



CSE 176 Introduction to Machine Learning

Lecture 15: Ensemble models

Some materials from Olga Veksler, Robin Dhamankar, Vandi Verma & Sebastian Thrun

What we have learnt so far...

- ❑ Machine learning models

- ❑ Bayesian classifier
 - ❑ K-nearest neighbor
 - ❑ Perceptron (Linear classifier)
 - ❑ Neural Network
 - ❑ Decision tree

- ❑ There are single models

Ensemble of multiple models

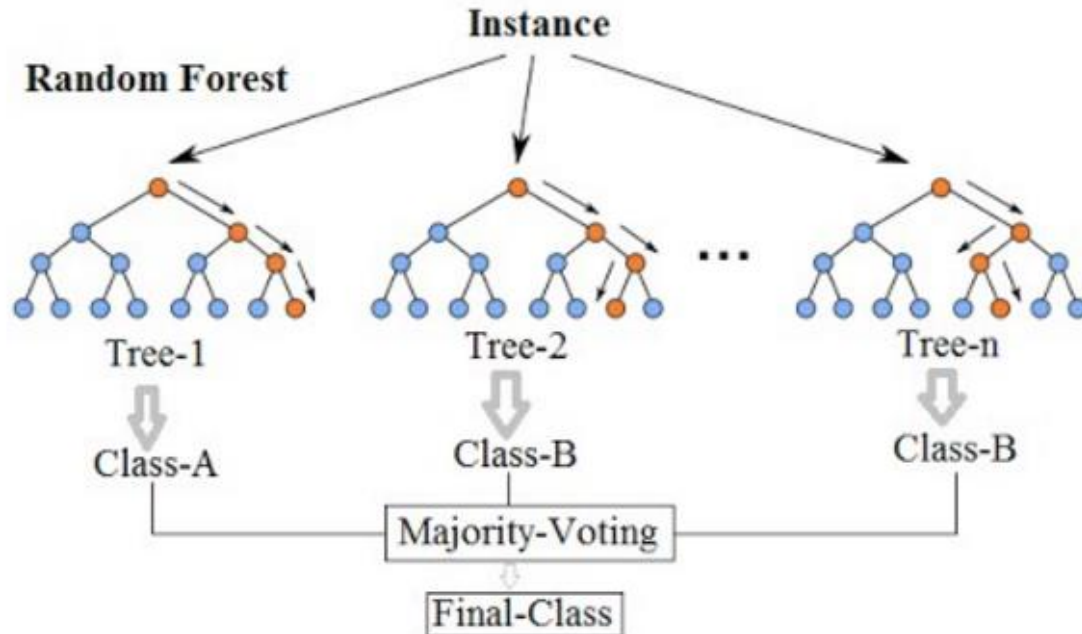
- ❑ Use multiple models, and “average” the predictions
- ❑ From statistics, it is good to average your predictions, reduces variance
- ❑ Consider L i.i.d. random variables y_1, \dots, y_L with expected value $E\{y_l\} = \mu$ and variance $\text{var}\{y_l\} = \sigma^2$

$$E\{y\} \stackrel{\text{pencil}}{=} E\left\{\frac{1}{L} \sum_{l=1}^L y_l\right\} = \frac{1}{L} \sum_{l=1}^L E\{y_l\} = \mu$$

$$\text{var}\{y\} \stackrel{\text{pencil}}{=} \text{var}\left\{\frac{1}{L} \sum_{l=1}^L y_l\right\} = \frac{1}{L^2} \text{var}\left\{\sum_{l=1}^L y_l\right\} = \frac{1}{L^2} \sum_{l=1}^L \text{var}\{y_l\} = \frac{1}{L} \sigma^2.$$

Example: Random forest

- ❑ Train an ensemble of L decision trees on L different subsets of the training set
- ❑ Define the ensemble output for a test instance as the majority vote (for classification) or the average (for regression) of the L trees

