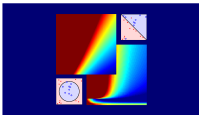


Machine Learning Foundations

(機器學習基石)



Lecture 2: Learning to Answer Yes/No

Hsuan-Tien Lin (林軒田)

htlin@csie.ntu.edu.tw

Department of Computer Science
& Information Engineering

National Taiwan University
(國立台灣大學資訊工程系)



Roadmap

1 When Can Machines Learn?

Lecture 1: The Learning Problem

\mathcal{A} takes \mathcal{D} and \mathcal{H} to get g

Lecture 2: Learning to Answer Yes/No

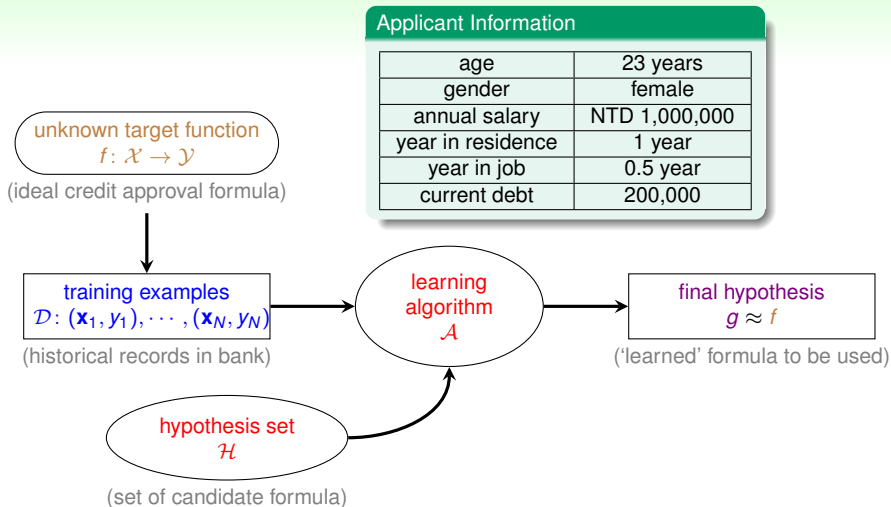
- Perceptron Hypothesis Set
- Perceptron Learning Algorithm (PLA)
- Guarantee of PLA

2 Why Can Machines Learn?

3 How Can Machines Learn?

4 How Can Machines Learn Better?

Credit Approval Problem Revisited



what hypothesis set can we use?

A Simple Hypothesis Set: the ‘Perceptron’

age	23 years
annual salary	NTD 1,000,000
year in job	0.5 year
current debt	200,000

- For $\mathbf{x} = (x_1, x_2, \dots, x_d)$ ‘features of customer’, compute a weighted ‘score’ and

approve credit if $\sum_{i=1}^d w_i x_i > \text{threshold}$

deny credit if $\sum_{i=1}^d w_i x_i < \text{threshold}$

- \mathcal{Y} : $\{+1(\text{good}), -1(\text{bad})\}$, 0 ignored—linear formula $h \in \mathcal{H}$ are

$$h(\mathbf{x}) = \text{sign} \left(\left(\sum_{i=1}^d w_i x_i \right) - \text{threshold} \right)$$

called ‘perceptron’ hypothesis historically

Vector Form of Perceptron Hypothesis

$$\begin{aligned}
 h(\mathbf{x}) &= \text{sign} \left(\left(\sum_{i=1}^d \mathbf{w}_i x_i \right) - \text{threshold} \right) \\
 &= \text{sign} \left(\left(\sum_{i=1}^d \mathbf{w}_i x_i \right) + \underbrace{(-\text{threshold})}_{w_0} \cdot \underbrace{(+1)}_{x_0} \right) \\
 &= \text{sign} \left(\sum_{i=0}^d \mathbf{w}_i x_i \right) \\
 &= \text{sign} \left(\mathbf{w}^T \mathbf{x} \right)
 \end{aligned}$$

- each ‘taller’ \mathbf{w} represents a hypothesis h & is multiplied with ‘taller’ \mathbf{x} — will use taller versions to simplify notation

what do perceptrons h ‘look like’?

