

The Perceptron Algorithm

①

An online algorithm for learning threshold functions.

$$f: \mathbb{R}^n \rightarrow \{+1, -1\}$$

$$\vec{x} \in \mathbb{R}^n \quad f_{\vec{w}, b}(\vec{x}) = \begin{cases} +1 & \text{if } \vec{w} \cdot \vec{x} \geq b \\ -1 & \text{if } \vec{w} \cdot \vec{x} < b \end{cases}$$

This defines a hyperplane in \mathbb{R}^n with \vec{w} as the normal to the hyperplane. The goal is to learn which hyperplane separates the positive labelled points from the negative ~~labeled~~ labelled points. Without loss of generality we can set $b=0$. We can redefine the problem in \mathbb{R}^{n+1}

$$\bar{f}: \mathbb{R}^{n+1} \rightarrow \{+1, -1\}$$

and change $\vec{x} \in \mathbb{R}^n$ to $\vec{x}' = \begin{pmatrix} \vec{x} \\ 1 \end{pmatrix} \in \mathbb{R}^{n+1}$

and \vec{w} to $\vec{w}' \in \mathbb{R}^{n+1}$ where $\vec{w}' = \begin{pmatrix} \vec{w} \\ -b \end{pmatrix}$

Now $\vec{x}' \cdot \vec{w}' \geq 0 \Leftrightarrow \vec{x} \cdot \vec{w} \geq b$.

We want to learn \vec{w} from examples. The examples are presented to us in online fashion. So we have to make a prediction and then the true label is revealed. We seek bounds on our mistakes.

BASIC PERCEPTRON ALGORITHM:

Init: $\vec{w}_1 = 0$

LOOP
 $t=1 \dots T$

(PREDICT) $\left\{ \begin{array}{l} \text{Obtain } \vec{x}_t \\ \text{if } \vec{w} \cdot \vec{x}_t \geq 0 \text{ then } \hat{y}_t = +1 \text{ else } \hat{y}_t = -1 \end{array} \right.$

(UPDATE)

$\left\{ \begin{array}{l} \text{Obtain } y_t \\ \text{if } \hat{y}_t \neq y_t \text{ then } \vec{w}_{t+1} \leftarrow \vec{w}_t + y_t \vec{x}_t \text{ else } \vec{w}_{t+1} = \vec{w}_t \end{array} \right.$

When it makes a mistake on a positive example it shifts the weight vector towards that point if it makes a mistake on a negative example it shifts the weight vector away.

Crucial concept: How close is a point to the hyperplane
Let f be a threshold function determined by a hyperplane passing through the origin with unit normal \vec{u} .

Def The margin of f at a labelled point (\vec{x}, y) is defined to be $y(\vec{x} \cdot \vec{u})$.

Remark If the margin is positive then f classifies (\vec{x}, y) correctly and if it is negative then f classifies the point incorrectly. The numerical value of the margin gives the distance to the decision boundary. The margin of a ^{finite} set of points is the minimum margin of any point in the set. The margin measures robustness of the classification.

Thm Suppose $(\vec{x}_1, y_1), (\vec{x}_2, y_2) \dots (\vec{x}_T, y_T)$ are such that $\|\vec{x}_t\| \leq D$ for some $D > 0$ and all $t \in \{1, \dots, T\}$.

Suppose $\exists \vec{u}$ a unit vector & $\gamma > 0$ such that $\forall t \quad y_t(\vec{x}_t \cdot \vec{u}) \geq \gamma$. Then the perceptron algorithm makes at most $(D/\gamma)^2$ mistakes.

Proof Let m_t be the number of mistakes just before round t . Thus $m_1 = 0$. We proceed by breaking the proof into 2 lemmas.

Lemma 1 $\forall t \in \{1, \dots, T\} \quad \vec{w}_t \cdot \vec{u} \geq m_t \gamma$

Proof By induction on t .

Base case $t=1, m_1=0$ so clearly $\vec{w}_1 \cdot \vec{u} = 0 = m_1 \gamma$

Induction step: suppose there is a mistake in round t . Then we have

$$\begin{aligned} \vec{w}_{t+1} \cdot \vec{u} &= (\vec{w}_t + y_t \vec{x}_t) \cdot \vec{u} = \vec{w}_t \cdot \vec{u} + y_t \vec{x}_t \cdot \vec{u} \\ &\geq \underbrace{m_t \gamma}_{\text{Ind Hyp}} + \underbrace{\gamma}_{\text{Margin assumption}} = m_{t+1} \gamma \end{aligned}$$

3

So lemma 1 shows that as mistakes are made the projection of \vec{w}_t on \vec{u} gets longer.

lemma 2 $\|\vec{w}_t\|^2 \leq m_t D^2$

Proof Induction on t ; base case is trivial as is the case when there is no mistake. Suppose there is a mistake in round t . Then we have

$$\begin{aligned} \|\vec{w}_{t+1}\|^2 &= \|\vec{w}_t + y_t \vec{x}_t\|^2 \\ &= \|\vec{w}_t\|^2 + \|\vec{x}_t\|^2 + \underbrace{2y_t (\vec{w}_t \cdot \vec{x}_t)}_{\text{Negative}} \end{aligned}$$

$$\leq \|\vec{w}_t\|^2 + \|\vec{x}_t\|^2$$

$$\leq m_t D^2 + D^2 = m_{t+1} D^2$$

Proof of Thm $\|\vec{w}_t\|^2 \leq m_t D^2$ (lemma 2)

$$\text{or } D \sqrt{m_t} \geq \|\vec{w}_t\| = \|\vec{w}_t\| \|\vec{u}\|$$

$$\geq \vec{w}_t \cdot \vec{u} \quad [\text{Cauchy-Schwartz}]$$

$$\geq m_t r$$

$$\text{or } D/r \geq \sqrt{m_t}$$

$$\text{or } m_t \leq (D/r)^2 \quad \blacksquare$$

APPLICATION Suppose there are n financial analysts who predict every day whether the market will go up or down. We represent each prediction as a vector in $\{+1, -1\}^n$. We would like to use perceptron to figure out whom to follow. Suppose there are k "experts" within the group s.t. a majority vote among these k always gives the right answer.

Define $\vec{u} = \frac{1}{\sqrt{k}} (0, \dots, 0, 1, 1, 0, \dots, 0, 1, \dots)$ where the entry is 1 if it corresponds to one of the k experts.

Then $y_t (\vec{x}_t \cdot \vec{u}) \geq \frac{1}{\sqrt{k}}$ since this subset is always right.

Thus $\|\vec{x}_t\| \leq \sqrt{n} = D$ & $r = \frac{1}{\sqrt{k}}$. So the mistake bound is $n k$.