

① Perceptron Convergence Theorem (PCT)

② Geometric Margin (γ_{geom})

Learning linear classifiers through the origin

$$\mathcal{F} = \left\{ f: \mathbb{R}^d \rightarrow \{-1, +1\}, f_{\underline{\theta}}(\underline{x}) = \text{sgn}(\langle \underline{\theta}, \underline{x} \rangle) \right\}$$

↗ feature space ↗ for some $\underline{\theta} \in \mathbb{R}^d$

all linear classifiers

$\theta \rightarrow$ parametrize the function function $f \Rightarrow f_\theta$ or $f(x; \theta)$

We have a labelled dataset $\mathcal{D} = \{(\underline{x}_t, y_t) \in \mathbb{R}^d \times \{+1, -1\} : 1 \leq t \leq n\}$.

Perceptron Algorithm : 1) $k=0$ initialize $\underline{\Theta}^{(k)}$ arbitrarily
(for the proof, set $\underline{\Theta}^{(0)} = \underline{0}$)

2) cycle through the dataset

If $y_t < \underline{x}_t$, $\underline{\theta}^{(k)} > \leq 0$, then

Update $\underline{\Theta}^{(k+1)} = \underline{\Theta}^{(k)} + y_t \underline{x}_t$

If $y_t < x_t, \underline{\Theta}^{(k)} > 0$, do nothing!

We don't need to cycle through the whole dataset when we have new data!

3) Out!

The disadvantage of Perception \Rightarrow only find one line to separate the 2 class!

- ① cannot find a 'good' line!
 - ② not unique \rightarrow depends on the sequence of data

Assumptions: ① (R-bounded)

finite dataset

↳ infinite (countable) dataset

② γ -linearly separable $\Rightarrow \exists \underline{\theta}^* \in \mathbb{R}^d, \exists \gamma > 0$

$$y_t < \underline{\theta}^* \cdot \underline{x}_t > \geq \gamma \quad \forall t=1,2,\dots,n$$

$\rightarrow \gamma$ is defined by

$$\gamma = \min_{t \in T} y_t < \underline{\theta}^* \cdot \underline{x}_t >$$

Thm: [Minkoff (1962)] start the perceptron alg. at $\underline{\theta}^{(0)} = \underline{0}$

then maximum # of update required

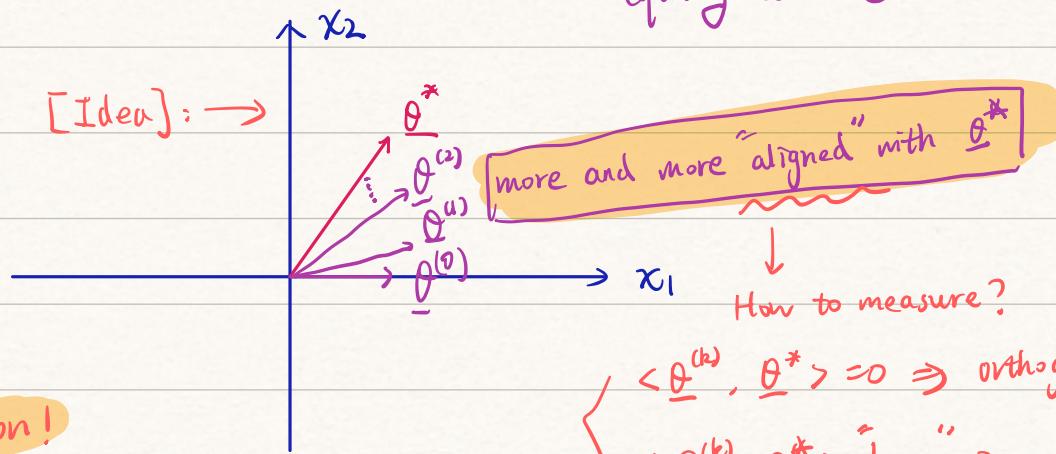
s.t. $\{\underline{x}_t\}_{t=1}^n$ are classified correctly k.

satisfies $k \leq R^2 \cdot \frac{\|\underline{\theta}^*\|^2}{\gamma^2}$

Not so good!

