

# The Boosting Approach to Machine Learning

Maria-Florina Balcan

03/16/2015

# Boosting

- General method for improving the accuracy of any given learning algorithm.
- Works by creating a series of challenge datasets s.t. even modest performance on these can be used to produce an overall high-accuracy predictor.
  - Works amazingly well in practice --- Adaboost and its variations one of the top 10 algorithms.
  - Backed up by solid foundations.

# Readings:



- The Boosting Approach to Machine Learning: An Overview. Rob Schapire, 2001
- Theory and Applications of Boosting. NIPS tutorial.  
<http://www.cs.princeton.edu/~schapire/talks/nips-tutorial.pdf>

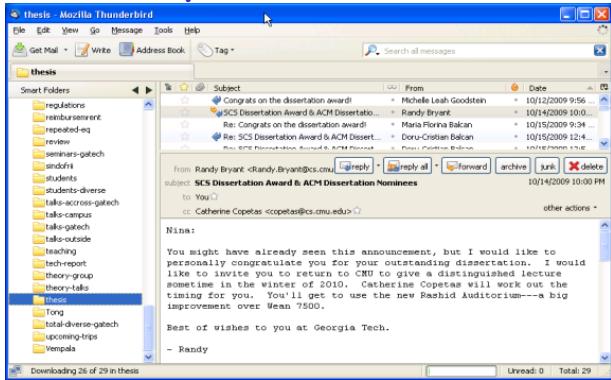
# Plan for today:

- Motivation.
- A bit of history.
- Adaboost: algo, guarantees, discussion.
- Focus on supervised classification.

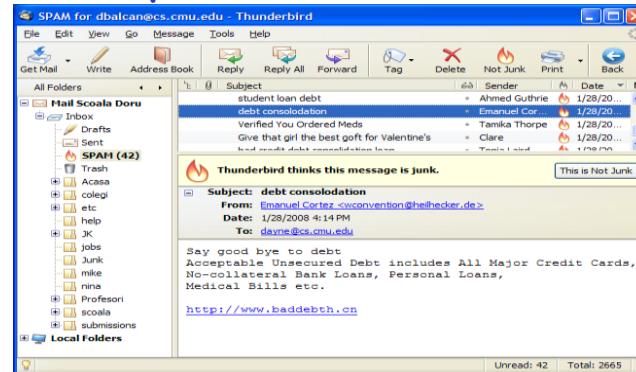
# An Example: Spam Detection

- E.g., classify which emails are spam and which are important.

Not spam



spam

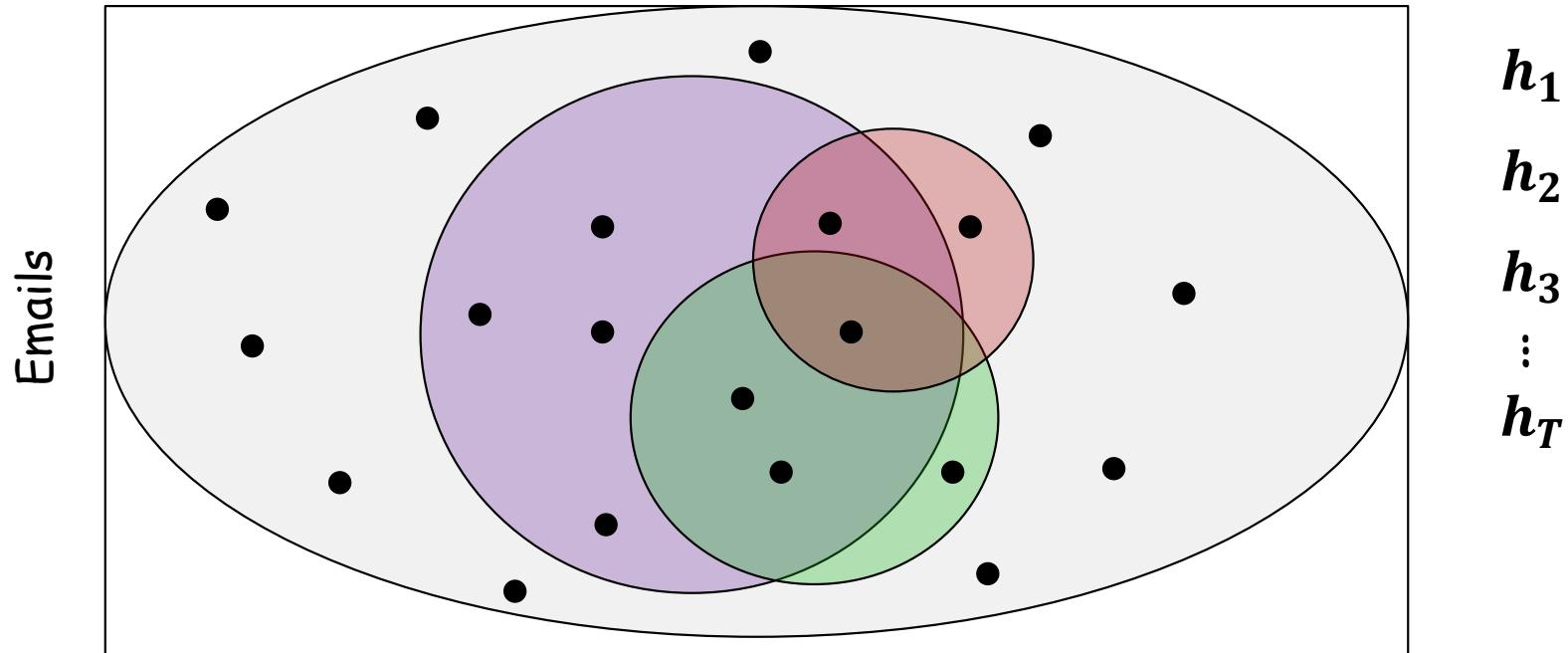


Key observation/motivation:

- Easy to find rules of thumb that are often correct.
  - E.g., "If buy now in the message, then predict spam."
  - E.g., "If say good-bye to debt in the message, then predict spam."
- Harder to find single rule that is very highly accurate.

# An Example: Spam Detection

- Boosting: meta-procedure that takes in an algo for finding rules of thumb (weak learner). Produces a highly accurate rule, by calling the weak learner repeatedly on cleverly chosen datasets.



- apply weak learner to a subset of emails, obtain rule of thumb
- apply to 2nd subset of emails, obtain 2nd rule of thumb
- apply to 3<sup>rd</sup> subset of emails, obtain 3rd rule of thumb
- repeat T times; combine weak rules into a single highly accurate rule.

→ 改变分类的权重.

# Boosting: Important Aspects

How to choose examples on each round?

- Typically, concentrate on “hardest” examples (those most often misclassified by previous rules of thumb)

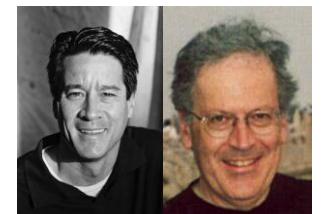
How to combine rules of thumb into single prediction rule?

- take (weighted) majority vote of rules of thumb

Historically....

# Weak Learning vs Strong/PAC Learning

- [Kearns & Valiant '88]: defined weak learning: being able to predict better than random guessing ( $\text{error} \leq \frac{1}{2} - \gamma$ ) , consistently. 稍少类错。  
 $\gamma \rightarrow 0$
- Posed an open pb: "Does there exist a boosting algo that turns a weak learner into a strong PAC learner (that can produce arbitrarily accurate hypotheses)?"  
 $\hookrightarrow \text{error} < \epsilon.$
- Informally, given "weak" learning algo that can consistently find classifiers of  $\text{error} \leq \frac{1}{2} - \gamma$ , a boosting algo would provably construct a single classifier with  $\text{error} \leq \epsilon$ .



# Weak Learning vs Strong/PAC Learning

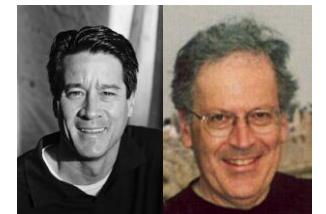
## Strong (PAC) Learning

- $\exists$  algo  $A$
  - $\forall c \in H$
  - $\forall D$
  - $\forall \epsilon > 0$
  - $\forall \delta > 0$
  - $A$  produces  $h$  s.t.:  
$$\Pr[\text{err}(h) \geq \epsilon] \leq \delta$$
- $\Pr[\text{err}(h) \leq \delta] > 1 - \delta$

## Weak Learning

- $\exists$  algo  $A$
- $\exists \gamma > 0$
- $\forall c \in H$
- $\forall D$
- $\forall \epsilon > \frac{1}{2} - \gamma$
- $\forall \delta > 0$
- $A$  produces  $h$  s.t.  
$$\Pr[\text{err}(h) \geq \epsilon] \leq \delta$$

- [Kearns & Valiant '88]: defined weak learning & posed an open pb of finding a boosting algo.



# Surprisingly....

## Weak Learning = Strong (PAC) Learning

Original Construction [Schapire '89]:

- poly-time boosting algo, exploits that we can learn a little on **every** distribution.
- A modest booster obtained via calling the weak learning algorithm on 3 distributions.  
$$\text{Error} = \beta < \frac{1}{2} - \gamma \rightarrow \text{error } 3\beta^2 - 2\beta^3$$

多次投票 → 合成 ->
- Then amplifies the modest boost of accuracy by running this somehow recursively.
- Cool conceptually and technically, not very practical.



