

- Last time:
- discuss Winnow 1 $\rightarrow F(\delta)$ "δ-sep. over LTFs" LTFs
 - Winnow 2 alg: OLMAB alg. for certain LTF's (with positive weights, $X = \{0, 1\}^n$)
 - Perceptron alg. for LTFs "with a margin" over \mathbb{R}^n

- Today:
- finish PCT proof; example of using PCT
 - dual Perceptron, Kernelization *End of our "concrete examples"*
 - start generic bounds: "halving algorithm"

Questions?

Recall PCT setup:

- $v = \text{target vector}$ $\|v\|=1$
- $w = \text{hyp. vector}$
- $\hookrightarrow \text{initially } 0^n$

Theorem (Perc. Convergence Theorem):

Suppose run Perc. on seq of ex labeled by $\text{sign}(v \cdot x)$, where $\|x\|$ each $\frac{\|x\|}{\|v\|}$ are ≥ 1 .

Let $\delta := \min_{\substack{\|x\| \\ \text{ex. } x \\ \text{weight}}} |v \cdot x|$. ($= \min_{\substack{\|x\| \\ \text{ex. } x \\ \text{weight}}} \text{Euclidean dist. from } x \text{ to } v^\perp$)

Then #mist Perc. makes is $\leq \frac{1}{\delta^2}$.

L1: After M mistakes, $v \cdot w \geq \delta \cdot M$

L2: After M mist., $w \cdot w \leq M$.

$$(\|w\| = \sqrt{w \cdot w} \leq \sqrt{M})$$

Proof of PCT:

$$\begin{aligned} \text{length of proj of } w \text{ in dir. } v &\leq (\text{length of } w) \\ \delta M \stackrel{(L1)}{\leq} v \cdot w &\leq \|w\| = \sqrt{v \cdot w} \stackrel{(L2)}{\leq} \sqrt{M}, \text{ so} \\ \delta \sqrt{M} \leq 1, \text{ so } M &\leq \frac{1}{\delta^2}. \end{aligned}$$

$v \cdot w = \|w\| \cdot \cos(\alpha)$

b/c v = unit vector



Discussion:

- + simple!
- + noise-tolerant
- + fast
- + "Kernelizable" (more soon...)

(better MB)

- $\frac{1}{\delta^2}$: not great.

Faster LTF learning known, with m.b. $\approx \log \frac{1}{\delta}$
but these lack nice properties.

Ex: Class of "reoriented n-var MAJ's"

(n odd) $\mathcal{C} = \{ \text{MAJ}(l_1, l_2, \dots, l_n) : l_i \text{ either } x_i \text{ or } \bar{x}_i \}$

$|\mathcal{C}| = 2^n$

$X = \{0, 1\}^n$

$$\text{ex: } n=3 \quad \text{MAJ}(x, \bar{x}_2, \bar{x}_3)$$

To use PCT: want each ex $\|x\|=1$,
target $\|v\|=1$, $\text{sign}(v \cdot x)$

Remap: $x \in \{0, 1\}^n$ $0 \rightarrow -\frac{1}{\sqrt{n}}$
 \downarrow
 $z \in \left\{-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right\}^n$. $1 \rightarrow \frac{1}{\sqrt{n}}$

every $z \in \left\{-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right\}^n$ has $\|z\|=1$.

Have $\underline{\text{sign}(v \cdot z)} = \text{MAJ}(l_1, \dots, l_n)$

by $x \mapsto z$ as above,

$$z_i = \begin{cases} 1 & \text{if } l_i = x_i \\ -1 & \text{if } l_i = \bar{x}_i \end{cases}$$

$v \cdot z$ adds up $\pm z_1, \pm z_2, \dots, \pm z_n$

$$+ \text{ iff } \text{MAJ}(l_1, \dots, l_n) = 1$$

$v \leftarrow \frac{1}{\sqrt{n}}$ this v :

$$v_i = \pm \frac{1}{\sqrt{n}} \text{ so } \|v\|=1$$

What's σ ? $\sigma = \min |v \cdot z|$

$$v \cdot z = \left(\frac{\pm 1}{\sqrt{n}} \cdot \frac{\pm 1}{\sqrt{n}} \right) + \dots$$

n of these

$|v \cdot z| \geq \frac{1}{\gamma}$ so $\delta = \frac{1}{\gamma} + \text{Perc. will make } \leq \gamma^2 \text{ mistakes.}$

① Dual Perceptron + ② Kernelization

↳ Different, \equiv way to run Perceptron

Clunkier... but enables " - way to run

Perc. over very high-dim feature space super-efficiently!