→ Many choice of $\underline{\theta}^*$

equality can be achieved?



Motivation!

中文意义 增长很快(关于k)

$$\frac{<\underline{\theta}^{(k)}, \underline{\theta}^*>}{\|\underline{\theta}^{(k)}\| \|\underline{\theta}^*\|} \leq A$$

增长比较慢(关于k)

$\Rightarrow k$ 不能太大!

- 1) $<\underline{\theta}^{(k)}, \underline{\theta}^*> = 0 \Rightarrow$ orthogonal
2) $<\underline{\theta}^{(k)}, \underline{\theta}^*> \text{ "large" } \Rightarrow$ more aligned!
3) this "large" require to control $\|\underline{\theta}^{(k)}\|$

$$\frac{<\underline{\theta}^{(k)}, \underline{\theta}^*>}{\|\underline{\theta}^{(k)}\| \|\underline{\theta}^*\|}$$

want to show:

- 1) $<\underline{\theta}^{(k)}, \underline{\theta}^*>$ increases at least linearly in k ✓
2) $\|\underline{\theta}^{(k)}\|^2$ increases at most linearly fast in k
3) Cauchy-Schwarz.

Pf: Part I:

update whenever make error

$$\begin{aligned}\Rightarrow \langle \underline{\theta}^*, \underline{\theta}^{(k)} \rangle &= \langle \underline{\theta}^*, \underline{\theta}^{(k-1)} + y_t \underline{x}_t \rangle \\ &= \langle \underline{\theta}^*, \underline{\theta}^{(k-1)} \rangle + y_t \underbrace{\langle \underline{\theta}^*, \underline{x}_t \rangle}_{\text{unnormalized alignment}} \\ &\geq \langle \underline{\theta}^*, \underline{\theta}^{(k-1)} \rangle + r.\end{aligned}$$

$$\text{since } \underline{\theta}^{(0)} = \underline{\theta}, \quad \langle \underline{\theta}^*, \underline{\theta}^{(k)} \rangle \geq k\delta.$$

\uparrow
"unnormalized alignment"

Part II:

$$\begin{aligned}\|\underline{\theta}^{(k)}\|^2 &= \langle \underline{\theta}^{(k-1)} + y_t \underline{x}_t, \underline{\theta}^{(k-1)} + y_t \underline{x}_t \rangle \\ &= \|\underline{\theta}^{(k-1)}\|^2 + \|\underline{x}_t\|^2 + 2y_t \langle \underline{x}_t, \underline{\theta}^{(k-1)} \rangle \\ &\leq \|\underline{\theta}^{(k-1)}\|^2 + R^2.\end{aligned}$$

(D is R-bounded) + (Update Property)

$$\leq kR^2$$

Part III

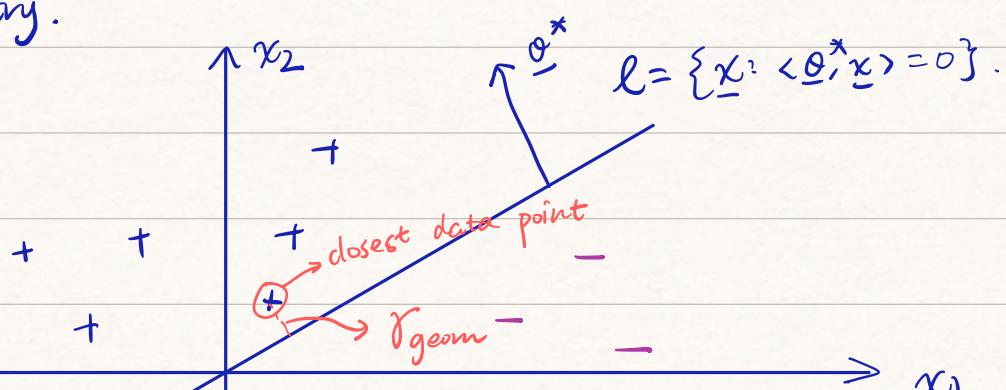
{ PART I
PART II.

