

Learning Halfspaces

Wednesday, October 6, 2021 8:05 AM

yours

$$w, b : \{x : w^T x \geq b\}.$$

- disjunctions
- conjunctions
- decision lists

How many halfspaces? infinite?

Effectively $\leq 2^n$ over $\{0, 1\}^n$. Why?

What about over Reals / Rationals? Later...

we can assume $b=0$. $w^T x - b \geq 0$

$$(w, b)^T \begin{pmatrix} x \\ -1 \end{pmatrix} \geq 0.$$

and assume $\|w\|_2 = 1$

and $\|x\| \leq 1$ (or $= 1$) .

Perception

Perception

1. $\omega = 0$
 2. On next example x , Predict $\text{sign}(\omega \cdot x)$
if mistake, $\omega \leftarrow \omega + l(x)x$
-

Thm. Given data classifiable by a halfspace ω^* , $\|\omega^*\| = 1$, with margin γ ($\min_x |\omega^* \cdot x|$), Perceptron makes at most $\frac{1}{\gamma^2}$ mistakes.

pf. consider $\cos(\omega, \omega^*) = \frac{\omega^T \omega^*}{\|\omega\|}$

starts at 0.

At each mistake

$$\omega \leftarrow \omega + l(x)x$$

$$\omega \cdot \omega^* \leftarrow \omega \cdot \omega^* + \underbrace{l(x)}_{\text{same sign}} \underbrace{\omega^* \cdot x}_{\omega \cdot \omega^* \text{ goes up by } \geq \gamma}$$

$\omega \cdot \omega^*$ goes up by $\geq \gamma$.

After t mistakes

$$\omega \cdot \omega^* \geq \gamma t.$$

$$\|\omega\|^2 \leq (\omega + l(x)x)^T (\omega + l(x)x)$$

- ... - T v

$$\begin{aligned}
 \|w\| &\leftarrow (w + \gamma(x)x) \cdot (w^T x + \dots) \\
 &= \|w\|^2 + \|x\|^2 + 2 \underbrace{\gamma(x)}_{\text{opp. sign}} \frac{w^T x}{\gamma} \\
 &\leq \|w\|^2 + 1.
 \end{aligned}$$

After t steps $\|w\|^2 \leq t$.

$$\therefore \cos(w, x) \geq \frac{\gamma t}{\sqrt{t}}$$

$$\text{Since } \cos \approx 1, \quad \gamma \sqrt{t} \leq 1 \Rightarrow t \leq \frac{1}{\gamma^2}.$$

More generally, # mistakes $\leq \frac{\|w^*\| \cdot \|x\|}{\gamma^2}$.

What if γ is super tiny?

- Modified Perceptron: correctly classifies all x with $|w^T x| > \gamma$, makes $O(\frac{\log n}{\gamma^2})$ mistakes.
- Using Linear Programming, can get $\log \frac{1}{\gamma}$ dependence.

Kernels.

mapping $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^m$ nonlinear.

$$k(x, y) = \langle \phi(x), \phi(y) \rangle = \phi(x)^T \phi(y).$$

legal kernel if such mapping exists.

i.e. $k \geq 0$. PSD.

e.g. $\phi(x) = x$. $k(x, y) = x^T y$

$$k(x, y) = \langle x, y \rangle^d$$

$$k(x, y) = (1 + \langle x, y \rangle)^d$$

$$k(x, y) = e^{-\|x-y\|^2}$$

legal ✓.

Suppose $\not\exists$ halfspace matching $l(x)$

but \exists halfspace in some nonlinear map.

Then if you know the map $x \rightarrow \phi(x)$

apply Perceptron to $\phi(x)$. . .

, ... ,

apply Perceptron

What if you have access to $K(x, y)$
 or ϕ is in very high (or infinite) dimension?

$$\text{we maintain } w = \sum_i l(x^{(i)}) x^{(i)}$$

$$w \cdot x = \sum_i l(x^{(i)}) (x^{(i)} \cdot x)$$

Instead, Kernel Perceptron:

$$w = \sum_i l(x) \phi(x^{(i)}) \quad \text{implicit}$$

$$w \cdot x = \sum_i l(x^{(i)}) K(x^{(i)}, x)$$

keep all examples on which a mistake was made.
 output as prediction on next x ,
 $\text{sign} \left(\sum_i l(x^{(i)}) K(x^{(i)}, x) \right)$.

Modified Perceptron

$w \leftarrow$ random unit vector

On mistake, if $|w \cdot x| > \sigma \|w\|$

$$w \leftarrow w + l(x) (w \cdot x) X$$

$$\text{Thm } \# \text{ mistakes} = O\left(\frac{\log n}{\sigma^2}\right)$$

Pf.: assume $l(x) = +$ (use flip $x \rightarrow -x$)

$$w \leftarrow w - (w \cdot x) X$$

$$w^* \cdot w_{\text{initial}} \geq \frac{1}{\sqrt{n}} \quad \text{with Prob} \geq \frac{1}{8}.$$

$$w^* \cdot w \geq w^* \cdot w - \frac{(w \cdot x)(w^* \cdot x)}{\text{opp}} \geq w^* \cdot w \geq \frac{1}{\sqrt{n}}.$$

$$w \cdot w \geq w \cdot w - (w \cdot x)^2 \geq \|w\|^2(1-\sigma^2)$$

$$\text{after } t \text{ steps } \cos(w, w^*) \geq \frac{\frac{1}{\sqrt{n}}}{(1-\sigma^2)^{t/2}}$$

$$\Rightarrow t \leq \frac{\ln n}{\sigma^2}.$$