# An explosion of subsequent work

## Background (cont.)

- [Freund & Schapire '95]:
  - introduced “AdaBoost” algorithm
  - strong practical advantages over previous boosting algorithms
- experiments and applications using AdaBoost:

[Drucker & Cortes '96]

[Jackson & Craven '96]

[Freund & Schapire '96]

[Quinlan '96]

[Breiman '96]

[Maclin & Opitz '97]

[Bauer & Kohavi '97]

[Schwenk & Bengio '98]

[Schapire, Singer & Singhal '98]

[Abney, Schapire & Singer '99]

[Haruno, Shirai & Ooyama '99]

[Cohen & Singer '99]

[Dietterich '00]

[Schapire & Singer '00]

[Collins '00]

[Escudero, Márquez & Rigau '00]

[Iyer, Lewis, Schapire, Singer & Singhal '00]

[Onoda, Rätsch & Müller '00]

[Tieu & Viola '00]

[Walker, Rambow & Rogati '01]

[Rochery, Schapire, Rahim & Gupta '01]

[Merler, Furlanello, Larcher & Sboner '01]

:

- continuing development of theory and algorithms:

[Breiman '98, '99]

[Schapire, Freund, Bartlett & Lee '98]

[Grove & Schuurmans '98]

[Mason, Bartlett & Baxter '98]

[Schapire & Singer '99]

[Cohen & Singer '99]

[Freund & Mason '99]

[Domingo & Watanabe '99]

[Mason, Baxter, Bartlett & Frean '99, '00]

[Duffy & Helmbold '99, '02]

[Freund & Mason '99]

[Ridgeway, Madigan & Richardson '99]

[Kivinen & Warmuth '99]

[Friedman, Hastie & Tibshirani '00]

[Rätsch, Onoda & Müller '00]

[Rätsch, Warmuth, Mika, Onoda, Lemm &

Müller '00]

[Allwein, Schapire & Singer '00]

[Friedman '01]

[Koltchinskii, Panchenko & Lozano '01]

[Collins, Schapire & Singer '02]

[Demiriz, Bennett & Shawe-Taylor '02]

[Lebanon & Lafferty '02]

:

# Adaboost (Adaptive Boosting)

"A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting"

[Freund-Schapire, JCSS'97]

Godel Prize winner 2003

# Informal Description Adaboost

- Boosting: turns a weak algo into a strong (PAC) learner.

Input:  $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ ;  $x_i \in X, y_i \in Y = \{-1, 1\}$

weak learning algo  $\xrightarrow{\text{用强弱可解决}} A$  (e.g., Naïve Bayes, decision stumps)

- For  $t=1, 2, \dots, T$   $\epsilon_t = \Pr_{D_t} [h_t(x_i) \neq y_i]$   $\xleftarrow{h_t}$
  - Construct  $D_t$  on  $\{x_1, \dots, x_m\} = \sum_{i=1}^m D_t(x_i) \mathbb{I}[h_t(x_i) \neq y_i]$   $\xrightarrow{+}$
  - Run  $A$  on  $D_t$  producing  $h_t: X \rightarrow \{-1, 1\}$  (weak classifier)
  - Output  $H_{\text{final}}(x) = \text{sign}(\sum_{t=1} \alpha_t h_t(x))$
- 
- $$\epsilon_t = P_{x_i \sim D_t}(h_t(x_i) \neq y_i) \text{ error of } h_t \text{ over } D_t \Rightarrow \epsilon_t = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h_t(x_i) \neq y_i]$$

Roughly speaking  $D_{t+1}$  increases weight on  $x_i$  if  $h_t$  incorrect on  $x_i$  ; decreases it on  $x_i$  if  $h_t$  correct.

# Adaboost (Adaptive Boosting)

- Weak learning algorithm A.
- For  $t=1, 2, \dots, T$ 
  - Construct  $D_t$  on  $\{x_1, \dots, x_m\}$
  - Run  $A$  on  $D_t$  producing  $h_t$

## Constructing $D_t$

- $D_1$  uniform on  $\{x_1, \dots, x_m\}$  [i.e.,  $D_1(i) = \frac{1}{m}$ ]

- Given  $D_t$  and  $h_t$  set

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{-\alpha_t\}} \text{ if } y_i = h_t(x_i)$$

权重减少 (因为已正确)

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{\alpha_t\}} \text{ if } y_i \neq h_t(x_i)$$

分类错误  $\Rightarrow$  提高权重

积分因子 使  $\sum P(x) = 1$

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{-\alpha_t y_i h_t(x_i)\}}$$

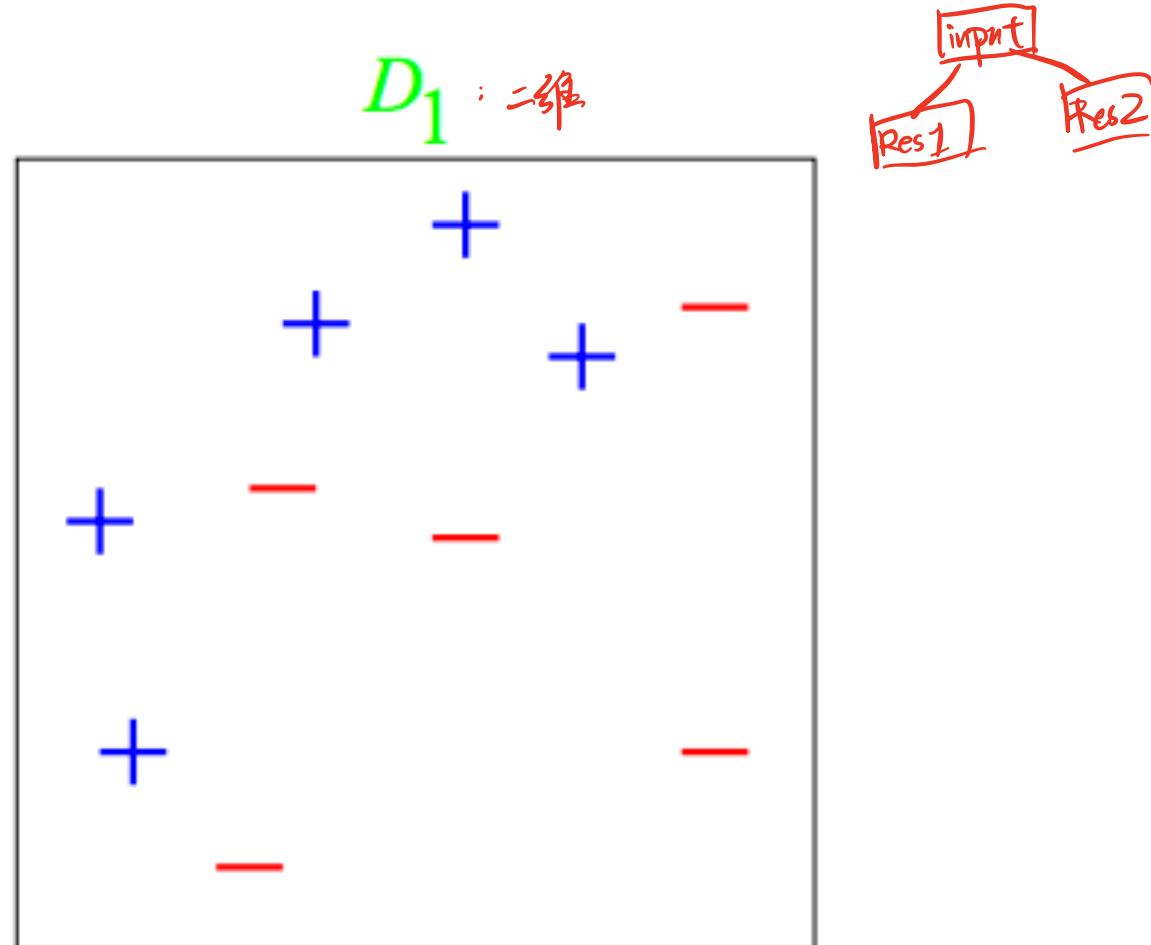
$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right) > 0$$

$D_{t+1}$  puts half of weight on examples  $x_i$  where  $h_t$  is incorrect & half on examples where  $h_t$  is correct

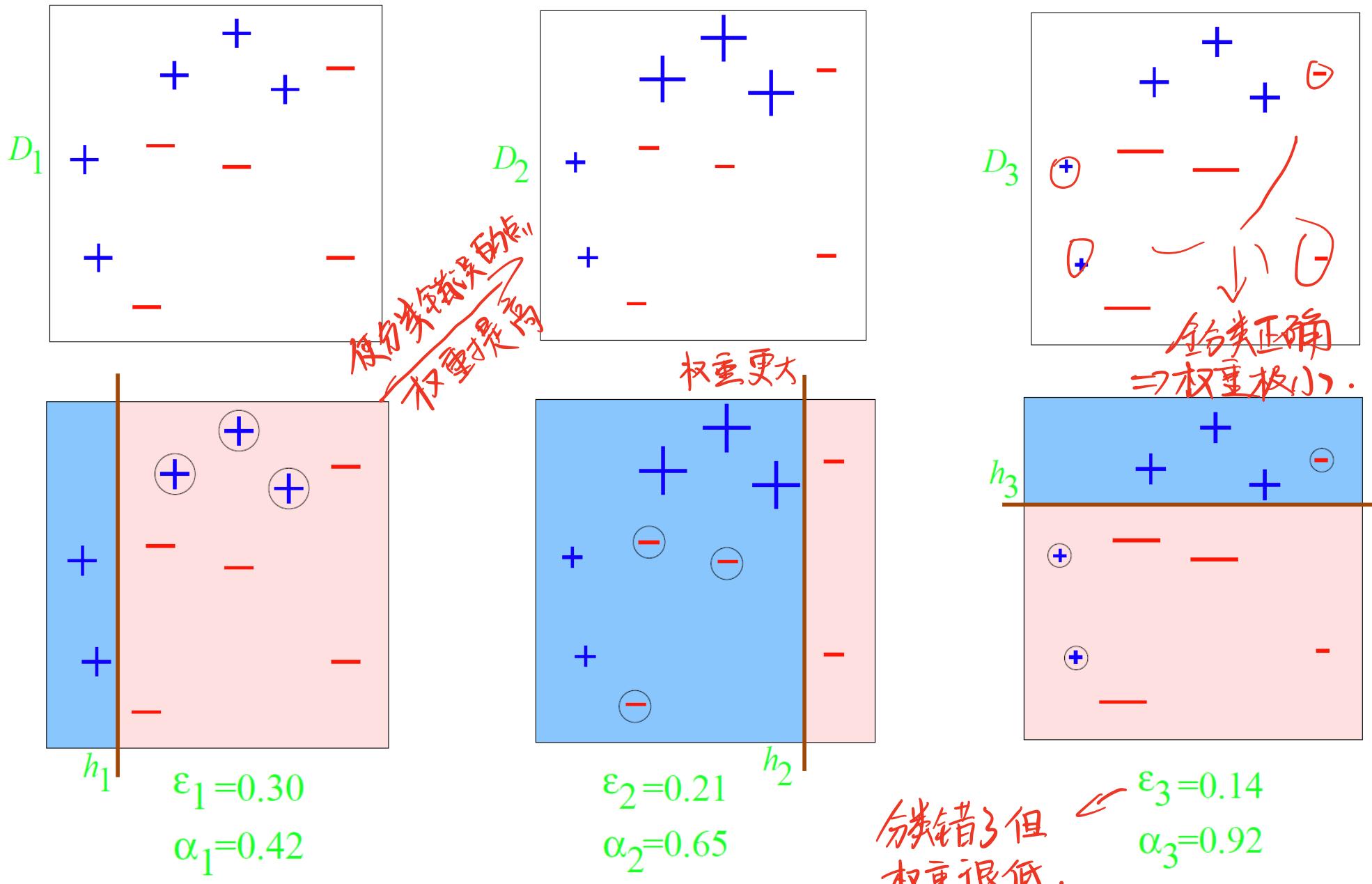
Final hyp:  $H_{\text{final}}(x) = \text{sign}(\sum_t \alpha_t h_t(x))$