$$\frac{\langle \underline{\theta}^{(k)}, \underline{\theta}^* \rangle}{\|\underline{\theta}^{(k)}\| \|\underline{\theta}^*\|} \leq 1 \Rightarrow 1 \geq \frac{\langle \underline{\theta}^{(k)}, \underline{\theta}^* \rangle}{\|\underline{\theta}^{(k)}\| \|\underline{\theta}^*\|} \geq \frac{k\delta}{\sqrt{kR^2 \|\underline{\theta}^*\|}}$$

$$\Rightarrow k \leq R^2 \frac{\|\underline{\theta}^*\|^2}{\delta^2}$$

□.

Margin & Geometry.





Claim: $r_{\text{geom}} = \frac{\gamma}{\|\underline{\theta}^*\|}$ is the smallest distance from any training sample to the decision boundary ℓ .

$$\text{perception alg. } k \leq R^2 \frac{\|\underline{\theta}\|}{\gamma^2} \triangleq \frac{R^2}{r_{\text{geom}}^2}$$

It shows that: the closer those training samples are to the decision boundary, the worst the upper bound of k is!

(r_{geom} is smaller in this sense)

$$\text{Rmk: (If dataset is normalized)} \Rightarrow k \leq \frac{1}{r_{\text{geom}}^2}$$

$$\stackrel{\Downarrow}{R=1}$$

So r_{geom} represents the difficulty of classification problem!

Target: $\gamma \rightarrow$ Actually we just prove when \underline{x}_t attains γ .
its distance is $\frac{\gamma}{\|\underline{\theta}^*\|} = r_{\text{geom}}$

Pf: To calculate the distance of closest point \underline{x}_t to the decision boundary $\ell = \{x : \langle \underline{\theta}^*, \underline{x} \rangle = 0\}$.

Note that \underline{x}_t satisfies $y_t \langle \underline{x}_t, \underline{\theta}^* \rangle = \gamma$.

Method \rightarrow define a line segment $\underline{x}(s)$ with $\underline{x}(0) = \underline{x}_t$

parallel to $\underline{\theta}^*$ towards the boundary.

$$\boxed{\underline{x}(s) = \underline{x}(0) - s \cdot \underline{y}_t \cdot \frac{\underline{\theta}^*}{\|\underline{\theta}^*\|}} \Rightarrow \langle \underline{x}(s), \underline{\theta}^* \rangle = 0 \text{ to solve } s.$$

\downarrow

s : length of segment from $\underline{x}(s)$ to $\underline{x}(0)$

$$\langle \underline{x}_t - s \cdot \underline{y}_t \cdot \frac{\underline{\theta}^*}{\|\underline{\theta}^*\|}, \underline{\theta}^* \rangle = 0$$

$$\rightarrow \langle \underline{x}_t, \underline{\theta}^* \rangle = s \cdot \underline{y}_t \cdot \|\underline{\theta}^*\|$$

$$\rightarrow s = \frac{\underline{y}_t \langle \underline{x}_t, \underline{\theta}^* \rangle}{\|\underline{\theta}^*\|} = \frac{\gamma}{\|\underline{\theta}^*\|}$$

Another Pf:

$$\begin{aligned}s &= \left| \langle \underline{x}_t, \frac{\underline{\theta}^*}{\|\underline{\theta}^*\|} \rangle \right| \\&= \frac{1}{\|\underline{\theta}^*\|} \left| \langle \underline{x}_t, \underline{\theta}^* \rangle \right| \\&= \frac{1}{\|\underline{\theta}^*\|} \cdot y_t \langle \underline{x}_t, \underline{\theta}^* \rangle \\&= \frac{y_t}{\|\underline{\theta}^*\|}\end{aligned}$$

$\|\underline{\theta}\|$ the distance of \underline{x}_t to $\underline{\theta}$.

for the definition

$$\underline{x}(t) = \underline{x}(0) - s \cdot \underline{y}_t \frac{\underline{\theta}^*}{\|\underline{\theta}^*\|}$$

to guarantee that $s > 0$.

