

学习问题

刘 新 旺

Email: xinwangliu@nudt.edu.cn

国防科技大学 计算机学院
计算科学系 人工智能与大数据教研室

2019 年 10 月 15 日



- 1 学习与机器学习
- 2 机器学习组成要素及与其它领域的关系
- 3 感知机假说集及感知机学习算法
- 4 感知机学习算法的理论保证
- 5 非可分数据

从“学习”到“机器学习”

Learning: Acquiring **skill**

With experience accumulated from **observations**

observations → **learning** → **skill**

从“学习”到“机器学习”

Learning: Acquiring skill

With experience accumulated from observations

observations → **learning** → skill

Machine Learning (ML): Acquiring skill

With experience accumulated/**computed** from data

data → **machine learning** → skill

从“学习”到“机器学习”

Learning: Acquiring **skill**

With experience accumulated from **observations**

observations → **learning** → **skill**

Machine Learning (ML): Acquiring **skill**

With experience accumulated/**computed** from **data**

data → **machine learning** → **skill**

what is **skill** ?

skill

⇔ improve some performance measure (e.g., prediction accuracy)

skill

⇔ improve some performance measure (e.g., prediction accuracy)

Machine Learning: improve some performance measure

With experience **computed** from data

data → machine learning → improved performance measure

An application in Computational Finance

stock data → machine learning → more investment gain

- 1 学习与机器学习
- 2 机器学习组成要素及与其它领域的关系
- 3 感知机假说集及感知机学习算法
- 4 感知机学习算法的理论保证
- 5 非可分数据

机器学习的组成要素：以“信用卡申请”为例

申请者信息

age	23 years
gender	female
annual salary	1,000,000
year in residence	1 year
year in job	0.5 year
current debt	200,000

unknown pattern to be learned: “approve credit card good for bank?”

Basic Notations

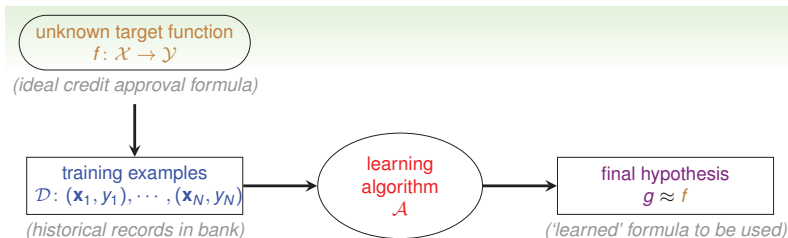
- input: $\mathbf{x} \in \mathcal{X}$ (customer application)
- output: $y \in \mathcal{Y}$ (good/bad after approving credit card)
- unknown pattern to be learned \Leftrightarrow target function:
 $f : \mathcal{X} \mapsto \mathcal{Y}$ (ideal credit approval formula)
- data \Leftrightarrow training examples: $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
(historical records in bank)
- hypothesis \Leftrightarrow skill with hopefully good performance:
 $g : \mathcal{X} \mapsto \mathcal{Y}$ (“learned” formula to be used)

Basic Notations

- input: $\mathbf{x} \in \mathcal{X}$ (customer application)
- output: $y \in \mathcal{Y}$ (good/bad after approving credit card)
- unknown pattern to be learned \Leftrightarrow target function:
 $f : \mathcal{X} \mapsto \mathcal{Y}$ (ideal credit approval formula)
- data \Leftrightarrow training examples: $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
(historical records in bank)
- hypothesis \Leftrightarrow skill with hopefully good performance:
 $g : \mathcal{X} \mapsto \mathcal{Y}$ (“learned” formula to be used)

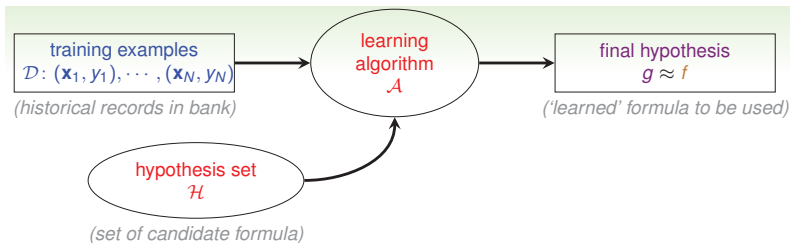
$\{(\mathbf{x}_n, y_n)\}$ from $f \longrightarrow$ **machine learning** $\longrightarrow g$

信用卡申请的学习流程图

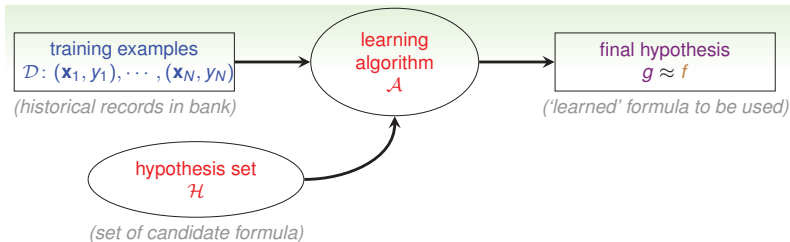


- target f **unknown**
(i.e. no programmable definition)
- hypothesis g hopefully $\approx f$
but possibly **different** from f
(perfection ‘impossible’ when f unknown)

What does g look like?