# Adaboost: A toy example

Weak classifiers: vertical or horizontal half-planes (a.k.a. decision stumps) 决策树桩



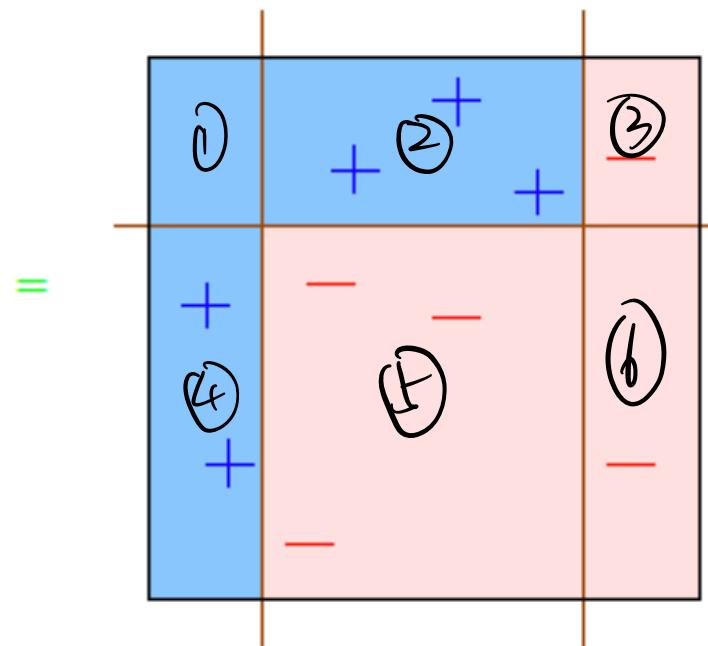
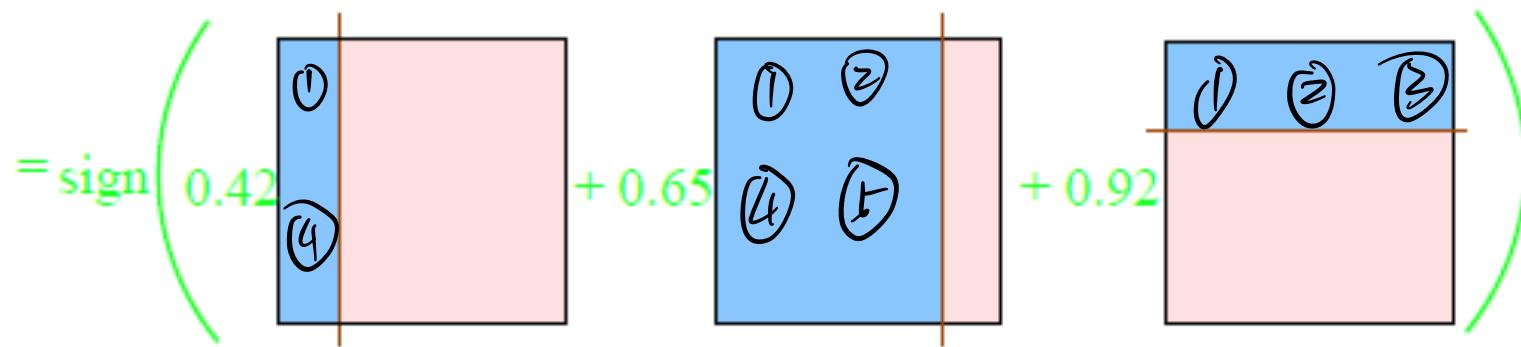
# Adaboost: A toy example



# Adaboost: A toy example

假设设有T=3个分类器，那么对这三者进行多数投票.

$H_{\text{final}}$



把每一块的权重相加去看.

# Adaboost (Adaptive Boosting)

- Weak learning algorithm  $A$ .
- For  $t=1, 2, \dots, T$ 
  - Construct  $D_t$  on  $\{x_1, \dots, x_m\}$
  - Run  $A$  on  $D_t$  producing  $h_t$

## Constructing $D_t$

- $D_1$  uniform on  $\{x_1, \dots, x_m\}$  [i.e.,  $D_1(i) = \frac{1}{m}$ ]
- Given  $D_t$  and  $h_t$  set

$$\left. \begin{array}{l} D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{-\alpha_t\}} \text{ if } y_i = h_t(x_i) \\ D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{\alpha_t\}} \text{ if } y_i \neq h_t(x_i) \end{array} \right\}$$

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{-\alpha_t y_i h_t(x_i)\}}$$

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right) > 0$$

$D_{t+1}$  puts half of weight on examples  $x_i$  where  $h_t$  is incorrect & half on examples where  $h_t$  is correct

Final hyp:  $H_{\text{final}}(x) = \text{sign}(\sum_t \alpha_t h_t(x))$

缺点:①对异常点敏感,②对样本数量要求高.

# Nice Features of Adaboost

用强分类器也可以但是提升太慢

- Very general: a meta-procedure, it can use any weak learning algorithm!!! (e.g., Naive Bayes, decision stumps)
- Very fast (single pass through data each round) & simple to code, no parameters to tune.
- Shift in mindset: goal is now just to find classifiers a bit better than random guessing.
- Grounded in rich theory.
- Relevant for big data age: quickly focuses on "core difficulties", well-suited to distributed settings, where data must be communicated efficiently [Balcan-Blum-Fine-Mansour COLT'12].

# Analyzing Training Error

**Theorem**  $\epsilon_t = 1/2 - \gamma_t$  (error of  $h_t$  over  $D_t$ )

$$err_S(H_{final}) \leq \exp \left[ -2 \sum_t \gamma_t^2 \right]$$

训练误差.

第t次的  $\gamma$

So, if  $\forall t, \gamma_t \geq \gamma > 0$ , then  $err_S(H_{final}) \leq \exp[-2\gamma^2 T]$

有一个下确界.

假设  $err_S(H_{final}) \leq e^{-2\gamma^2 T} \leq \epsilon$

The training error drops exponentially in  $T$ !!!

To get  $err_S(H_{final}) \leq \epsilon$ , need only  $T = O\left(\frac{1}{\gamma^2} \log\left(\frac{1}{\epsilon}\right)\right)$  rounds

## Adaboost is adaptive

- Does not need to know  $\gamma$  or  $T$  a priori
- Can exploit  $\gamma_t \gg \gamma$

# Understanding the Updates & Normalization

**Claim:**  $D_{t+1}$  puts half of the weight on  $x_i$  where  $h_t$  was incorrect and half of the weight on  $x_i$  where  $h_t$  was correct.

$$\text{Recall } D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{-\alpha_t y_i h_t(x_i)\}}$$

每次分类，将一半权重分给错误分类的，另一半给分类正确的  
Probabilities are equal!

$$\Pr_{D_{t+1}} [y_i \neq h_t(x_i)] = \sum_{i:y_i \neq h_t(x_i)} \frac{D_t(i)}{Z_t} e^{\alpha_t} = \epsilon_t \frac{1}{Z_t} e^{\alpha_t} = \frac{\epsilon_t}{Z_t} \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} = \frac{\sqrt{\epsilon_t(1 - \epsilon_t)}}{Z_t}$$

$$\Pr_{D_{t+1}} [y_i = h_t(x_i)] = \sum_{i:y_i = h_t(x_i)} \frac{D_t(i)}{Z_t} e^{-\alpha_t} = \frac{1 - \epsilon_t}{Z_t} e^{-\alpha_t} = \frac{1 - \epsilon_t}{Z_t} \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} = \frac{\sqrt{(1 - \epsilon_t)\epsilon_t}}{Z_t}$$

$P(\text{correct}) = P(\text{error})$  且有  $Z_t \Rightarrow P(\text{correct}) + P(\text{error}) = 1$  因为分类正确的点更多  
故每个点的概率也会偏小。

$$\begin{aligned} Z_t &= \sum_{i:y_i = h_t(x_i)} D_t(i) e^{-\alpha_t y_i h_t(x_i)} = \sum_{i:y_i = h_t(x_i)} D_t(i) e^{-\alpha_t} + \sum_{i:y_i \neq h_t(x_i)} D_t(i) e^{\alpha_t} \\ &= (1 - \epsilon_t) e^{-\alpha_t} + \epsilon_t e^{\alpha_t} = 2\sqrt{\epsilon_t(1 - \epsilon_t)} \end{aligned}$$

规范化因子

# Analyzing Training Error: Proof Intuition

**Theorem**  $\epsilon_t = 1/2 - \gamma_t$  (error of  $h_t$  over  $D_t$ )

$$err_S(H_{final}) \leq \exp \left[ -2 \sum_t \gamma_t^2 \right]$$

- On round  $t$ , we increase weight of  $x_i$  for which  $h_t$  is wrong.
- If  $H_{final}$  incorrectly classifies  $x_i$ ,
  - Then  $x_i$  incorrectly classified by (wtd) majority of  $h_t$ 's.
  - Which implies final prob. weight of  $x_i$  is large.
- Since sum of prob. = 1, can't have too many of high weight.

Can show # incorrectly classified  $\leq m (\prod_t Z_t)$ .

And  $(\prod_t Z_t) \rightarrow 0$ .

# Analyzing Training Error: Proof Math

Step 1: unwrapping recurrence:  $D_{T+1}(i) = \frac{1}{m} \left( \frac{\exp(-y_i f(x_i))}{\prod_t z_t} \right)$   
where  $f(x_i) = \sum_t \alpha_t h_t(x_i)$ . [Unthresholded weighted vote of  $h_i$  on  $x_i$  ]

Step 2:  $\text{err}_S(H_{final}) \leq \prod_t Z_t$ .

Step 3:  $\prod_t Z_t = \prod_t 2\sqrt{\epsilon_t(1 - \epsilon_t)} = \prod_t \sqrt{1 - 4\gamma_t^2} \leq e^{-2 \sum_t \gamma_t^2}$

# Analyzing Training Error: Proof Math

Step 1: unwrapping recurrence:  $D_{T+1}(i) = \frac{1}{m} \left( \frac{\exp(-y_i f(x_i))}{\prod_t z_t} \right)$   
where  $f(x_i) = \sum_t \alpha_t h_t(x_i)$ .