Some Notation Conventions

fonts

- normal x : just a scalar
- bold \mathbf{x} : a vector (x_0, x_1, \dots, x_d)
bold \mathbf{w} : a vector (w_0, w_1, \dots, w_d)
- normal x_i : the i -th component in \mathbf{x}
normal w_i : the i -th component in \mathbf{w}
- bold \mathbf{x}_n : the n -th vector (in the data)
bold \mathbf{w}_t : the t -th vector (we will see)
- normal $x_{n,i}$ (rarely used): the i -th component in \mathbf{x}_n
normal $w_{t,i}$ (rarely used): the i -th component in \mathbf{w}_t
- caligraphic as sets: input \mathcal{X} , output \mathcal{Y} , data \mathcal{D} , hypothesis \mathcal{H} , except algorithm \mathcal{A}

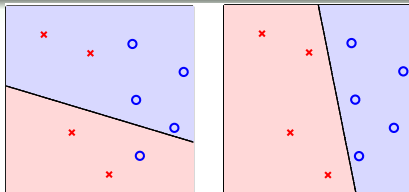
two important numbers

- N examples (\mathbf{x}_n, y_n) , indexed by $n = 1, 2, \dots, N$
- d features, indexed by $i = 0, 1, 2, \dots, d$

important to follow the notations
from the very beginning

Perceptrons in \mathbb{R}^2

$$h(\mathbf{x}) = \text{sign}(w_0 + w_1 x_1 + w_2 x_2)$$



- customer features \mathbf{x} : points on the plane (or points in \mathbb{R}^d)
- labels y : $\circ (+1)$, $\times (-1)$
- hypothesis h : visually **lines** $w_0 + w_1 x_1 + w_2 x_2 = 0$
(or hyperplanes in \mathbb{R}^d)
—**positive** on one side of a line, **negative** on the other side
- different line classifies customers differently

perceptrons \Leftrightarrow **linear (binary) classifiers**

Fun Time

Consider using a perceptron to detect spam messages.

Assume that each email is represented by the frequency of keyword occurrence, and output $+1$ indicates a spam. Which keywords below shall have large positive weights in a **good perceptron** for the task?

- ① coffee, tea, hamburger, steak
- ② free, drug, fantastic, deal
- ③ machine, learning, statistics, textbook
- ④ national, Taiwan, university, coursera

Fun Time

Consider using a perceptron to detect spam messages.

Assume that each email is represented by the frequency of keyword occurrence, and output $+1$ indicates a spam. Which keywords below shall have large positive weights in a **good perceptron** for the task?

- ① coffee, tea, hamburger, steak
- ② free, drug, fantastic, deal
- ③ machine, learning, statistics, textbook
- ④ national, Taiwan, university, coursera

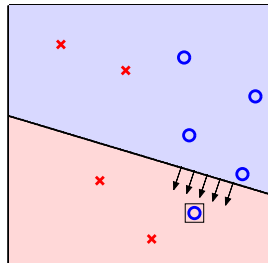
Reference Answer: ②

The occurrence of keywords with positive weights increase the 'spam score', and hence those keywords should often appear in spams.