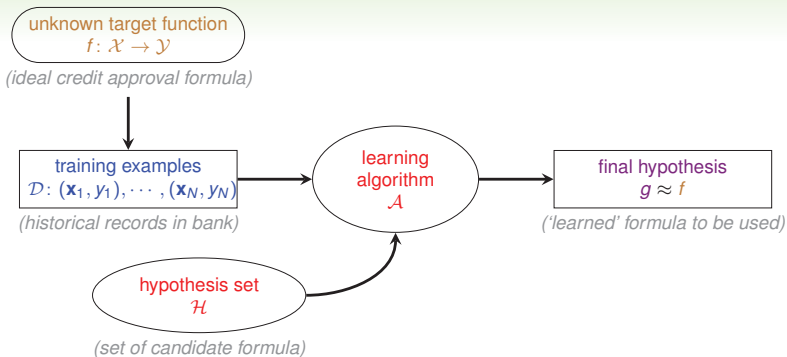
- assume $g \in \mathcal{H} = \{h_k\}$, i.e. approving if
 - h_1 : annual salary > NTD 800,000
 - h_2 : debt > NTD 100,000 (really?)
 - h_3 : year in job ≤ 2 (really?)
- hypothesis set \mathcal{H} :
 - can contain **good or bad hypotheses**
 - up to \mathcal{A} to pick the 'best' one as g



- assume $g \in \mathcal{H} = \{h_k\}$, i.e. approving if
 - h_1 : annual salary > NTD 800,000
 - h_2 : debt > NTD 100,000 (really?)
 - h_3 : year in job ≤ 2 (really?)
- hypothesis set \mathcal{H} :
 - can contain **good or bad hypotheses**
 - up to \mathcal{A} to pick the 'best' one as g

learning model = \mathcal{A} and \mathcal{H}

机器学习的实用定义



machine learning:
use **data** to compute **hypothesis** g
that approximates **target** f

How to use the four sets below to form a learning problem for song recommendation?

$$\mathcal{S}_1 = [0, 100]$$

$$\mathcal{S}_2 = \text{all possible (userid, songid) pairs}$$

$$\mathcal{S}_3 = \text{all formula that 'multiplies' user factors \& song factors, indexed by all possible combinations of such factors}$$

$$\mathcal{S}_4 = 1,000,000 \text{ pairs of ((userid, songid), rating)}$$

$$\textcircled{1} \mathcal{S}_1 = \mathcal{X}, \mathcal{S}_2 = \mathcal{Y}, \mathcal{S}_3 = \mathcal{H}, \mathcal{S}_4 = \mathcal{D}$$

$$\textcircled{2} \mathcal{S}_1 = \mathcal{Y}, \mathcal{S}_2 = \mathcal{X}, \mathcal{S}_3 = \mathcal{H}, \mathcal{S}_4 = \mathcal{D}$$

$$\textcircled{3} \mathcal{S}_1 = \mathcal{D}, \mathcal{S}_2 = \mathcal{H}, \mathcal{S}_3 = \mathcal{Y}, \mathcal{S}_4 = \mathcal{X}$$

$$\textcircled{4} \mathcal{S}_1 = \mathcal{X}, \mathcal{S}_2 = \mathcal{D}, \mathcal{S}_3 = \mathcal{Y}, \mathcal{S}_4 = \mathcal{H}$$

Machine Learning

use data to compute hypothesis g
that approximates target f

Data Mining

use **(huge)** data to **find property**
that is interesting

- if 'interesting property' **same as** 'hypothesis that approximate target'
— **ML = DM** (usually what KDDCup does)
- if 'interesting property' **related to** 'hypothesis that approximate target'
— **DM can help ML, and vice versa** (often, but not always)
- traditional DM also focuses on **efficient computation in large database**

difficult to distinguish ML and DM in reality

Machine Learning

use data to compute hypothesis g
that approximates target f

Artificial Intelligence

compute **something**
that shows intelligent behavior

- $g \approx f$ is something that shows intelligent behavior
— **ML can realize AI**, among other routes
- e.g. chess playing
 - traditional AI: game tree
 - ML for AI: 'learning from board data'

ML is one possible route to realize AI

Machine Learning

use data to compute hypothesis g
that approximates target f

Statistics

use data to **make inference**
about an unknown process

- g is an inference outcome; f is something unknown
—statistics **can be used to achieve ML**
- traditional statistics also focus on **provable results with math assumptions**, and care less about computation

statistics: many useful tools for ML

Which of the following claim is not totally true?

- ① machine learning is a route to realize artificial intelligence
- ② machine learning, data mining and statistics all need data
- ③ data mining is just another name for machine learning
- ④ statistics can be used for data mining

Which of the following claim is not totally true?