Recall  $\underbrace{D_1(i)}_{\text{初始}} = \frac{1}{m}$  and  $\underbrace{D_{t+1}(i)}_{\text{递归关系}} = D_t(i) \frac{\exp(-y_i \alpha_t h_t(x_i))}{z_t}$

$$\begin{aligned} D_{T+1}(i) &= \frac{\exp(-y_i \alpha_T h_T(x_i))}{Z_T} \times D_T(i) \\ &= \frac{\exp(-y_i \alpha_T h_T(x_i))}{Z_T} \times \frac{\exp(-y_i \alpha_{T-1} h_{T-1}(x_i))}{Z_{T-1}} \times D_{T-1}(i) \end{aligned}$$

.....

$$\begin{aligned} &= \frac{\exp(-y_i \alpha_T h_T(x_i))}{Z_T} \times \dots \times \frac{\exp(-y_i \alpha_1 h_1(x_i))}{Z_1} \frac{1}{m} \\ &= \frac{1}{m} \frac{\exp(-y_i(\alpha_1 h_1(x_i) + \dots + \alpha_T h_T(x_T)))}{Z_1 \cdots Z_T} = \frac{1}{m} \frac{\exp(-y_i f(x_i))}{\prod_t z_t} \end{aligned}$$

# Analyzing Training Error: Proof Math

Step 1: unwrapping recurrence:  $D_{T+1}(i) = \frac{1}{m} \left( \frac{\exp(-y_i f(x_i))}{\prod_t Z_t} \right)$

where  $f(x_i) = \sum_t \alpha_t h_t(x_i)$ .

$$f(x_i) = \sum_{t=1}^T \alpha_t h_t(x_i)$$

$$H(x_i) = \text{sign}(f(x_i))$$

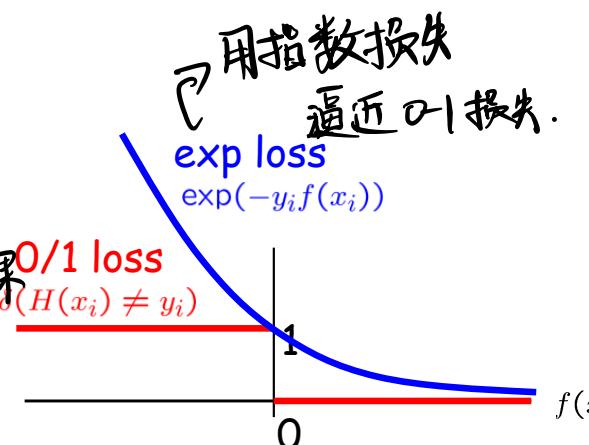
Step 2:  $\text{err}_S(H_{final}) \leq \prod_t Z_t$ .

$$\begin{aligned} \text{err}_S(H_{final}) &= \frac{1}{m} \sum_i \mathbf{1}_{y_i \neq H_{final}(x_i)} \\ &= \frac{1}{m} \sum_i \mathbf{1}_{y_i f(x_i) \leq 0} \end{aligned}$$

$\rightarrow$  预测结果与观测结果不一致

$$\leq \frac{1}{m} \sum_i \underbrace{\exp(-y_i f(x_i))}_{\text{是上界.}}$$

$$= \sum_i D_{T+1}(i) \prod_t Z_t = \prod_t Z_t.$$



# Analyzing Training Error: Proof Math

Step 1: unwrapping recurrence:  $D_{T+1}(i) = \frac{1}{m} \left( \frac{\exp(-y_i f(x_i))}{\prod_t Z_t} \right)$   
where  $f(x_i) = \sum_t \alpha_t h_t(x_i)$ .

Step 2:  $\text{err}_S(H_{final}) \leq \prod_t Z_t$ .

$$fx \leq e^{-x} \Rightarrow \sqrt{fx} \leq e^{-\frac{x}{2}}$$

Step 3:  $\prod_t Z_t = \prod_t 2\sqrt{\epsilon_t(1-\epsilon_t)} = \prod_t \sqrt{1-4\gamma_t^2} \leq e^{-2 \sum_t \gamma_t^2}$

Note: recall  $Z_t = (1-\epsilon_t)e^{-\alpha_t} + \epsilon_t e^{\alpha_t} = 2\sqrt{\epsilon_t(1-\epsilon_t)}$

$\alpha_t$  minimizer of  $\alpha \rightarrow (1-\epsilon_t)e^{-\alpha} + \epsilon_t e^{\alpha}$

# Analyzing Training Error: Proof Intuition

**Theorem**  $\epsilon_t = 1/2 - \gamma_t$  (error of  $h_t$  over  $D_t$ )

$$err_S(H_{final}) \leq \exp \left[ -2 \sum_t \gamma_t^2 \right]$$

- On round  $t$ , we increase weight of  $x_i$  for which  $h_t$  is wrong.
- If  $H_{final}$  incorrectly classifies  $x_i$ ,
  - Then  $x_i$  incorrectly classified by (wtd) majority of  $h_t$ 's.
  - Which implies final prob. weight of  $x_i$  is large.
- Since sum of prob. = 1, can't have too many of high weight.

Can show probability  $\geq \frac{1}{m} \left( \frac{1}{\prod_t Z_t} \right)$   $D_{t+1}(i) = \frac{1}{m} \frac{\exp(-y_i f(x_i))}{\prod_t Z_t}$

Can show # incorrectly classified  $\leq m (\prod_t Z_t)$ .

And  $(\prod_t Z_t) \rightarrow 0$ .

↳ 否则概率之和超过1

# Generalization Guarantees

**Theorem**  $\text{err}_S(H_{final}) \leq \exp\left[-2 \sum_t \gamma_t^2\right]$  where  $\epsilon_t = 1/2 - \gamma_t$

How about generalization guarantees?



Original analysis [Freund&Schapire'97]

- $H$  space of weak hypotheses;  $d = \text{VCdim}(H)$

$H_{final}$  is a weighted vote, so the hypothesis class is:

$G = \{\text{all fns of the form } \text{sign}(\sum_{t=1}^T \alpha_t h_t(x))\}$

Theorem [Freund&Schapire'97]

$$\forall g \in G, \text{err}(g) \leq \text{err}_S(g) + \tilde{O}\left(\sqrt{\frac{Td}{m}}\right)$$

描述算法复杂度  
样本个数

T= # of rounds

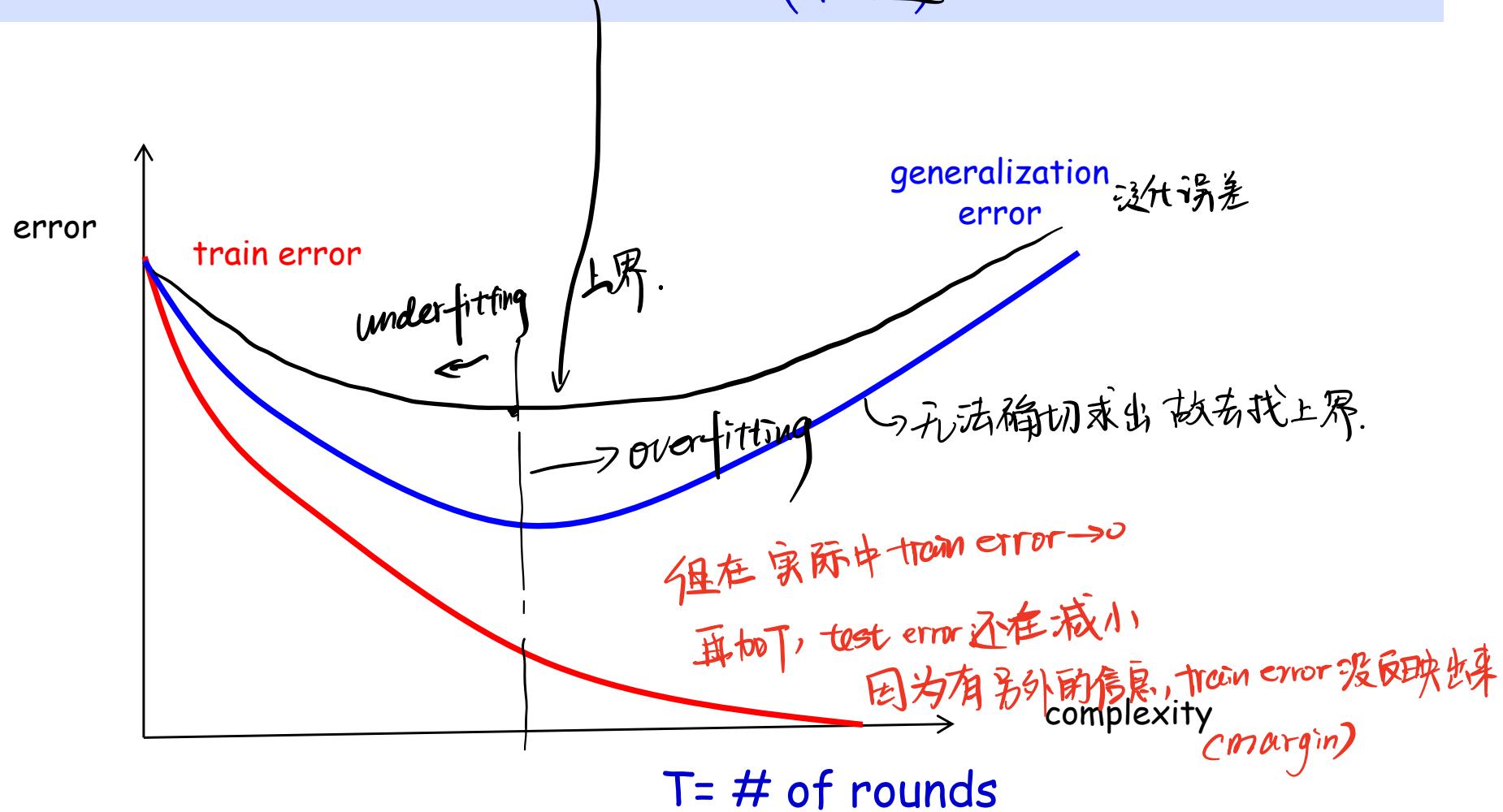
Key reason:  $\text{VCdim}(G) = \tilde{O}(dT)$  plus typical VC bounds.