Select g from \mathcal{H}

\mathcal{H} = all possible perceptrons, $g = ?$

- want: $g \approx f$ (hard when f unknown)
- almost necessary: $g \approx f$ on \mathcal{D} , ideally $g(\mathbf{x}_n) = f(\mathbf{x}_n) = y_n$
- difficult: \mathcal{H} is of **infinite** size
- idea: start from some g_0 , and ‘correct’ its mistakes on \mathcal{D}



will represent g_0 by its weight vector \mathbf{w}_0

Perceptron Learning Algorithm

start from some \mathbf{w}_0 (say, $\mathbf{0}$), and 'correct' its mistakes on \mathcal{D}

For $t = 0, 1, \dots$

- 1 find a **mistake** of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$

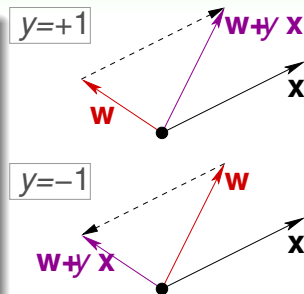
$$\text{sign} \left(\mathbf{w}_t^T \mathbf{x}_{n(t)} \right) \neq y_{n(t)}$$

- 2 (try to) correct the mistake by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}$$

... until **no more mistakes**

return **last \mathbf{w}** (called \mathbf{w}_{pla}) as g



That's it!

—A fault confessed is half redressed. :-)

Handling $\text{sign}(\cdot) = 0$

Perceptron Learning Algorithm

start from some \mathbf{w}_0 (say, $\mathbf{0}$), and 'correct' its mistakes on \mathcal{D}

When $\mathbf{w}_0 = \mathbf{0}$, technically $\text{sign}(\mathbf{w}_0^T \mathbf{x}_{n(0)}) = 0$, shall we update?

- convention -1: $\text{sign}(0) = -1$ (update if $y_{n(0)} = +1$)
- convention +1: $\text{sign}(0) = +1$ (update if $y_{n(0)} = -1$)
- convention 0: $\text{sign}(0) = 0$ (always update)
- convention r: $\text{sign}(0) = \text{random flip}$ (50% chance of update)

—usually does not matter much, as long as \mathbf{w}_1 often becomes non-zero

$\mathbf{w}_t^T \mathbf{x}_{n(t)} = 0$ rarely happens in practice

How to Update w_0 ?

Perceptron Learning Algorithm

For $t = 0, 1, \dots$

- ① find a **mistake** of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$: $\text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) \neq y_{n(t)}$
- ② (try to) correct the mistake by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}, \text{ i.e.,}$$

$$\begin{bmatrix} w_{t+1,0} \\ w_{t+1,1} \\ \dots \\ w_{t+1,d} \end{bmatrix} = \begin{bmatrix} w_{t,0} \\ w_{t,1} \\ \dots \\ w_{t,d} \end{bmatrix} + y_{n(t)} \begin{bmatrix} x_0 (= \text{what?}) \\ x_{n(t),1} \\ \dots \\ x_{n(t),d} \end{bmatrix}$$

... until **no more mistakes**

return **last \mathbf{w}** (called \mathbf{w}_{pla}) as g

each update changes $w_{t,0}$ by $y_{n(t)}$

Practical Implementation of PLA

start from some \mathbf{w}_0 (say, $\mathbf{0}$), and 'correct' its mistakes on \mathcal{D}

Cyclic PLA

For $t = 0, 1, \dots$

- 1 find the next mistake of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$

$$\text{sign} \left(\mathbf{w}_t^T \mathbf{x}_{n(t)} \right) \neq y_{n(t)}$$

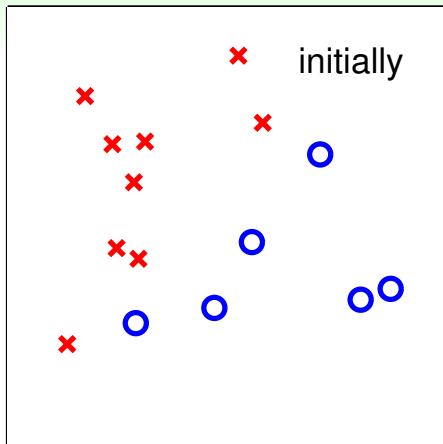
- 2 correct the mistake by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}$$

