}$$

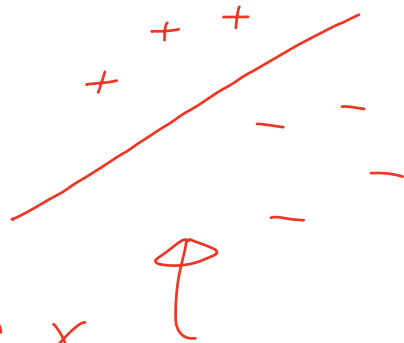
→ i.e., inner product in  $X'$ -space.

Why do feature exp?

Consider LTFs over  $\mathbb{R}^2$

$$X = \begin{matrix} \text{all} \\ (x_1, x_2) \\ \text{pairs} \end{matrix}$$

$(x_1, x_2)$  orig features



Consider  $X' = \text{feat. exp. of } X$

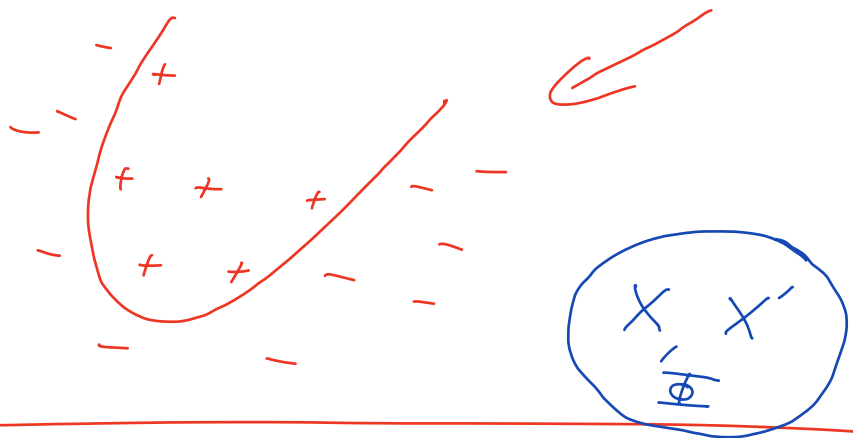
obt. by allowing monom. of degree up to 2:

$$X' = \text{tuples } (1, x_1, x_2, x_1 x_2, x_1^2, x_2^2)$$

$$\Phi(x) = \Phi(x_1, x_2) = (1, x_1, x_2, x_1 x_2, x_1^2, x_2^2)$$

LTFs over  $X'$ : ← richer!

$$\mathbb{1} [3x_1^2 - 2x_1 x_2 + 5x_2 \geq 3]$$



Suppose you have data in  $X$ , but your dream is to run Perceptron on the  $\Phi$ -feature-expanded data in  $X'$ .

Obvious approach: do it directly.

Each  $x$  you get: write down  $\Phi(x)$ .

Maintain  $w'$  in  $X'$  domain.

Legit, but could be super-slow:  
time  $N$  to write down  $\Phi(x)$ .

---

Better way: use kernel  $K$  for  $\Phi$ ,  
run dual Perc., replacing every  $a \cdot b$  (inner  
prod) by  $K(a,b)$ . This exactly simulates  
running Perceptron over the  $\Phi$ -expanded  
examples! GREAT!!!.....

---

→ if can efficiently compute  $K(a,b)$ .....

---

**Good news:** sometimes (for some  
interesting feat. exp.  $\Phi$ ), there is  
a computationally efficient way to  
compute  $K(a,b)$ , much faster than  
writing  $\Phi(a)$ ,  $\Phi(b)$  & doing it explicitly.



Ex  
Back to

$\Phi(x) = \text{all-conj feat. exp.}$

$$X = \{0, 1\}^n, \quad X' = \{0, 1\}^{3^n} = N$$

Claim:  $K(a, b)$  for this can be computed  
in  $O(n)$  time!

PF: Fix  $a, b \in \{0, 1\}^n$

Let  $\text{SAME}(a, b)$  denote # coords in  $\{1, \dots, n\}$   
where  $a_i = b_i$

$n = 16$   
 $\downarrow a_8 \neq b_8$

$$a = 1111111100000000$$
$$b = 1100000011111000$$

$\text{SAME}(a, b) = 5$

each coord: 0 or 1

$$\Phi(a) = (1, a_1, \bar{a}_1, a_2, \bar{a}_2, \dots, a_1, \bar{a}_1, a_2, \bar{a}_2, \dots)$$

$$\Phi(b) = (1, b_1, \bar{b}_1, b_2, \bar{b}_2, \dots, b_1, \bar{b}_1, b_2, \bar{b}_2, \dots)$$

$\Phi(a) \cdot \Phi(b) = \# \text{ of the } 3^n \text{ positions with a 1}$   
in both vectors.

since  $a_8 \neq b_8$ , any conj involving  $x_8$  or  $\bar{x}_8$   
won't be sat. by both.

Only conj that could give  $1 \cdot 1$  in inner prod. must be <sup>over</sup> vars where  $a_i = b_i$ :

in above ex, these are conj over only

$$x_1, x_2, x_{14}, x_{15}, x_{16}.$$

$$\underbrace{\quad} \quad \underbrace{\quad} \quad S \subseteq \{1, 2, 14, 15, 16\}$$

$$S = \{1, 14\} : \quad \underbrace{x_1, \wedge \bar{x}_{14}} \quad \underbrace{a+b \text{ both eval. to } 1}$$

$$S = \{14, 15, 16\} : \quad \bar{x}_{14} \wedge \bar{x}_{15} \wedge \bar{x}_{16}$$

$$\text{So } K(a, b) = 2^{\text{SAME}(a, b)}$$

So can run Perc. over  $\xrightarrow{\text{(dual)}} \underline{\underline{3^7\text{-dim}}}$  space of all conj!

⊢ If Perc makes  $K$  mistakes, runtime per trial is  $\leq \underline{\underline{\text{poly}(n, k)}}$

---

---

Next time: generic  $\mathcal{C}$ 's.

---