- ① machine learning is a route to realize artificial intelligence
- ② machine learning, data mining and statistics all need data
- ③ data mining is just another name for machine learning
- ④ statistics can be used for data mining

Reference Answer: ③

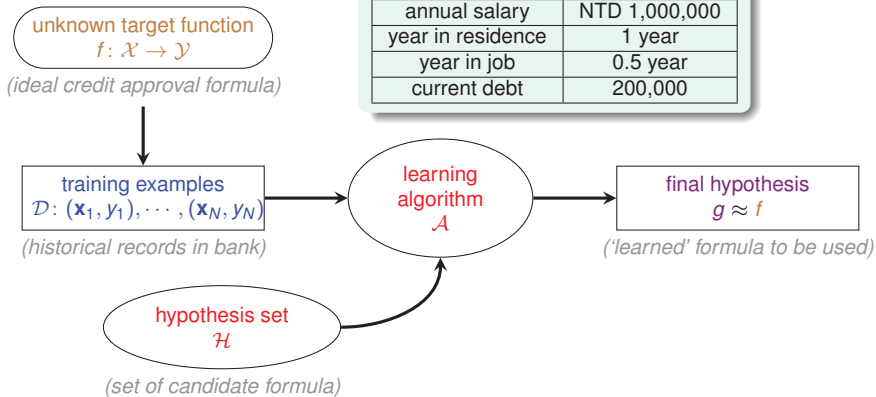
While data mining and machine learning do share a huge overlap, they are arguably not equivalent because of the difference of focus.

- 1 学习与机器学习
- 2 机器学习组成要素及与其它领域的关系
- 3 感知机假说集及感知机学习算法
- 4 感知机学习算法的理论保证
- 5 非可分数据

回顾信用卡申请案例

Applicant Information

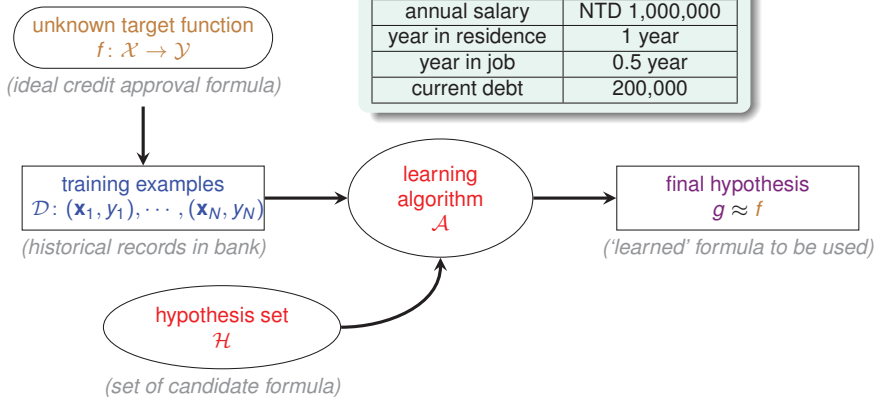
age	23 years
gender	female
annual salary	NTD 1,000,000
year in residence	1 year
year in job	0.5 year
current debt	200,000



回顾信用卡申请案例

Applicant Information

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



What hypothesis set can we use?

一类简单的假说集合：“感知机”

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

感知机假说的向量形式

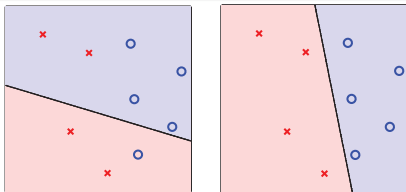
$$\begin{aligned}h(\mathbf{x}) &= \text{sign} \left(\left(\sum_{i=1}^d \mathbf{w}_i x_i \right) - \text{threshold} \right) \\&= \text{sign} \left(\left(\sum_{i=1}^d \mathbf{w}_i x_i \right) + \underbrace{(-\text{threshold})}_{\mathbf{w}_0} \cdot \underbrace{(+1)}_{x_0} \right) \\&= \text{sign} \left(\sum_{i=0}^d \mathbf{w}_i x_i \right) \\&= \text{sign} \left(\mathbf{w}^T \mathbf{x} \right)\end{aligned}$$

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

what do perceptrons h “look like”?