① Key insight: all it takes to run Perc. is being able to compute $x \cdot x'$ where x, x' are example vectors.

Recall how Perc works: $w_{\text{init}} = (0, \dots, 0)$

After 1st update, on example x' ,
w is either x' or $-x'$.

To ^{keep} runningPerceptron: next ex x , $w \cdot x$ is $x \cdot x'$ or $-x \cdot x'$

2nd mist is on, say, x^2 , & new w is (say) $x^1 - x^2$

So $w \cdot x = x^1 \cdot x - x^2 \cdot x$ Can keep going!!

More formal descr. of dual Perc:

Hyp. maintained as list of pairs

$\bigcirc \star \quad \left\{ \begin{array}{l} (x^1, y^1) \\ \vdots \\ (x^k, y^k) \end{array} \right\}$

 (x^i, y^i) = i^{th} example Pera. got wrong.
 "label $\in \{-1, 1\}$ "

Given new example x , compute

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

Perfectly simulates Perceptron.

Why?

- More memory-intensive, slower
- " " " We can replace " "
 (inner prod. over original ex)

with any "Kernel function".

② Kernel functions + "feature expansions"

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

X' = new "high dim" feature space
 $(\{0, 1\}^N \text{ or } \mathbb{R}^N)$

both spaces have inner products.

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

Ex: $X = \{0, 1\}^n$, $X' = \{0, 1\}^3$,

Φ = "all-conjunctions" feature expansion

$\Phi(x)$ has a coord for each of the 3ⁿ conj. over x_1, \dots, x_n .

$n=2$: $x = (x_1, x_2) \in \{0, 1\}^2$

$\Phi(x)$ = 9-bit string

$$= (1, \underset{\text{empty conj}}{\cancel{x_1}}, \cancel{x_1}, x_2, x_1 x_2, \cancel{x_1} x_2, \cancel{x_2}, x_1 \cancel{x_2}, \cancel{x_1} \cancel{x_2})$$

$\Phi(x)$ = "expanded" version of x .

A Kernel function K corresp. to Φ is a fn

$$K: X \times X \rightarrow \mathbb{R}, \quad K(a, b) = \underbrace{\Phi(a)}_{\in X} \cdot \underbrace{\Phi(b)}_{\in X} \in \mathbb{R}$$

Back to learning:

Suppose we get ex. in X' , but we want to run Perc. over X' using $\Phi(x)$ for each ex. x .

• Obvious approach: do it directly.
for each x we get, write down $\Phi(x)$,
maintain $w \in \mathbb{R}^N$ as hyp. vector, run Perc. on.

Drawback: takes time $\geq N$ per example - slow!

• Better way: use dual Perc., replacing each $a \cdot b$ with $K(a, b)$.

This exactly simulates running Perc. on $\Phi(\cdot)$ -versions of examples!

If we can evaluate $K(a, b)$ "fast", we can run Perc. over X' fast!