# Generalization Guarantees

Theorem [Freund&Schapire'97]

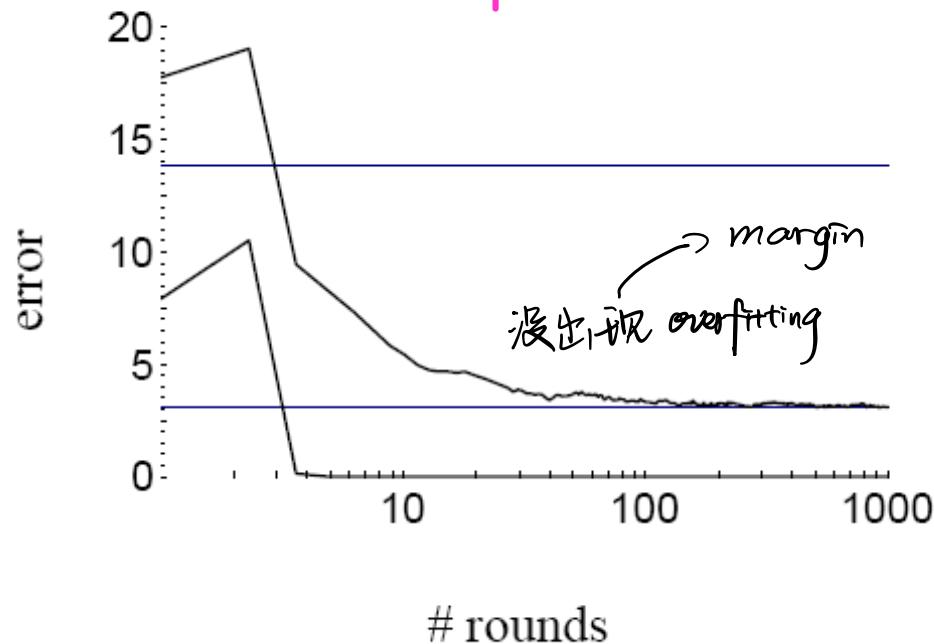
$$\forall g \in co(H), \text{err}(g) \leq \underbrace{\text{err}_s(g)}_{\text{train error}} + \tilde{O}\left(\sqrt{\frac{Td}{m}}\right) \text{ where } d = \text{VCdim}(H)$$

模型复杂度



# Generalization Guarantees

- Experiments with boosting showed that the test error of the generated classifier usually **does not increase** as its size becomes very large.
- Experiments showed that continuing to add new weak learners after **correct** classification of the training set had been achieved could further **improve** test set performance!!!



# Generalization Guarantees

- Experiments with boosting showed that the test error of the generated classifier usually **does not increase** as its size becomes very large.
- Experiments showed that continuing to add **weak** learners after **correct** classification of the training set had been achieved could further **improve** test set performance!!!
- These results seem to contradict FS'87 bound and Occam's razor. In order to achieve good test error the classifier should be as simple as



# How can we explain the experiments?

R. Schapire, Y. Freund, P. Bartlett, W. S. Lee. present in "Boosting the margin: A new explanation for the effectiveness of voting methods" a nice theoretical explanation.

**Key Idea:**

Training error does not tell the whole story.

We need also to consider the classification confidence!!

$$H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$$

但  $\sum_{t=1}^T \alpha_t h_t(x)$  可以是 1, 0.1, 0.001

而 0.001 有极大可能变成负的。

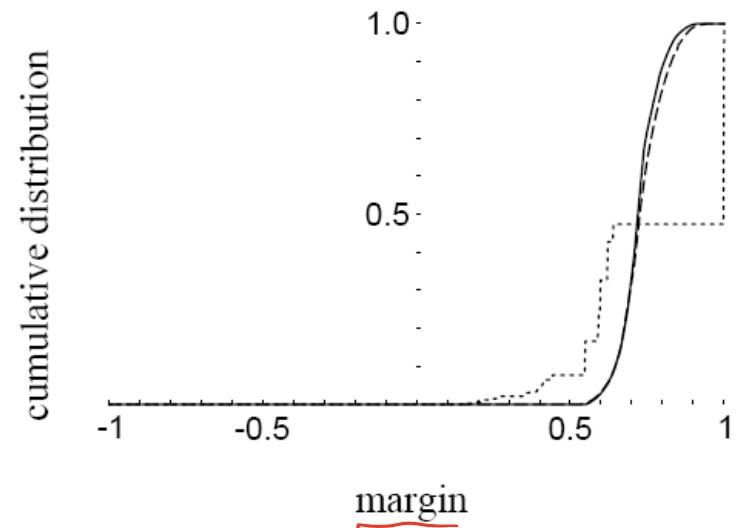
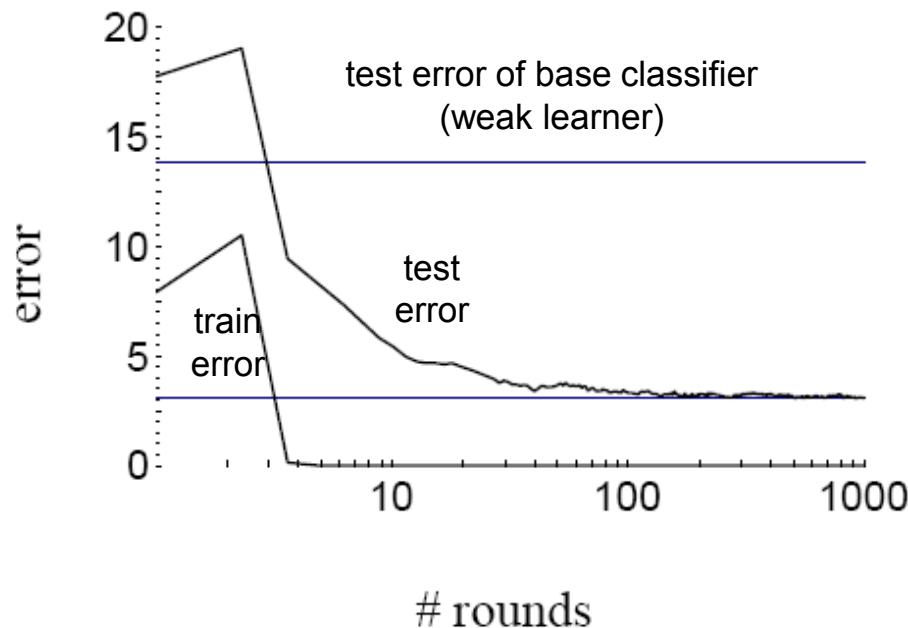
$|f(x)|$  越大，置信度越高。

$$\begin{aligned} \text{训练误差概率: } \text{err}_S &= P_{r_S}(y \neq H(x)) \\ &= P_{r_S}(y H(x) < 0) = P_{r_S}(y f(x) < 0) \end{aligned}$$

margin :  $y f(x)$  函数间隔。

Boosting didn't seem  
to overfit...(!)

...because it turned out to be  
increasing the margin of the  
classifier



Error Curve, Margin Distr. Graph - Plots from [SFBL98]

# Classification Margin

对  $\alpha_t$  归一化:  $\alpha_t = \frac{\alpha_t}{\sum_{t=1}^T \alpha_t}$  使其变成类似概率的形式.

- $H$  space of weak hypotheses. The **convex hull** of  $H$ :

$$co(H) = \{f = \sum_{t=1}^T \alpha_t h_t, \alpha_t \geq 0, \sum_{t=1}^T \alpha_t = 1, h_t \in H\}$$

- Let  $f \in co(H), f = \sum_{t=1}^T \alpha_t h_t, \alpha_t \geq 0, \sum_{t=1}^T \alpha_t = 1$ .

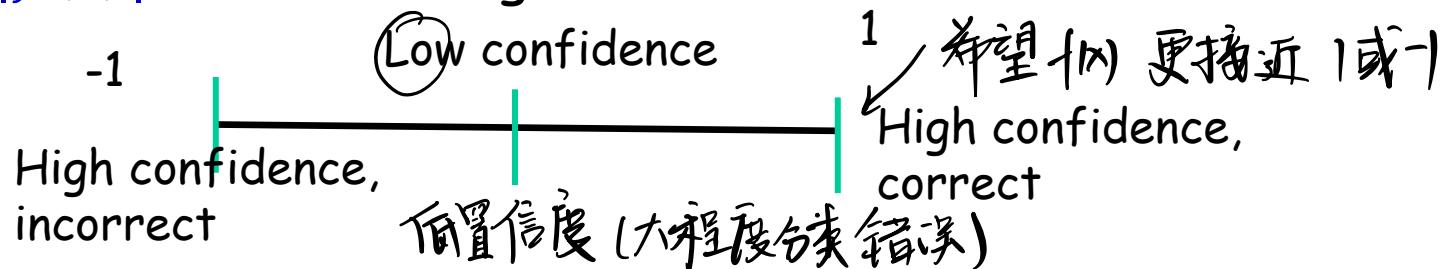
The majority vote rule  $H_f$  given by  $f$  (given by  $H_f = \text{sign}(f(x))$ ) predicts wrongly on example  $(x, y)$  iff  $yf(x) \leq 0$ .

Definition: **margin** of  $H_f$  (or of  $f$ ) on example  $(x, y)$  to be  $yf(x)$ .

$$yf(x) = y \sum_{t=1}^T [\alpha_t h_t(x)] = \sum_{t=1}^T [y \alpha_t h_t(x)] = \sum_{t:y=h_t(x)}^{\text{分类正确}} \alpha_t - \sum_{t:y \neq h_t(x)}^{\text{分类错误}} \alpha_t$$

The margin is positive iff  $y = H_f(x)$ .

See  $|yf(x)| = |f(x)|$  as the strength or the confidence of the vote.