\mathbb{R}^2 空间中的感知机

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

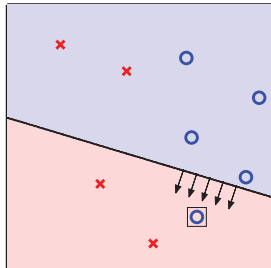


- customer features \mathbf{x} : points on the plane (or points in \mathbb{R}^d)
- labels y : $\circ (+1)$, $\times (-1)$
- hypothesis h : **lines** (or hyperplanes in \mathbb{R}^d)
—positive on one side of a line, negative on the other side
- different line classifies customers differently

perceptrons \Leftrightarrow linear (binary) classifiers

\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

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 $(\mathbf{x}_{n(t)}, y_{n(t)})$

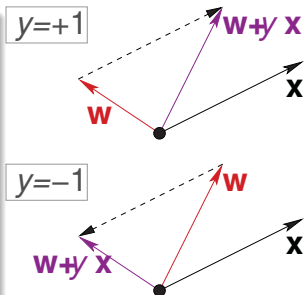
$$\text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) \neq y_{n(t)}$$

- 2 (try to) correct the mistake by

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

... until **no more mistakes**

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



感知机学习算法的实用实现

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 $(\mathbf{x}_{n(t)}, y_{n(t)})$

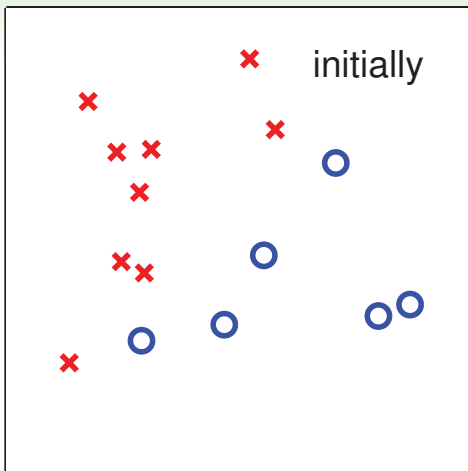
$$\text{sign} \left(\mathbf{w}_t^T \mathbf{x}_{n(t)} \right) \neq y_{n(t)}$$

- 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**

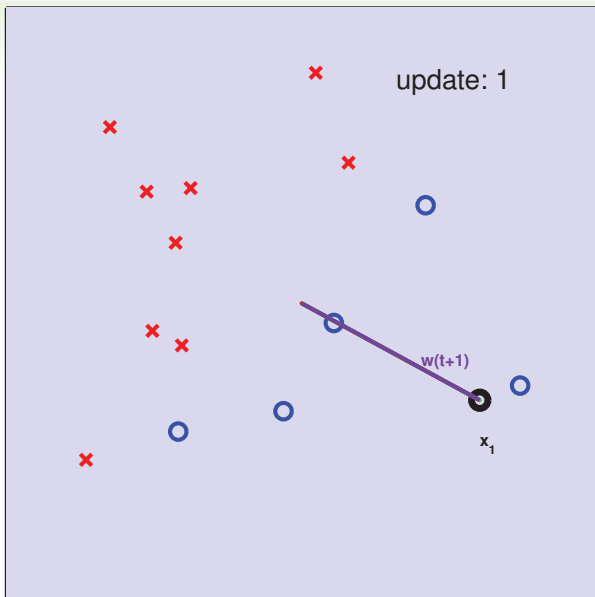
Seeing is Believing



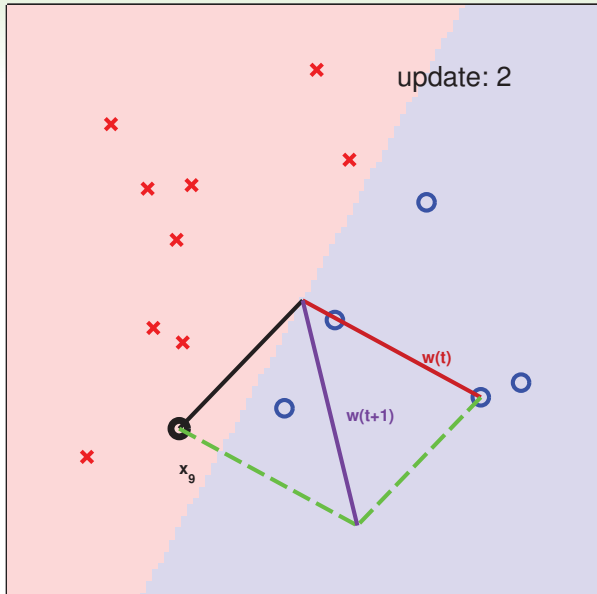
worked like a charm with < 20 lines!!

(note: made $\mathbf{x}_i \gg \mathbf{x}_0 = 1$ for visual purpose)