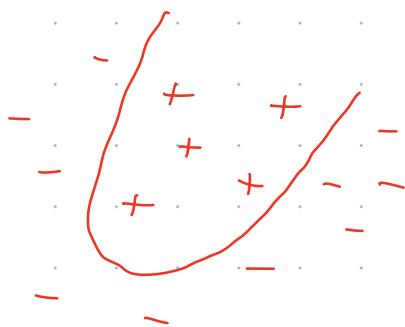
Nice to run Perc over feature-expanded space:

over \mathbb{R}^n , Perc. hyp. is like $\text{sign}(3x_1 + 4x_2 - 5x_3)$

If we did a feature exp. to

$$\Phi(x) = (1, x_1, \dots, x_n, x_1^2, x_1 x_2, \dots, x_1 x_n, x_2^2, x_2 x_3, \dots, x_n^2)$$

Perc hyp. can be $\text{sign}(5x_1^2 - 4x_1 x_4 + 3x_2^2)$



Good news: for some interesting Φ 's, can eval. the associated $K(a, b)$ very fast!

Ex: "all-conj feature exp" Φ from above.

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

Φ = "all-conjunctions" feature expansion

$\Phi(x)$ has a coord for each of the 3^n

conj. over x_1, \dots, x_n .

$$K(a, b) = \underline{\Phi(a) \cdot \Phi(b)}$$

Given $a, b \in \{0, 1\}^n$, define $SAME(a, b) \subseteq [n]$

$$SAME(a, b) = \{i \in [n] : a_i = b_i\}.$$

Ex: $n = 16$

$a = 1111111000000000$

$b = 1100000111110000$

$SAME(a, b) = \{1, 2, 14, 15, 16\}$

Diagram showing a and b as binary strings of length 16. Circles around $a_1, a_2, a_{14}, a_{15}, a_{16}$ and $b_1, b_2, b_{14}, b_{15}, b_{16}$ indicate they are equal. Brackets below the strings group them into pairs: (a_1, b_1) , (a_2, b_2) , \dots , (a_{14}, b_{14}) , (a_{15}, b_{15}) , (a_{16}, b_{16}) . A bracket above the last three groups indicates they are not equal.

Claim: $K(a, b) = 2^{|SAME(a, b)|}$.

Pf:

$$\Phi(a) = (1, \underbrace{a_1, \bar{a}_1, a_2, \bar{a}_2, \dots, a_{14}, \bar{a}_{14}, a_{15}, \bar{a}_{15}, a_{16}, \bar{a}_{16}}, \dots)$$

Diagram showing $\Phi(a)$ as a sequence of bits. Brackets group bits into pairs: $(a_1, \bar{a}_1), (a_2, \bar{a}_2), \dots, (a_{14}, \bar{a}_{14}), (a_{15}, \bar{a}_{15}), (a_{16}, \bar{a}_{16})$. The bit a_{15} is highlighted with a box and labeled "not in SAME".

$$\Phi(b) = (1, \underbrace{b_1, \bar{b}_1, b_2, \bar{b}_2, \dots, b_{14}, \bar{b}_{14}, b_{15}, \bar{b}_{15}, b_{16}, \bar{b}_{16}}, \dots)$$

Diagram showing $\Phi(b)$ as a sequence of bits. Brackets group bits into pairs: $(b_1, \bar{b}_1), (b_2, \bar{b}_2), \dots, (b_{14}, \bar{b}_{14}), (b_{15}, \bar{b}_{15}), (b_{16}, \bar{b}_{16})$. The bit b_{15} is highlighted with a box and labeled "not in SAME".

$\Phi(a) \cdot \Phi(b) = \# \text{ posit. where } \Phi(a) + \Phi(b) \text{ both are } 1 \text{ in that pos.}$

Any conj containing a var outside $SAME(a, b)$:

gives 1·0 or 0·1 or 0·0 $\rightarrow 0$

What abt conj. only over vars in $SAME$?

For each $i \in SAME$

Must either not include x_i or include

$$\begin{cases} x_i & \text{if } a_i = b_i = 1 \\ \bar{x}_i & \text{if } a_i = b_i = 0. \end{cases}$$

So $\sum |ISAME(a, b)|$ conj. that are sat.

both by a + by b ;

$$so K(a, b) = \sum |ISAME(a, b)| \rightarrow O(1) \text{ time}$$

So... can run Perc. over space of all conj.,
+ after K mistakes, it takes $O(K \cdot n)$
time per trial.

End of "specific" algos

Next time: generic algorithms