... until a full cycle of not encountering mistakes

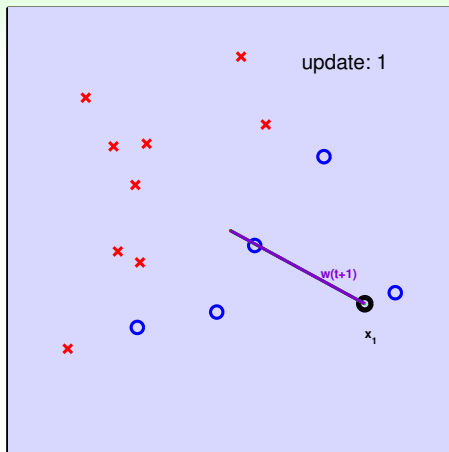
next can follow naïve cycle $(1, \dots, N)$
or precomputed random cycle

Seeing is Believing



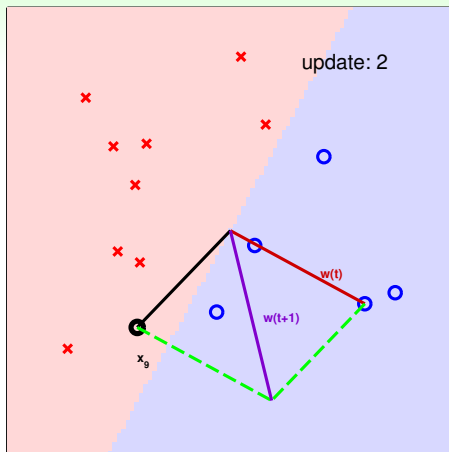
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



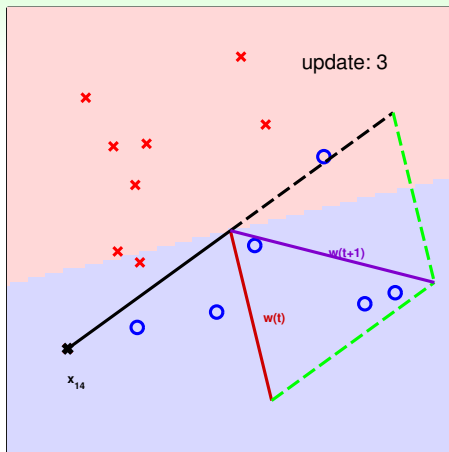
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



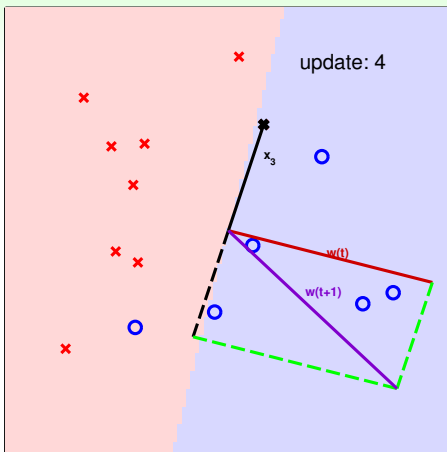
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



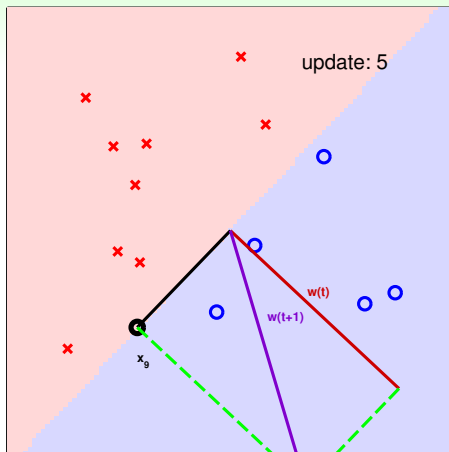
worked like a charm with < 20 lines!!
 (note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