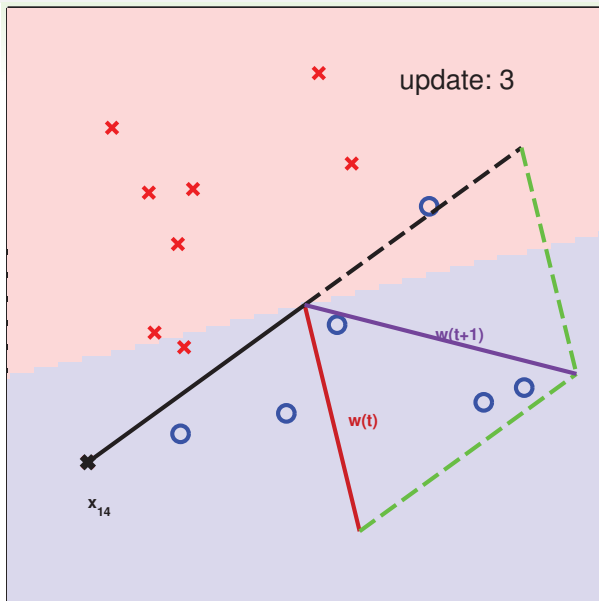
Seeing is Believing



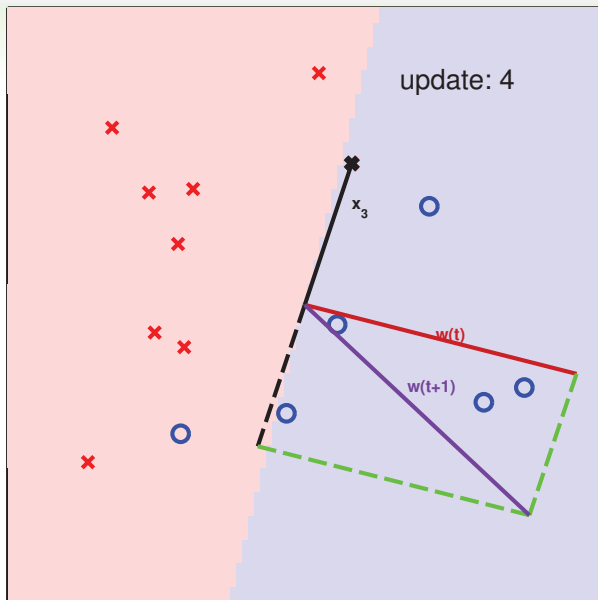
Seeing is Believing



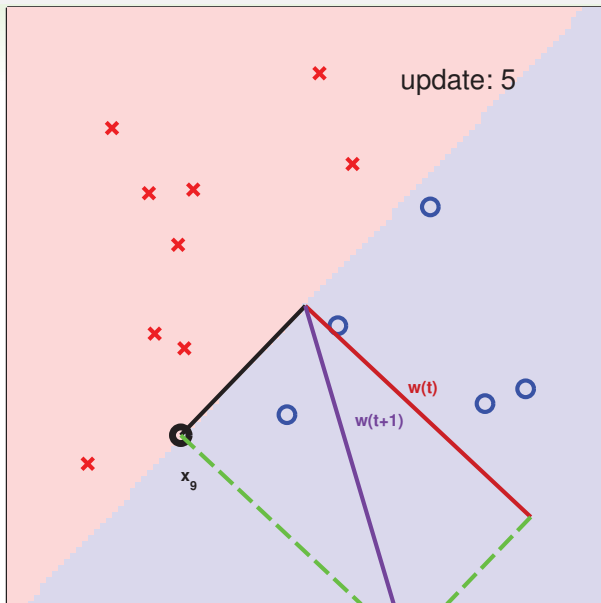
Seeing is Believing



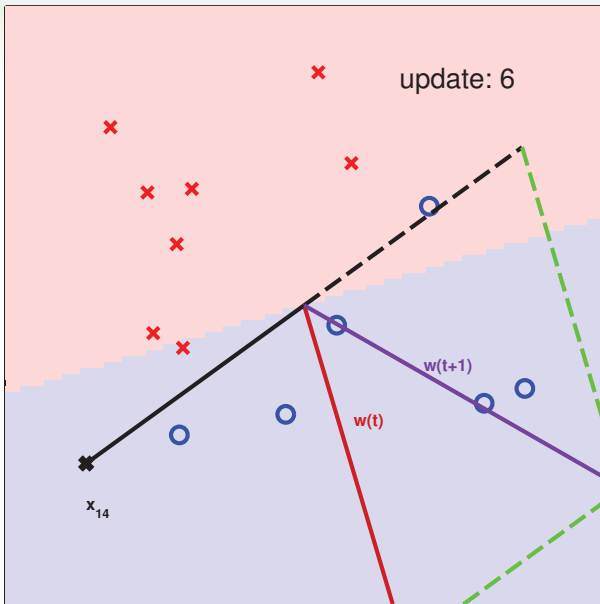
Seeing is Believing



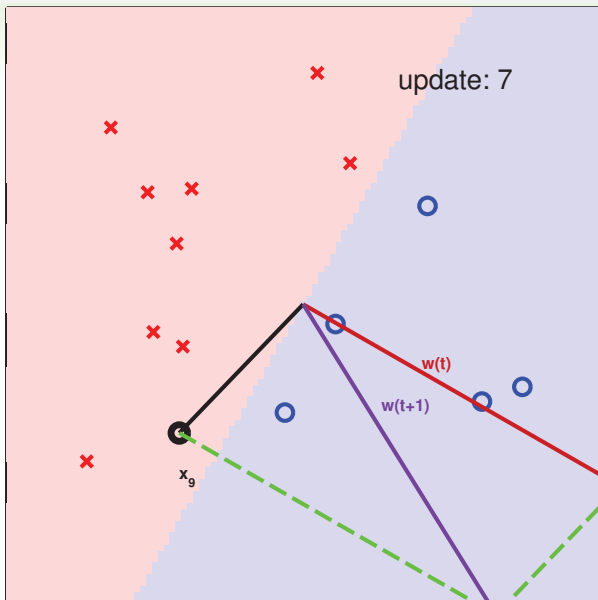
Seeing is Believing



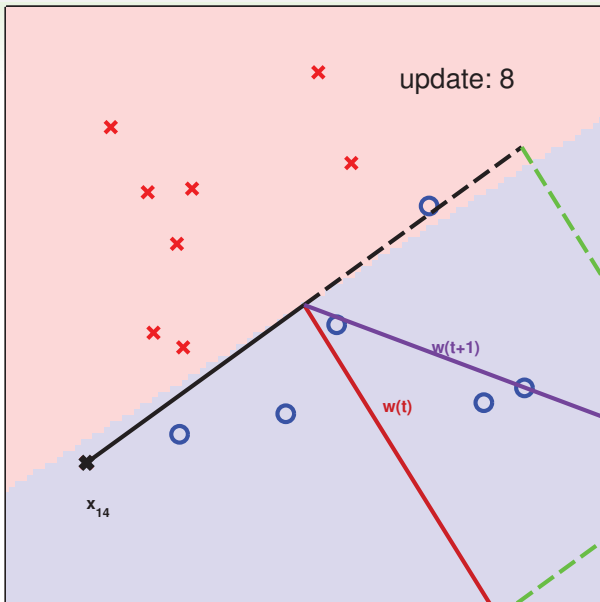
Seeing is Believing



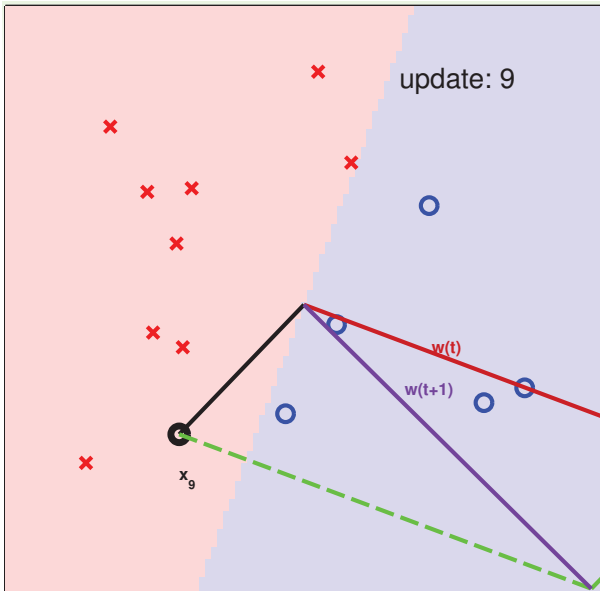
Seeing is Believing



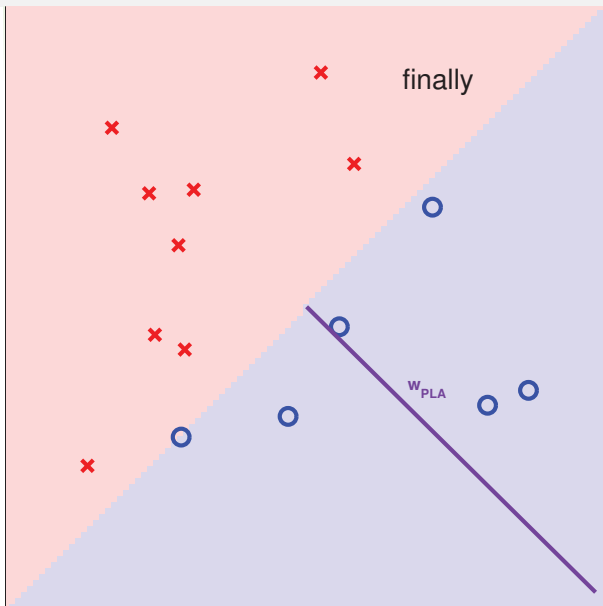
Seeing is Believing



Seeing is Believing



Seeing is Believing



“correct” mistakes on \mathcal{D} until no mistakes

Algorithmic: halt (with no mistake)?

- naive cyclic: ?
- random cyclic: ?
- other variant: ?

Learning: $g \approx f$?

- on \mathcal{D} , if halt, yes (no mistake)
- outside \mathcal{D} : ?
- if not halting: ?

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 $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n < y_n \mathbf{w}_t^T \mathbf{x}_n$

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?