Rmk: i). Bound on # of updates of perception alg.

$$\text{is } K \leq \left(\frac{R}{r_{\text{geom}}} \right)^2 - (*)$$

ii) $\frac{R}{r_{\text{geom}}}$ can be viewed as the difficulty of learning linear classification

→ related to VC dimension

$$\text{iii) } \underline{x}_t \in \mathbb{R}^d$$

The bound (*) is almost independent of d .

$$R = \max_i \|\underline{x}_i\|$$

increase with dimension d !

(R gets larger)

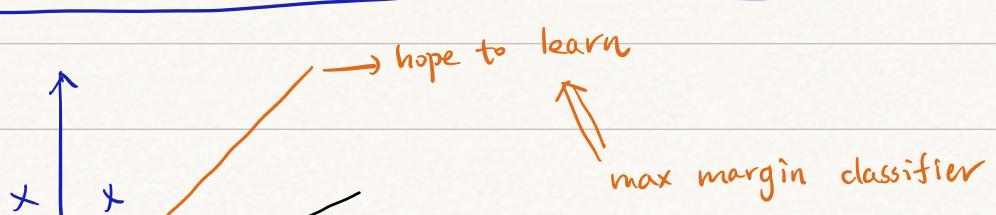
iv) the bound (*) is almost independent of n .

$$\text{Recall: } r = \min_{1 \leq t \leq n} y_t \langle \underline{x}_t, \underline{\theta}^* \rangle . \downarrow$$

when $n \uparrow$

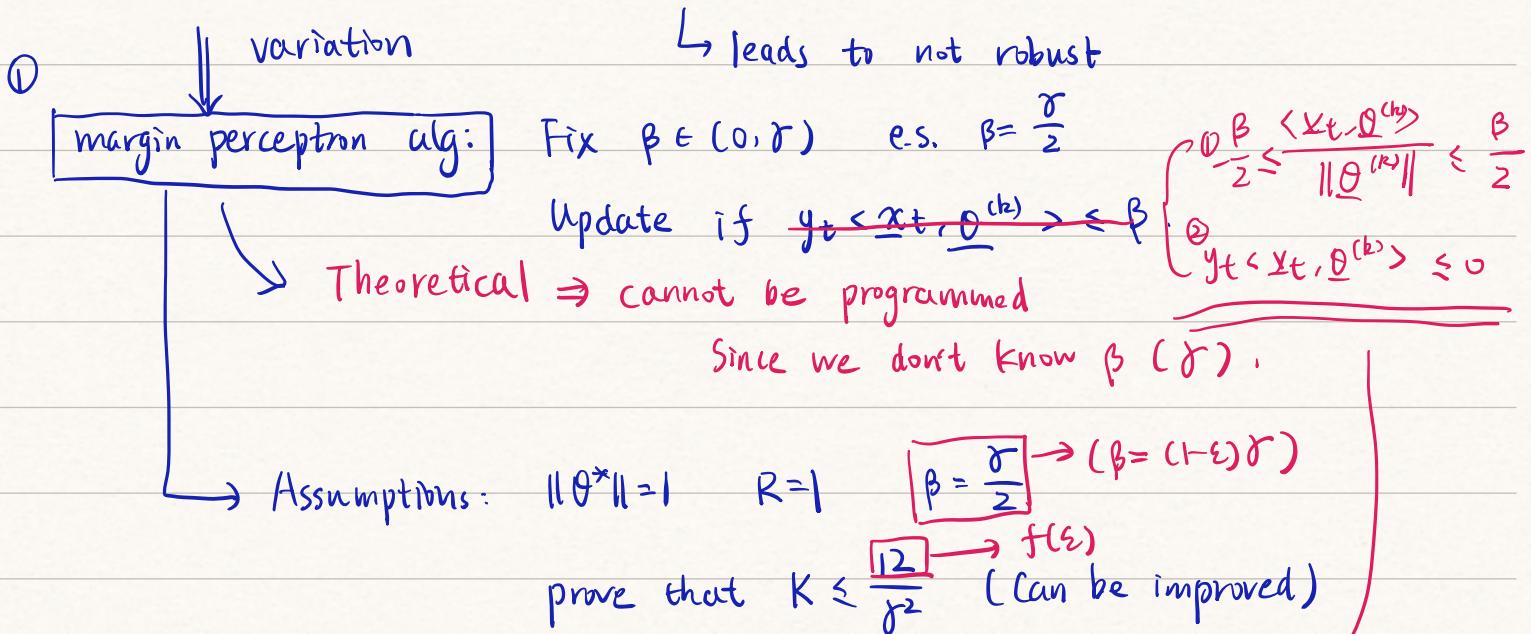
$$\Rightarrow r_{\text{geom}} \downarrow$$