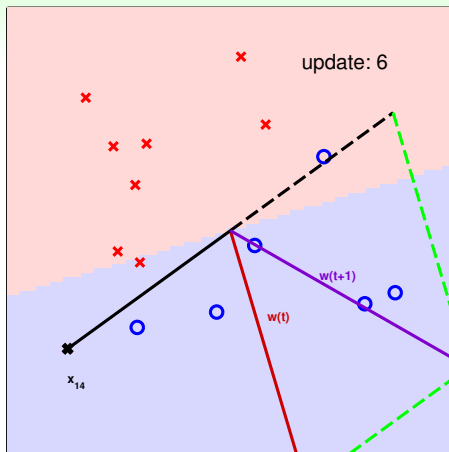
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



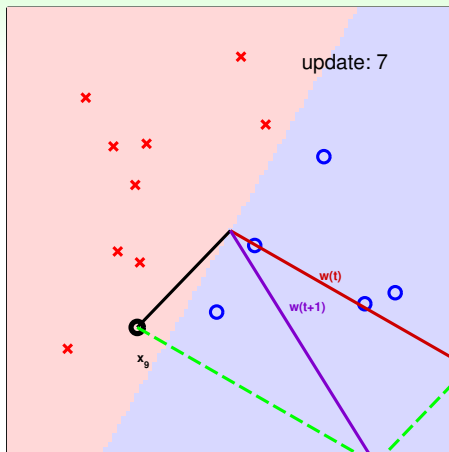
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



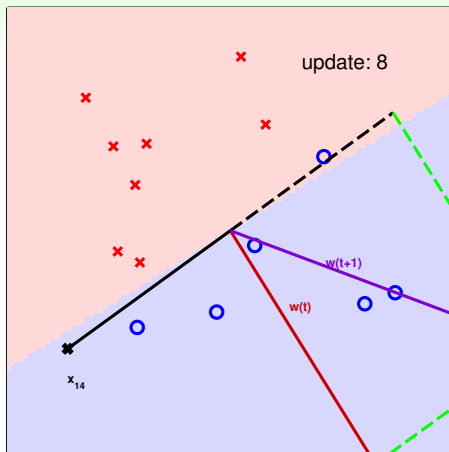
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



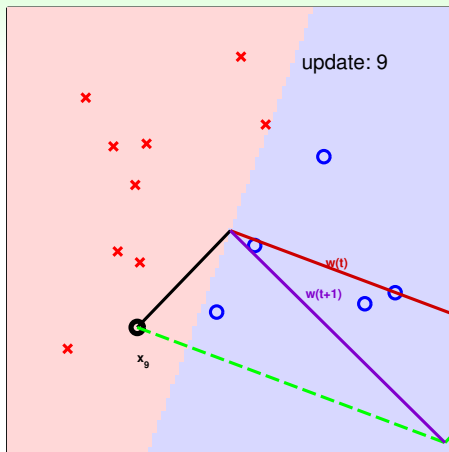
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



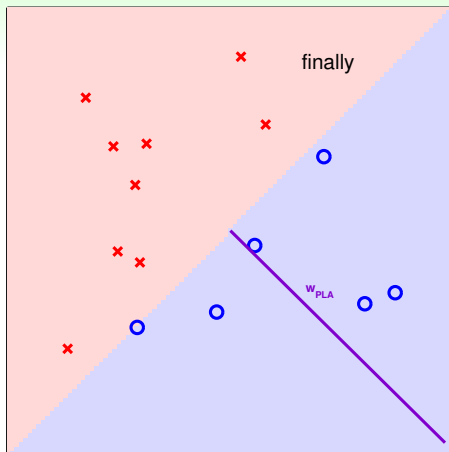
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Fun Time

Let's try to think about why PLA may work.

Let $n = n(t)$, according to the rule of PLA below, which formula is true?

$$\text{sign}(\mathbf{w}_t^T \mathbf{x}_n) \neq y_n, \quad \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_n \mathbf{x}_n$$

- ① $\mathbf{w}_{t+1}^T \mathbf{x}_n = y_n$
- ② $\text{sign}(\mathbf{w}_{t+1}^T \mathbf{x}_n) = y_n$
- ③ $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n \geq y_n \mathbf{w}_t^T \mathbf{x}_n$
- ④ $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n < y_n \mathbf{w}_t^T \mathbf{x}_n$

Fun Time

Let's try to think about why PLA may work.