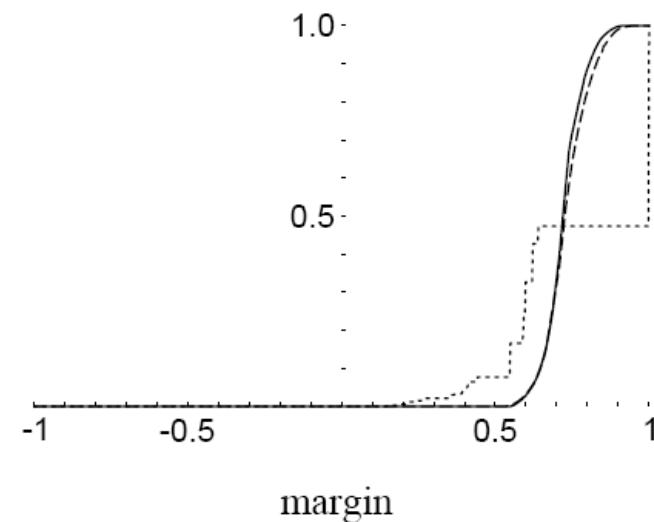
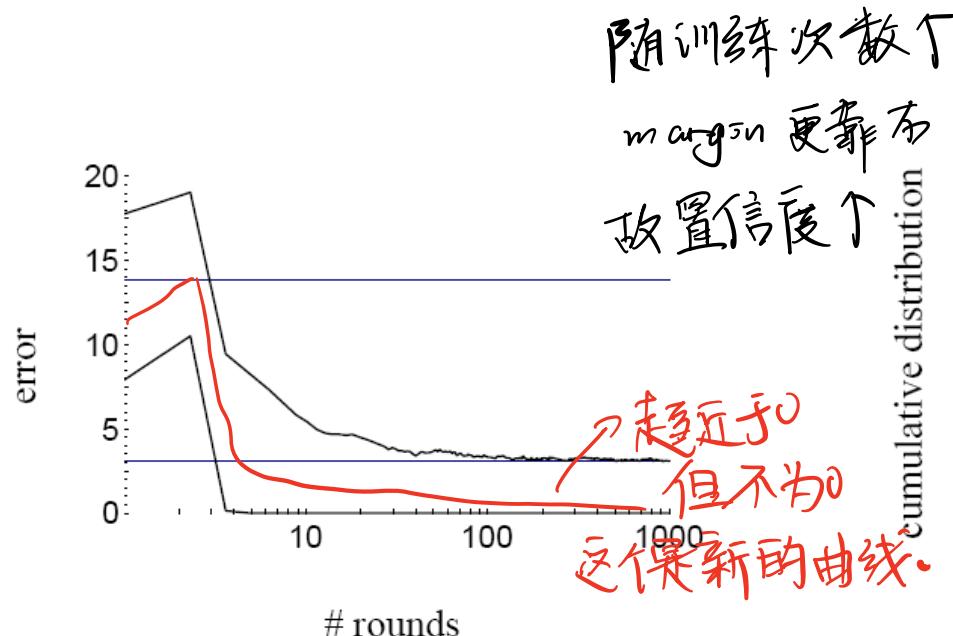
# Boosting and Margins

**Theorem:**  $\text{VCdim}(H) = d$ , then with prob.  $\geq 1 - \delta$ ,  $\forall f \in \text{co}(H)$ ,  $\forall \theta > 0$ ,

$$\Pr_D[yf(x) \leq 0] \leq \Pr_S[yf(x) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \sqrt{\frac{d \ln^2 \frac{m}{\delta}}{\theta^2} + \ln \frac{1}{\delta}}\right)$$

数据集  $\Rightarrow$  只保留线性部分(舍去非线性部分)  $\Rightarrow O(\sqrt{\frac{d}{m \theta^2}})$

**Note:** bound does **not** depend on  $T$  (the # of rounds of boosting), depends only on the complex. of the weak hyp space and the margin!



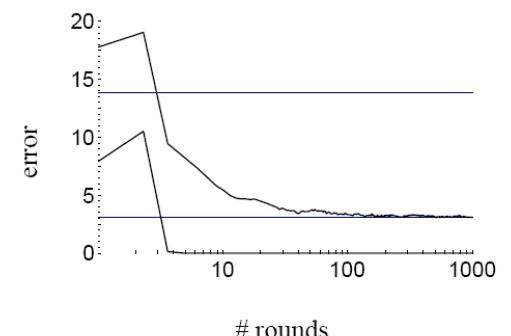
**Quiz:** according to this slide, explain why adaboost keeps decreasing testing error, even if training error equals to zero.

# Boosting and Margins

**Theorem:**  $\text{VCdim}(H) = d$ , then with prob.  $\geq 1 - \delta$ ,  $\forall f \in \text{co}(H)$ ,  $\forall \theta > 0$ ,

$$\Pr_D[yf(x) \leq 0] \leq \Pr_S[yf(x) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \sqrt{\frac{d \ln^2 \frac{m}{d}}{\theta^2} + \ln \frac{1}{\delta}}\right)$$

- If all training examples have **large margins**, then we can **approximate** the final classifier by a much smaller classifier.
- Can use this to prove that **better margin  $\rightarrow$  smaller test error**, regardless of the number of weak classifiers.
- Can also prove that **boosting tends to increase the margin** of training examples by concentrating on those of smallest margin.
- Although final classifier is getting **larger**, **margins are likely to be increasing**, so the final classifier is actually getting closer to a **simpler** classifier, driving **down** test error.

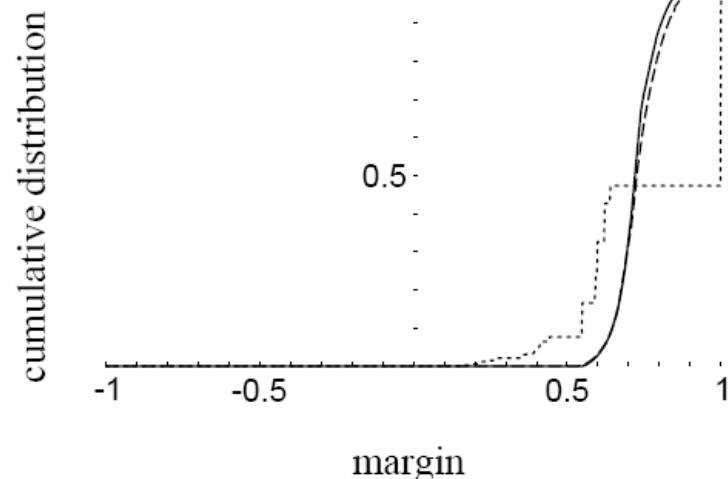
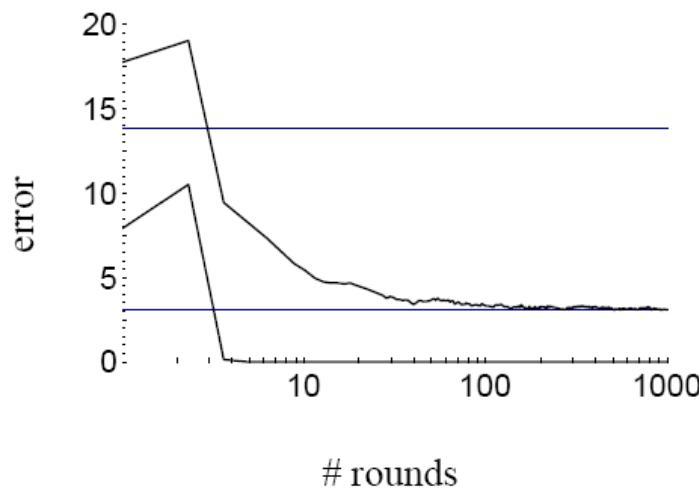


# Boosting and Margins

**Theorem:**  $\text{VCdim}(H) = d$ , then with prob.  $\geq 1 - \delta$ ,  $\forall f \in \text{co}(H)$ ,  $\forall \theta > 0$ ,

$$\Pr_D[yf(x) \leq 0] \leq \Pr_S[yf(x) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \sqrt{\frac{d \ln^2 \frac{m}{d}}{\theta^2} + \ln \frac{1}{\delta}}\right)$$

**Note:** bound does **not** depend on  $T$  (the # of rounds of boosting), depends only on the complex. of the weak hyp space and the margin!



# Boosting, Adaboost Summary

- Shift in mindset: goal is now just to find classifiers a bit better than random guessing.
- Backed up by solid foundations.
- Adaboost work and its variations well in practice with many kinds of data (one of the top 10 ML algos).
- More about classic applications in Recitation.
- Relevant for big data age: quickly focuses on “core difficulties”, so well-suited to distributed settings, where data must be communicated efficiently [Balcan-Blum-Fine-Mansour COLT’12].

Interestingly, the usefulness of margin  
recognized in Machine Learning since late 50's.

Perceptron [Rosenblatt'57] analyzed via geometric  
(aka  $L_2, L_2$ ) margin.

根据场景优化。  
在线学习。之前是Batch Learning。

Original guarantee in the online learning scenario.

# The Perceptron Algorithm

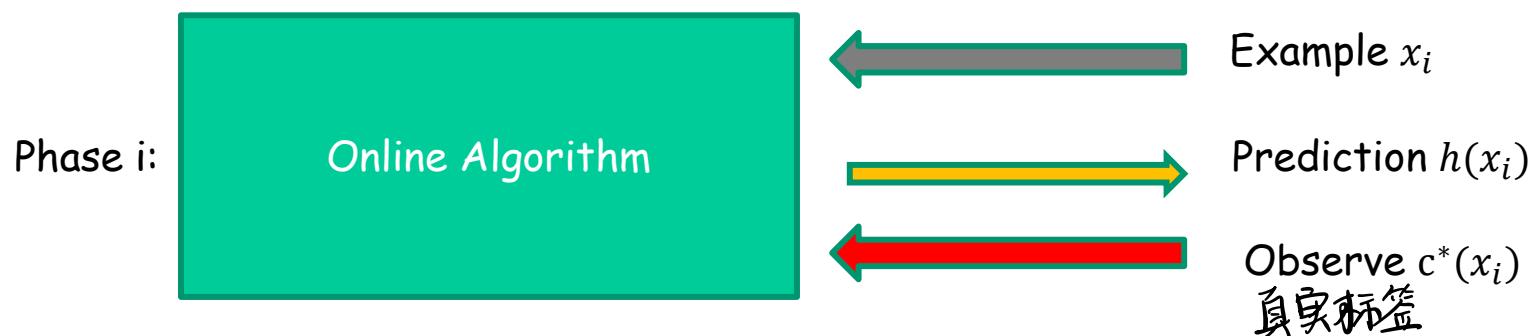
- Online Learning Model
- Margin Analysis
- Kernels

# The Online Learning Model

- Examples arrive **sequentially**.
- We need to make a prediction.

Afterwards observe the outcome.