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

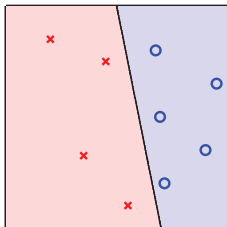
- ① $\mathbf{w}_{t+1}^T \mathbf{x}_n = y_n$
- ② $\text{sign}(\mathbf{w}_{t+1}^T \mathbf{x}_n) = y_n$
- ③ $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n \geq y_n \mathbf{w}_t^T \mathbf{x}_n$
- ④ $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n < y_n \mathbf{w}_t^T \mathbf{x}_n$

Reference Answer: ③

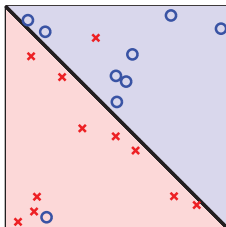
Simply multiply the second part of the rule by $y_n \mathbf{x}_n$. The result shows that the rule somewhat “tries to correct the mistake.”

- 1 学习与机器学习
- 2 机器学习组成要素及与其它领域的关系
- 3 感知机假说集及感知机学习算法
- 4 感知机学习算法的理论保证
- 5 非可分数据

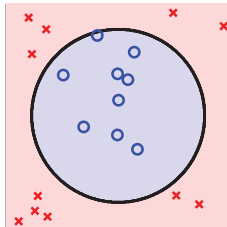
- if PLA halts (i.e. no more mistakes), (necessary condition)
 \mathcal{D} allows some w to make no mistake.
- call such \mathcal{D} **linear separable**



(linear separable)

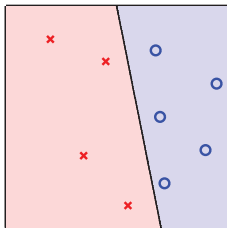


(not linear separable)

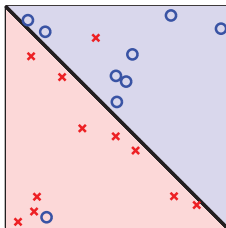


(not linear separable)

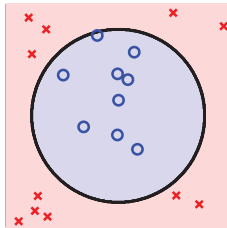
- if PLA halts (i.e. no more mistakes), (necessary condition)
 \mathcal{D} allows some w to make no mistake.
- call such \mathcal{D} **linear separable**



(linear separable)



(not linear separable)



(not linear separable)

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

观测一: \mathbf{w}_t 与 \mathbf{w}_f 愈加靠近

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

- \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 $(\mathbf{x}_{n(t)}, y_{n(t)})$

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

观测一: \mathbf{w}_t 与 \mathbf{w}_f 愈加靠近

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

- \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?

观测二: \mathbf{w}_t 并没有增长过快

\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,

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

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

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

$$T \leq R^2/\rho^2.$$

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

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

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
- not fully sure how long halting takes (ρ depends on \mathbf{w}_f)
 - though practically fast

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

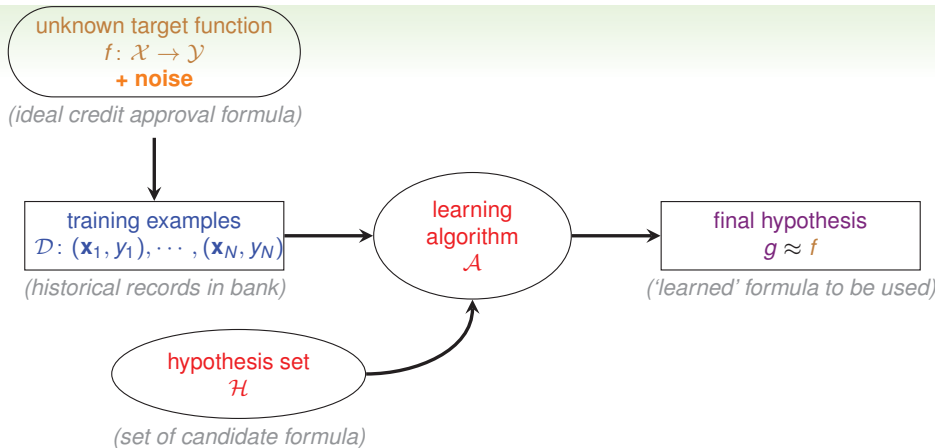
simple to implement, fast, works in any dimension d

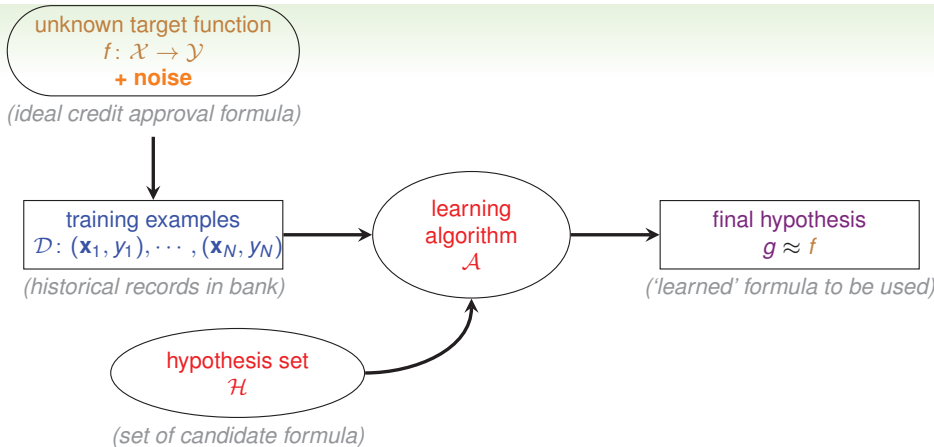
Cons

- “assumes” linear separable \mathcal{D} to halt
 - property unknown in advance
- not fully sure how long halting takes (ρ depends on \mathbf{w}_f)
 - though practically fast