Let $n = n(t)$, according to the rule of PLA below, which formula is true?

$$\text{sign}(\mathbf{w}_t^T \mathbf{x}_n) \neq y_n, \quad \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_n \mathbf{x}_n$$

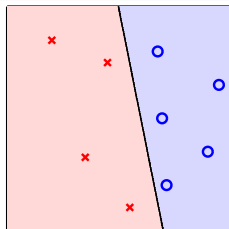
- ① $\mathbf{w}_{t+1}^T \mathbf{x}_n = y_n$
- ② $\text{sign}(\mathbf{w}_{t+1}^T \mathbf{x}_n) = y_n$
- ③ $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n \geq y_n \mathbf{w}_t^T \mathbf{x}_n$
- ④ $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n < y_n \mathbf{w}_t^T \mathbf{x}_n$

Reference Answer: ③

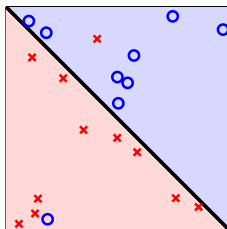
Simply multiply the second part of the rule by $y_n \mathbf{x}_n$. The result shows that the rule somewhat 'tries to correct the mistake.'

Linear Separability

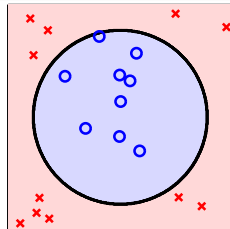
- if PLA halts (i.e. no more mistakes),
(necessary condition) \mathcal{D} allows some \mathbf{w} to make no mistake
- call such \mathcal{D} linear separable



(linear separable)



(not linear separable)



(not linear separable)

assume linear separable \mathcal{D} ,
does PLA always halt?

PLA Fact: \mathbf{w}_t Gets More Aligned with \mathbf{w}_f

linear separable $\mathcal{D} \Leftrightarrow$ exists perfect \mathbf{w}_f such that $y_n = \text{sign}(\mathbf{w}_f^T \mathbf{x}_n)$

- \mathbf{w}_f perfect hence every \mathbf{x}_n correctly away from line:

$$y_{n(t)} \mathbf{w}_f^T \mathbf{x}_{n(t)} \geq \min_n y_n \mathbf{w}_f^T \mathbf{x}_n > 0$$

- $\mathbf{w}_f^T \mathbf{w}_t \uparrow$ by updating with any $(\mathbf{x}_{n(t)}, y_{n(t)})$

$$\begin{aligned} \mathbf{w}_f^T \mathbf{w}_{t+1} &= \mathbf{w}_f^T (\mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}) \\ &\geq \mathbf{w}_f^T \mathbf{w}_t + \min_n y_n \mathbf{w}_f^T \mathbf{x}_n \\ &> \mathbf{w}_f^T \mathbf{w}_t + 0. \end{aligned}$$

\mathbf{w}_t appears more aligned with \mathbf{w}_f after update
(really?)

PLA Fact: \mathbf{w}_t Does Not Grow Too Fast

\mathbf{w}_t changed only when mistake

$$\Leftrightarrow \text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) \neq y_{n(t)} \Leftrightarrow y_{n(t)} \mathbf{w}_t^T \mathbf{x}_{n(t)} \leq 0$$