Variants





Standard perceptron : Update if $y_t \langle \underline{x}_t, \underline{\theta}^{(k)} \rangle \leq 0 \Leftrightarrow$ mislabel



② Batch Perception Alg.

Update: $Y_K := \{ t=1, 2, \dots, n : y_t \langle \underline{x}_t, \underline{\theta}^{(k)} \rangle \leq 0 \}$

$$\underline{\theta}^{(k+1)} = \underline{\theta}^{(k)} + \eta_k \sum_{t \in Y_K} y_t \underline{x}_t \quad - (*)$$

Trajectory of $\{\underline{\theta}^{(k)}\}_{k \in \mathbb{N}}$ is "smoother". (In some sense)

Run (*) until $Y_K = \emptyset$.

E.s. $\eta_k = \frac{1}{K}$.

Thm: (i) $\eta_k \geq 0$, (ii) $\lim_{m \rightarrow \infty} \sum_{k=1}^m \eta_k = \infty$, (iii) $\lim_{m \rightarrow \infty} \frac{\sum_{k=1}^m \eta_k^2}{\left(\sum_{k=1}^m \eta_k\right)^2} = 0$,

then Batch Perceptron Alg. converges!

Outline:

① Perceptron (linear classifier)

$$\text{PCT} \rightarrow R = R^2 \cdot \frac{\|\theta^*\|^2}{\gamma^2} := R^2 \cdot \frac{1}{r_{\text{geom}}} \quad \text{and} \quad r_{\text{geom}} = \frac{\gamma}{\|\theta^*\|}$$

Geometric Meaning: the
closest distance from every
dataset to the decision Boundary.

② Variants

{ Margin Perceptron Alg.
Batch Perceptron Alg.

i) Margin Perceptron [change the Update]

{ Perceptron: Update if $y_t \langle x_t, \underline{\theta}^{(k)} \rangle \leq 0$,
M-Perceptron: Update if $y_t \langle x_t, \underline{\theta}^{(k)} \rangle \leq \beta(\gamma) \|\underline{\theta}^{(k)}\|$

A[thm] like PCT

we are UNKNOWN

误分类+分类不清 theoretically!

ii) Batch Perceptron Alg.



① when we have $\underline{\theta}^{(k)}$: find and set I_k .

$$I_k = \{t : y_t < \underline{x}_t \cdot \underline{\theta}^{(k)} \leq 0\}$$

② Update: $\underline{\theta}^{(k+1)} = \underline{\theta}^{(k)} + \eta_k \sum_{i \in I_k} y_i \underline{x}_i$