For  $i=1, 2, \dots$ :



Mistake bound model

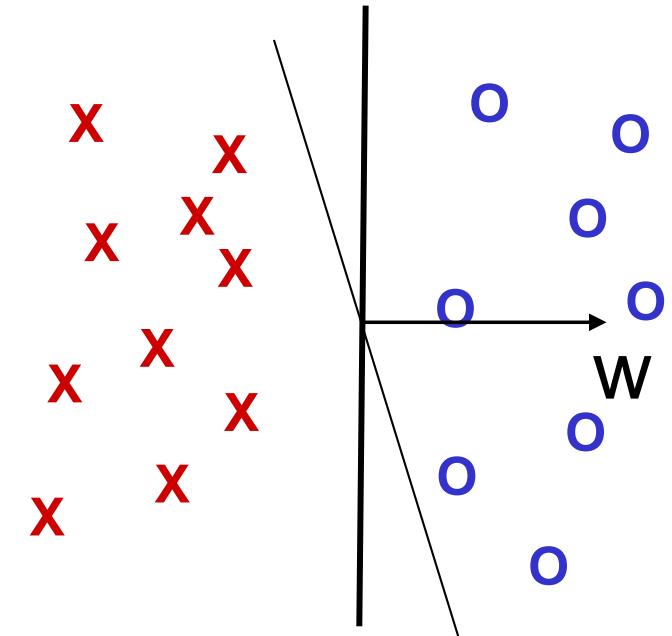
- Analysis wise, make **no distributional assumptions**.  
不对数据分布做假设  $\Rightarrow$  不能统计概率.
- Goal: **Minimize the number of mistakes.**  
最小化分类出错的次数.

# The Online Learning Model. Motivation

- Email classification (distribution of both spam and regular mail changes over time, but the target function stays fixed - last year's spam still looks like spam).
- Recommendation systems. Recommending movies, etc.
- Predicting whether a user will be interested in a new news article or not.
- Ad placement in a new market.

# Linear Separators

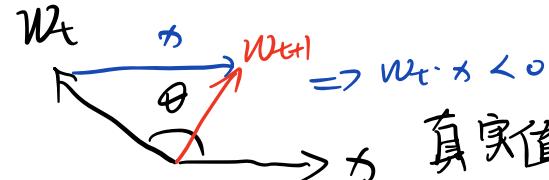
- Instance space  $X = \mathbb{R}^d$
- Hypothesis class of linear decision surfaces in  $\mathbb{R}^d$ .
- $w^T x + w_0$ :  $d$  维向量.
- $h(x) = w \cdot x + w_0$ , if  $h(x) \geq 0$ , then label  $x$  as +, otherwise label it as -



**Claim:** WLOG  $w_0 = 0$ .

**Proof:** Can simulate a non-zero threshold with a dummy input feature  $x_0$  that is always set up to 1.

- $x = (x_1, \dots, x_d) \rightarrow \tilde{x} = (x_1, \dots, x_d, 1)$
- $w \cdot x + w_0 \geq 0$  iff  $(w_1, \dots, w_d, w_0) \cdot \tilde{x} \geq 0$



where  $w = (w_1, \dots, w_d)$

真实值: >0 => mistake on positive  $w_{t+1} = w_t + s \Rightarrow w_{t+1} \cdot x > 0$

但经常是一次拉不回来，需多次迭代拉回来。

# Linear Separators: Perceptron Algorithm



- Set  $t=1$ , start with the all zero vector  $w_1$ .
- Given example  $x$ , predict positive iff  $w_t \cdot x \geq 0$
- On a mistake, update as follows: 如果没出左首:  $w_{t+1} \leftarrow w_t$ .
  - Mistake on positive, then update  $w_{t+1} \leftarrow w_t + x$
  - Mistake on negative, then update  $w_{t+1} \leftarrow w_t - x$   
或用  $w_{t+1} = w_t + y_t$ .  $y \in \{-1, 1\}$  是统合后的形式.

Note:  $w_t$  is weighted sum of incorrectly classified examples

$$w_t = a_{i_1}x_{i_1} + \cdots + a_{i_k}x_{i_k}$$

$$w_t \cdot x = a_{i_1}x_{i_1} \cdot x + \cdots + a_{i_k}x_{i_k} \cdot x$$

Important when we talk about kernels.



# Perceptron Algorithm: Example

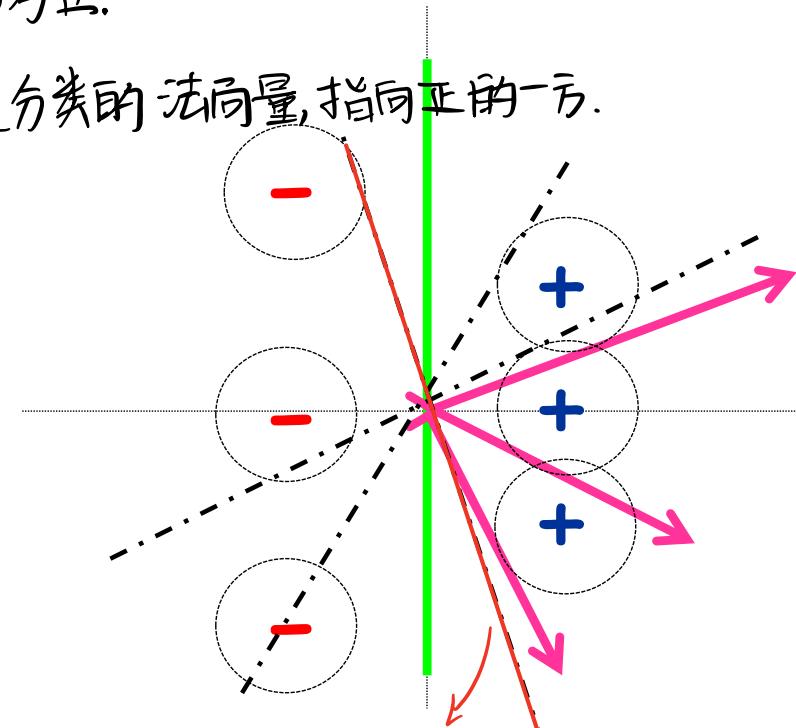
initialize  $w_i = (0,0)$  认为 0 的符号为正.

$x$	$y$	$\text{sign}(w^T x)$
(-1,2)	-	+
(1,0)	+	+
(1,1)	+	-
(-1,0)	-	-
(-1,-2)	-	+
(1,-1)	+	+

$w$  是分类的法向量, 指向正的一方.

Example:

$x$	$y$	$\text{sign}(w^T x)$
(-1,2)	-	+
(1,0)	+	+
(1,1)	+	-
(-1,0)	-	-
(-1,-2)	-	+
(1,-1)	+	+



## Algorithm:

- Set  $t=1$ , start with all-zeroes weight vector  $w_1$ .
- Given example  $x$ , predict positive iff  $w_t \cdot x \geq 0$ .
  - On a mistake, update as follows:
    - Mistake on positive, update  $w_{t+1} \leftarrow w_t + x$
    - Mistake on negative, update  $w_{t+1} \leftarrow w_t - x$

$$w_1 = (0,0)$$

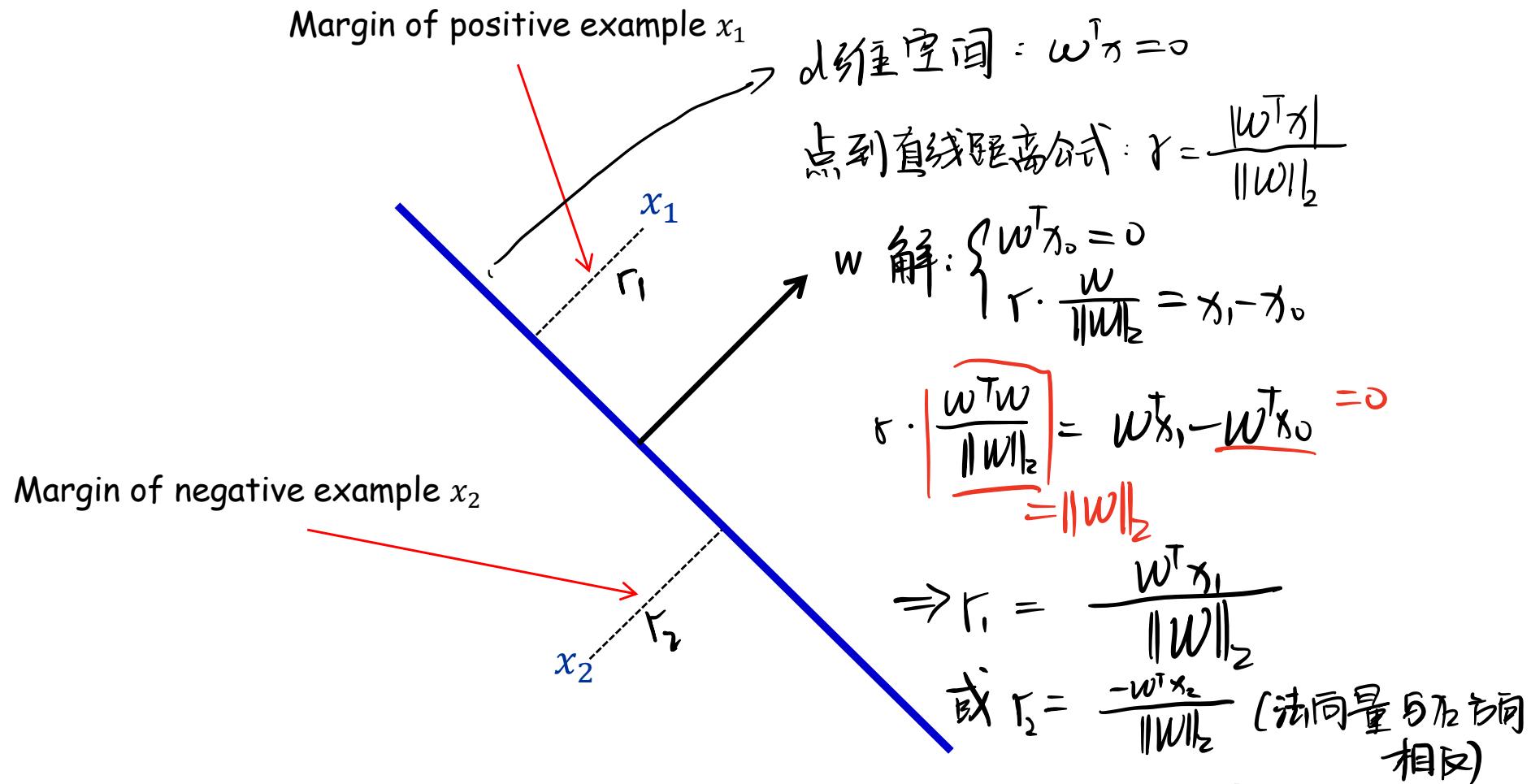
$$w_2 = w_1 - (-1,2) = (1,-2)$$

$$w_3 = w_2 + (1,1) = (2,-1)$$

$$w_4 = w_3 - (-1,-2) = (3,1)$$

# Geometric Margin

**Definition:** The margin of example  $x$  w.r.t. a linear sep.  $w$  is the distance from  $x$  to the plane  $w \cdot x = 0$  (or the negative if on wrong side)