- mistake 'limits' $\|\mathbf{w}_t\|^2$ growth, even when updating with 'longest' \mathbf{x}_n

$$\begin{aligned} \|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}\|^2 \\ &= \|\mathbf{w}_t\|^2 + 2y_{n(t)} \mathbf{w}_t^T \mathbf{x}_{n(t)} + \|y_{n(t)} \mathbf{x}_{n(t)}\|^2 \\ &\leq \|\mathbf{w}_t\|^2 + 0 + \|y_{n(t)} \mathbf{x}_{n(t)}\|^2 \\ &\leq \|\mathbf{w}_t\|^2 + \max_n \|y_n \mathbf{x}_n\|^2 \end{aligned}$$

start from $\mathbf{w}_0 = \mathbf{0}$, after T mistake corrections,

$$\frac{\mathbf{w}_f^T}{\|\mathbf{w}_f\|} \frac{\mathbf{w}_T}{\|\mathbf{w}_T\|} \geq \sqrt{T} \cdot \text{constant}$$

Fun Time

Let's upper-bound T , the number of mistakes that PLA 'corrects'.

$$\text{Define } R^2 = \max_n \|\mathbf{x}_n\|^2 \quad \rho = \min_n y_n \frac{\mathbf{w}_f^T \mathbf{x}_n}{\|\mathbf{w}_f\|}$$

We want to show that $T \leq \square$. Express the upper bound \square by the two terms above.

- 1 R/ρ
- 2 R^2/ρ^2
- 3 R/ρ^2
- 4 ρ^2/R^2

Fun Time

Let's upper-bound T , the number of mistakes that PLA 'corrects'.

$$\text{Define } R^2 = \max_n \|\mathbf{x}_n\|^2 \quad \rho = \min_n y_n \frac{\mathbf{w}_f^T \mathbf{x}_n}{\|\mathbf{w}_f\|}$$

We want to show that $T \leq \square$. Express the upper bound \square by the two terms above.

- ① R/ρ
- ② R^2/ρ^2
- ③ R/ρ^2
- ④ ρ^2/R^2

Reference Answer: ②

The maximum value of $\frac{\mathbf{w}_f^T \mathbf{w}_t}{\|\mathbf{w}_f\| \|\mathbf{w}_t\|}$ is 1. Since T mistake corrections **increase the inner product by \sqrt{T} · constant**, the maximum number of corrected mistakes is $1/\text{constant}^2$.

PLA Mistake Bound

inner product grows **fast**

$$\mathbf{w}_f^T \mathbf{w}_{t+1} \geq \mathbf{w}_f^T \mathbf{w}_t + \underbrace{\min_n y_n \mathbf{w}_f^T \mathbf{x}_n}_{\rho \cdot \|\mathbf{w}_f\|}$$

length² grows **slowly**

$$\|\mathbf{w}_{t+1}\|^2 \leq \|\mathbf{w}_t\|^2 + \underbrace{\max_n \|\mathbf{x}_n\|^2}_{R^2}$$

Magic Chain!

$$\mathbf{w}_f^T \mathbf{w}_1 \geq \mathbf{w}_f^T \mathbf{w}_0 + \rho \cdot \|\mathbf{w}_f\|$$

$$\mathbf{w}_f^T \mathbf{w}_2 \geq \mathbf{w}_f^T \mathbf{w}_1 + \rho \cdot \|\mathbf{w}_f\|$$

$$\mathbf{w}_f^T \mathbf{w}_3 \geq \mathbf{w}_f^T \mathbf{w}_2 + \rho \cdot \|\mathbf{w}_f\|$$

...

$$\mathbf{w}_f^T \mathbf{w}_T \geq \mathbf{w}_f^T \mathbf{w}_{T-1} + \rho \cdot \|\mathbf{w}_f\|$$

Magic Chain!

$$\|\mathbf{w}_1\|^2 \leq \|\mathbf{w}_0\|^2 + R^2$$

$$\|\mathbf{w}_2\|^2 \leq \|\mathbf{w}_1\|^2 + R^2$$

$$\|\mathbf{w}_3\|^2 \leq \|\mathbf{w}_2\|^2 + R^2$$

...