what if \mathcal{D} not linear separable?

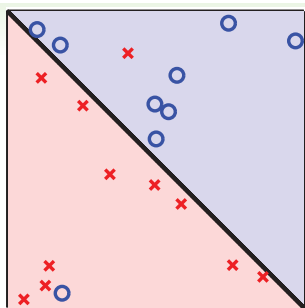
- 1 学习与机器学习
- 2 机器学习组成要素及与其它领域的关系
- 3 感知机假说集及感知机学习算法
- 4 感知机学习算法的理论保证
- 5 非可分数据





how to at least get $g \approx f$ on **noisy** data?

容忍噪声的线性分类器

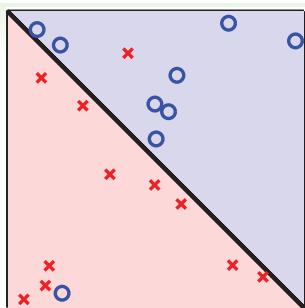


- assume 'little' noise: $y_n = f(\mathbf{x}_n)$ **usually**
- if so, $g \approx f$ on $\mathcal{D} \Leftrightarrow y_n = g(\mathbf{x}_n)$ **usually**
- how about

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

—**NP-hard to solve, unfortunately**

容忍噪声的线性分类器



- 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**

can we modify PLA to get an **approximately good** g ?

modify PLA algorithm by **keeping best weights in pocket**

initialize pocket weights $\hat{\mathbf{w}}$

For $t = 0, 1, \dots$

- 1 find a (random) mistake of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$
- 2 (try to) correct the mistake by

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

- 3 if \mathbf{w}_{t+1} makes fewer mistakes than $\hat{\mathbf{w}}$, replace $\hat{\mathbf{w}}$ by \mathbf{w}_{t+1}

...until **enough iterations**

return $\hat{\mathbf{w}}$ (called $\mathbf{w}_{\text{POCKET}}$) as g

modify PLA algorithm by **keeping best weights in pocket**

initialize pocket weights $\hat{\mathbf{w}}$

For $t = 0, 1, \dots$

- 1 find a (random) mistake of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$
- 2 (try to) correct the mistake by

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

- 3 if \mathbf{w}_{t+1} makes fewer mistakes than $\hat{\mathbf{w}}$, replace $\hat{\mathbf{w}}$ by \mathbf{w}_{t+1}

...until **enough iterations**

return $\hat{\mathbf{w}}$ (called $\mathbf{w}_{\text{POCKET}}$) as g

a simple modification of PLA to find (somewhat) “best” weights

Should we use pocket or PLA?

Since we do not know whether \mathcal{D} is linear separable in advance, we may decide to just go with pocket instead of PLA. If \mathcal{D} is actually linear separable, what's the difference between the two?

- 1 pocket on \mathcal{D} is slower than PLA
- 2 pocket on \mathcal{D} is faster than PLA
- 3 pocket on \mathcal{D} returns a better g in approximating f than PLA
- 4 pocket on \mathcal{D} returns a worse g in approximating f than PLA

Should we use pocket or PLA?

Since we do not know whether \mathcal{D} is linear separable in advance, we may decide to just go with pocket instead of PLA. If \mathcal{D} is actually linear separable, what's the difference between the two?

- ① pocket on \mathcal{D} is slower than PLA
- ② pocket on \mathcal{D} is faster than PLA
- ③ pocket on \mathcal{D} returns a better g in approximating f than PLA
- ④ pocket on \mathcal{D} returns a worse g in approximating f than PLA

Reference Answer: ①

Because pocket needs to check whether \mathbf{w}_{t+1} is better than $\hat{\mathbf{w}}$ at each iteration, it is slower than PLA. On linear separable \mathcal{D} , $\mathbf{w}_{\text{POCKET}}$ is the same as \mathbf{w}_{PLA} , both making no mistakes.

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
- Non-Separable Data
hold somewhat 'best' weights in pocket

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
- Non-Separable Data
hold somewhat 'best' weights in pocket

Homework 1: **Implement PLA (Pocket) algorithm.**

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
- Non-Separable Data
hold somewhat 'best' weights in pocket

Homework 1: **Implement PLA (Pocket) algorithm.**

Next: **Support Vector Machines and Kernel Methods.**