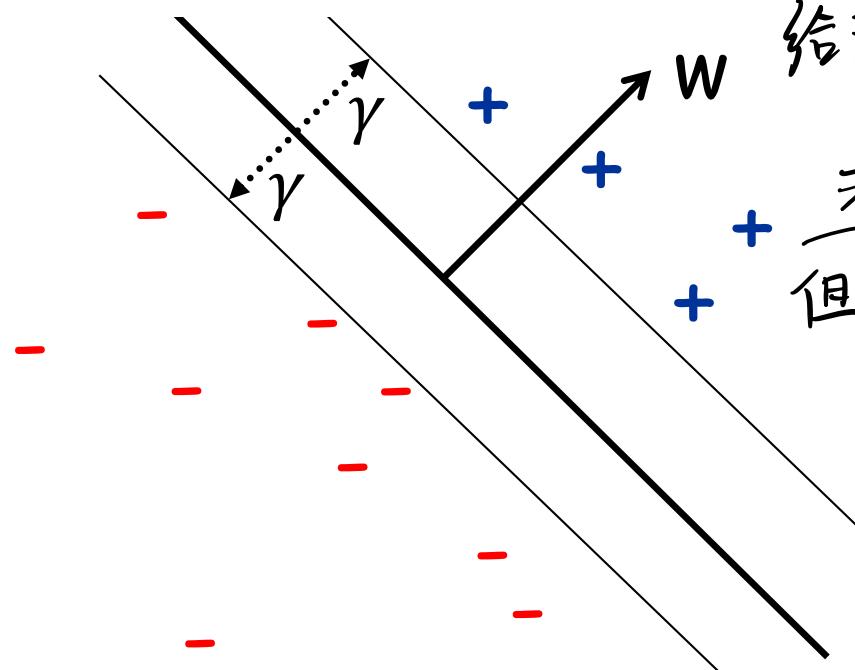


# Geometric Margin

**Definition:** The margin of example  $x$  w.r.t. a linear sep.  $w$  is the distance from  $x$  to the plane  $w \cdot x = 0$  (or the negative if on wrong side)

**Definition:** The margin  $\gamma_w$  of a set of examples  $S$  wrt a linear separator  $w$  is the smallest margin over points  $x \in S$ .

**Definition:** The margin  $\gamma$  of a set of examples  $S$  is the maximum  $\gamma_w$  over all linear separators  $w$ .



给定数据集,  $\gamma_w = \min_{x_i \in S} \frac{|w^T x_i|}{\|w\|_2}$

求最小间隔

但  $w$  会改变  $\Rightarrow \gamma = \max_w \gamma_w$

$\Rightarrow \gamma = \max_{w \in \mathbb{R}^d} \min_{x_i \in S} \frac{|w^T x_i|}{\|w\|_2}$

起到取绝对值作用

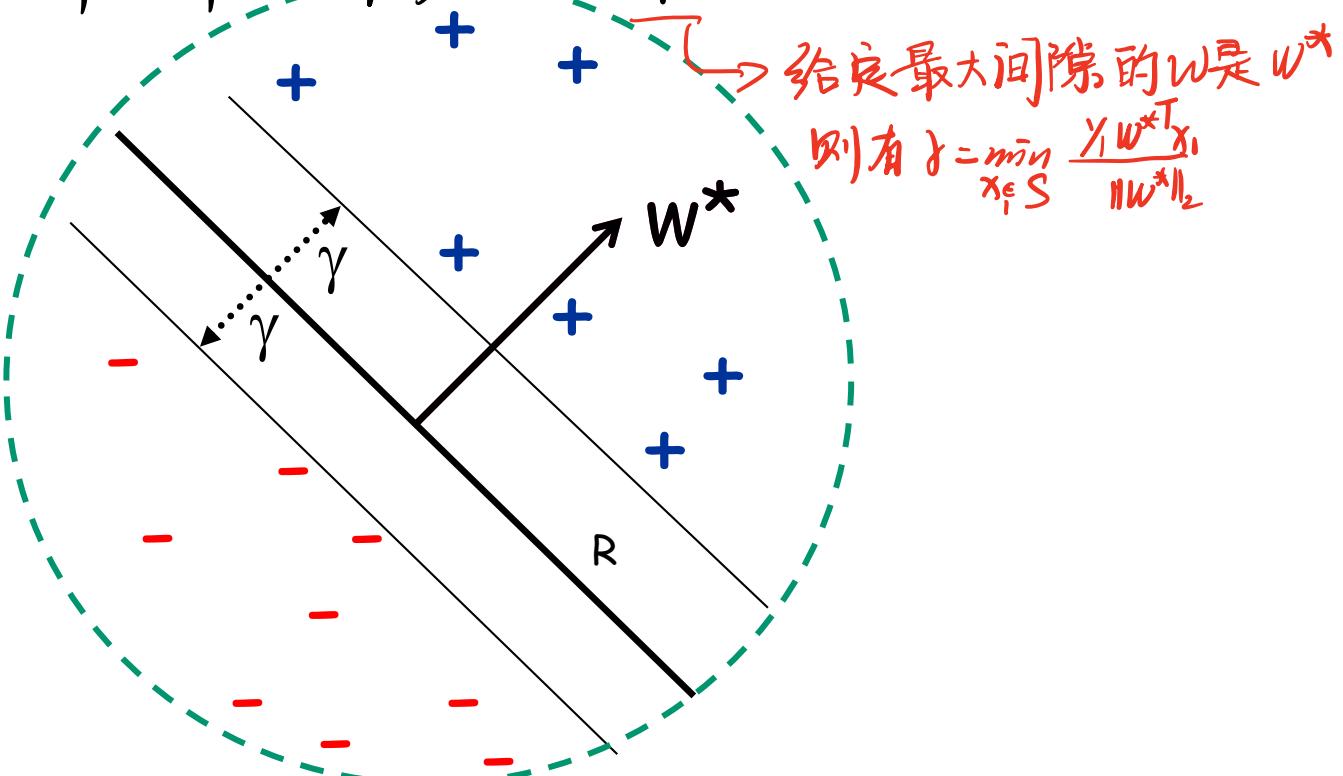
$y_i \in \{-1, 1\}$

# Perceptron: Mistake Bound

**Theorem:** If data has margin  $\gamma$  and all points inside a ball of radius  $R$ , then Perceptron makes  $\leq (R/\gamma)^2$  mistakes.

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes; algo is invariant to scaling.)

所有点在  $R$  半径球内 间隔  $\gamma$ , 则 错误最多不超过  $\frac{R}{\gamma}$  故找最大  $\gamma$



# Perceptron Algorithm: Analysis

**Theorem:** If data has margin  $\gamma$  and all points inside a ball of radius  $R$ , then Perceptron makes  $\leq (R/\gamma)^2$  mistakes.

**Proof:**

**Idea:** analyze  $w_t \cdot w^*$  and  $\|w_t\|$ , where  $w^*$  is the max-margin sep,  $\|w^*\| = 1$ .

**Claim 1:**  $w_{t+1} \cdot w^* \geq w_t \cdot w^* + \gamma$ . (because  $l(x)x \cdot w^* \geq \gamma$ )

**Claim 2:**  $\|w_{t+1}\|^2 \leq \|w_t\|^2 + R^2$ . (by Pythagorean Theorem)

$$\|w_t + yx\|^2 = \|w_t\|^2 + 2\|w_t yx\|_2 + \|yx\|^2 \leq \|w_t\|^2 + R^2$$

数据集范围

After  $M$  mistakes:  $\underbrace{< 0}_{(\text{预测错误})} \leq R^2$

$w_{M+1} \cdot w^* \geq \gamma M$  (by Claim 1)

$\|w_{M+1}\| \leq R\sqrt{M}$  (by Claim 2)

$w_{M+1} \cdot w^* \leq \|w_{M+1}\|$  (since  $w^*$  is unit length)

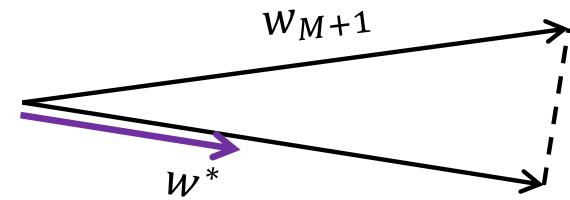
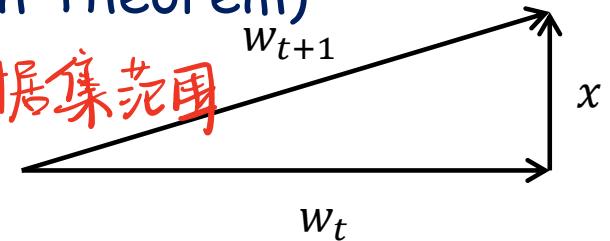
So,  $\gamma M \leq R\sqrt{M}$ , so  $M \leq \left(\frac{R}{\gamma}\right)^2$ .

Update rule:

- Mistake on positive:  $w_{t+1} \leftarrow w_t + x$
- Mistake on negative:  $w_{t+1} \leftarrow w_t - x$

$$w_{t+1} = w_t + yx \Rightarrow w_{t+1} \cdot w^* = w_t \cdot w^* + yxw^*$$

由于定义  $\gamma = \min \frac{y w^*}{\|w^*\|}$  故  $yxw^* \geq \gamma$



# Perceptron Extensions

- Can use it to find a consistent separator (by cycling through the data).



- One can convert the mistake bound guarantee into a distributional guarantee too (for the case where the  $x_i$ s come from a fixed distribution).

假设是线性可分  $\Rightarrow$  一定有一个一致的分类器，且循环不超过  $(\frac{R}{\gamma})^2$  次

→ 总错误数

- Can be adapted to the case where there is no perfect separator as long as the so called hinge loss (i.e., the total distance needed to move the points to classify them correctly large soft margin) is small. 若有少量误差，则若  $|W^T x_i| < \gamma$  则认为是一个正确的分类器。

- Can be kernelized to handle non-linear decision boundaries!

若线性不可分：使用内核化的方式拓展到高维空间使线性可分  
而映射回低维空间。

# Perceptron Discussion

- Simple online algorithm for learning linear separators with a nice guarantee that depends only on the geometric (aka  $L_2, L_2$ ) margin.
- It can be kernelized to handle non-linear decision boundaries --- see next class!
- Simple, but very useful in applications like Branch prediction; it also has interesting extensions to structured prediction.