$$\|\mathbf{w}_T\|^2 \leq \|\mathbf{w}_{T-1}\|^2 + R^2$$

start from $\mathbf{w}_0 = \mathbf{0}$, after T mistake corrections,

$$1 \geq \frac{\mathbf{w}_f^T}{\|\mathbf{w}_f\|} \frac{\mathbf{w}_T}{\|\mathbf{w}_T\|} \geq \frac{T\rho\|\mathbf{w}_f\|}{\|\mathbf{w}_f\|\sqrt{TR}} \implies T \leq \left(\frac{R}{\rho}\right)^2$$

More about PLA

Guarantee

as long as **linear separable** and **correct by mistake**

- inner product of \mathbf{w}_f and \mathbf{w}_t grows fast; length of \mathbf{w}_t grows slowly
- PLA 'lines' are more and more aligned with $\mathbf{w}_f \Rightarrow$ halts

Pros

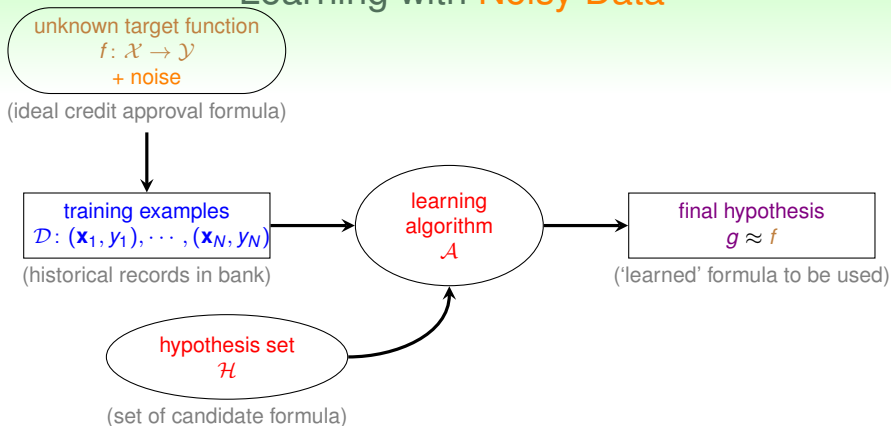
simple to implement, fast, works in any dimension d

Cons

- 'assumes' linear separable \mathcal{D} to halt
—property unknown in advance (no need for PLA if we know \mathbf{w}_f)
- not fully sure how long halting takes (ρ depends on \mathbf{w}_f)
—though practically fast

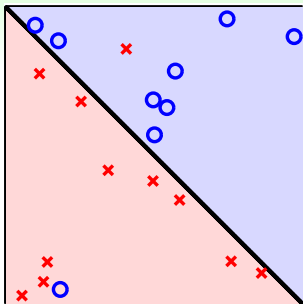
what if \mathcal{D} not linear separable?

Learning with Noisy Data



how to at least get $g \approx f$ on **noisy** \mathcal{D} ?

Line with Noise Tolerance



- assume ‘little’ noise: $y_n = f(\mathbf{x}_n)$ usually
- if so, $g \approx f$ on $\mathcal{D} \Leftrightarrow y_n = g(\mathbf{x}_n)$ usually
- how about

$$\mathbf{w}_g \leftarrow \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N \mathbb{I}[y_n \neq \operatorname{sign}(\mathbf{w}^T \mathbf{x}_n)]$$

—NP-hard to solve, unfortunately

will discuss other solutions for an
‘approximately good’ g later?

Summary

① When Can Machines Learn?

Lecture 1: The Learning Problem

Lecture 2: Learning to Answer Yes/No

- Perceptron Hypothesis Set
hyperplanes/linear classifiers in \mathbb{R}^d
- Perceptron Learning Algorithm (PLA)
correct mistakes and improve iteratively
- Guarantee of PLA
no mistake eventually if linear separable

• next: the zoo of learning problems

② Why Can Machines Learn?

③ How Can Machines Learn?

④ How Can Machines Learn Better?