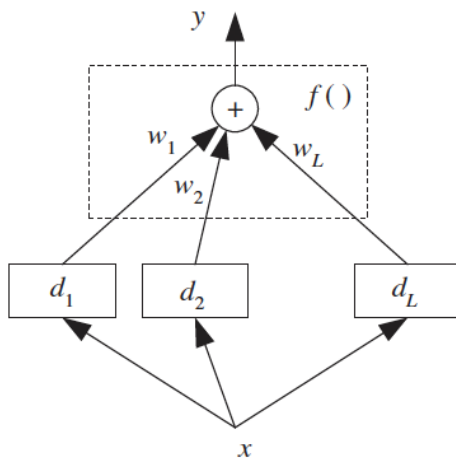


Mechanisms to generate diversity

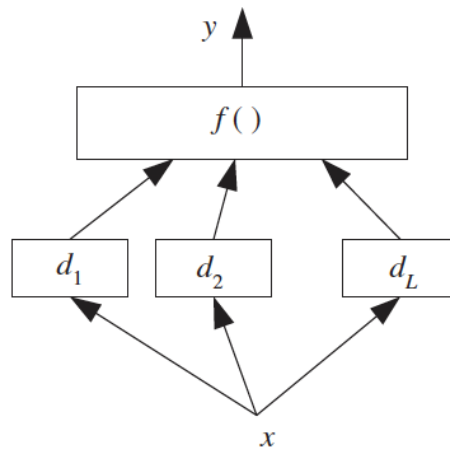
- ❑ Different models: linear, neural net, decision tree. . .
- ❑ Different hyperparameters
- ❑ Different optimization algorithm or initialization
- ❑ Different features
- ❑ Different training sets

Model combination scheme

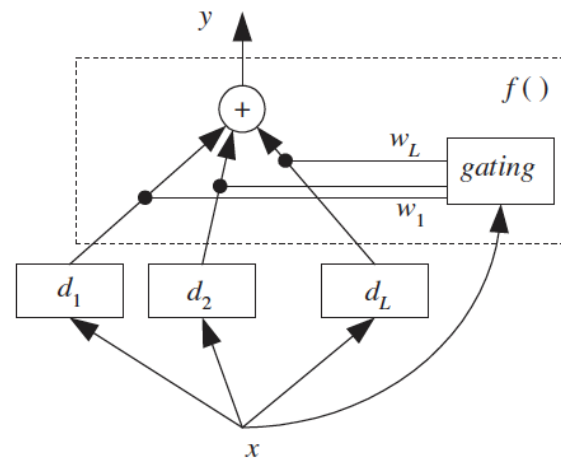
Voting or averaging



Stacking

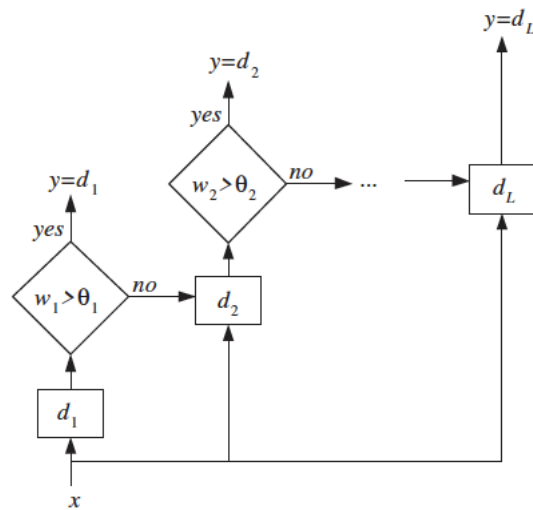


Mixture of experts



Model combination scheme

Boosting or cascading



Bagging

- ❑ We generate L (partly different) subsets of the training set
- ❑ We train L learners, each on a different subset
- ❑ The ensemble output is defined as the vote or average
- ❑ Random forest: a variation of bagging

Boosting

- ❑ Weak learner: a learner that has probability of error $< 1/2$ (i.e., better than random guessing on binary classification).
 - ❑ Ex: decision trees with only 1 or 2 levels.
- ❑ Strong learner: a learner that can have arbitrarily small probability of error.
 - ❑ Ex: neural net
- ❑ Boosting combines many weak learners to a strong learner

Ada Boost

- Assume 2-class problem, with labels +1 and -1
 - y^i in $\{-1, 1\}$

- Ada boost produces a discriminant function:

$$g(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) = \alpha_1 h_1(\mathbf{x}) + \alpha_2 h_2(\mathbf{x}) + \dots + \alpha_T h_T(\mathbf{x})$$

- Where $h_t(\mathbf{x})$ is a weak classifier, for example:

$$h_t(\mathbf{x}) = \begin{cases} -1 & \text{if email has word "money"} \\ 1 & \text{if email does not have word "money"} \end{cases}$$

- The final classifier is the sign of the discriminant function

$$f_{\text{final}}(\mathbf{x}) = \text{sign}[g(\mathbf{x})]$$

Idea Behind Ada Boost

- Maintains distribution of weights over the training examples
- Initially weights are equal
- Main Idea: at successive iterations, the weight of misclassified examples is increased



Idea Behind Ada Boost

- Examples of high weight are shown more often at later rounds
- Face/nonface classification problem:

Round 1

best weak classifier:

change weights:

						
1/7	1/7	1/7	1/7	1/7	1/7	1/7
✓	✗	✓	✓	✗	✓	✗
1/16	1/4	1/16	1/16	1/4	1/16	1/4

Round 2

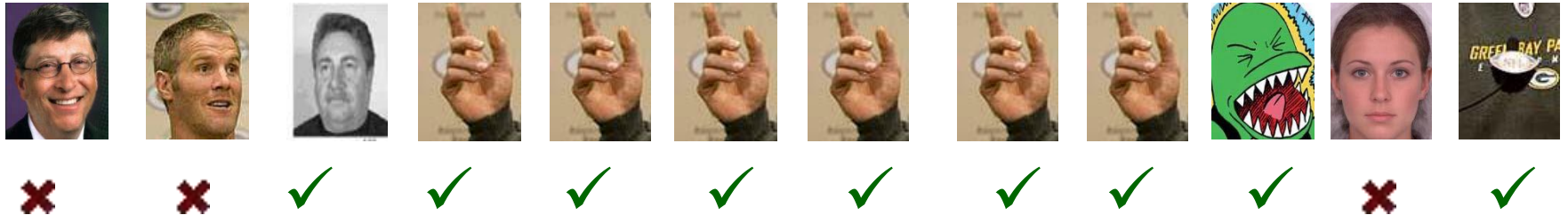
best weak classifier:

change weights:

									
✓	✓	✓	✗	✗	✗	✓	✓	✓	✓
	1/8	1/32	11/32		1/2		1/8	1/32	1/32

Idea Behind Ada Boost

Round 3




- out of all available weak classifiers, we choose the one that works best on the data we have at round 3

Idea Behind Ada Boost

Round 3



- out of all available weak classifiers, we choose the one that works best on the data we have at round 3
- we assume there is always a weak classifier better than random (better than 50% error)
-  image is half of the data given to the classifier
- chosen weak classifier **has to** classify this image correctly

More Comments on Ada Boost

- Ada boost is simple to implement, provided you have an implementation of a “weak learner”

More Comments on Ada Boost

- Ada boost is simple to implement, provided you have an implementation of a “weak learner”
- Will work as long as the “basic” classifier $h_t(\mathbf{x})$ is at least slightly better than random
 - will work if the error rate of $h_t(\mathbf{x})$ is less than 0.5
 - 0.5 is the error rate of a random guessing for 2-class problem
- Can be applied to boost any classifier, not necessarily weak
 - but there may be no benefits in boosting a “strong” classifier

Ada Boost for 2 Classes

Initialization step: for each example \mathbf{x} , set
$$\mathbf{D}(\mathbf{x}) = \frac{1}{N}$$
, where N is the number of examples

Iteration step (for $t = 1 \dots T$):

1. Find best weak classifier $h_t(\mathbf{x})$ using weights $\mathbf{D}(\mathbf{x})$
2. Compute the error rate ϵ_t as

$$\epsilon_t = \sum_{i=1}^N \mathbf{D}(\mathbf{x}^i) \cdot \mathbb{I}[y^i \neq h_t(\mathbf{x}^i)]$$

3. compute weight α_t of classifier h_t

$$\alpha_t = \log\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)$$

4. For each \mathbf{x}^i , $\mathbf{D}(\mathbf{x}^i) = \mathbf{D}(\mathbf{x}^i) \cdot \exp(\alpha_t \cdot \mathbb{I}[y^i \neq h_t(\mathbf{x}^i)])$

5. Normalize $\mathbf{D}(\mathbf{x}^i)$ so that
$$\sum_{i=1}^N \mathbf{D}(\mathbf{x}^i) = 1$$

$$\mathbf{f}_{\text{final}}(\mathbf{x}) = \text{sign} \left[\sum \alpha_t h_t(\mathbf{x}) \right]$$

Ada Boost: Step 1

1. Find best weak classifier $h_t(\mathbf{x})$ using weights $\mathbf{D}(\mathbf{x})$
 - some classifiers accept weighted samples, not all
 - if classifier does not take weighted samples, sample from the training samples according to the distribution $\mathbf{D}(\mathbf{x})$



1/16



1/4



1/16



1/16



1/4



1/16



1/4

Ada Boost: Step 1

1. Find best weak classifier $h_t(\mathbf{x})$ using weights $\mathbf{D}(\mathbf{x})$

- some classifiers accept weighted samples, not all
- if classifier does not take weighted samples, sample from the training samples according to the distribution $\mathbf{D}(\mathbf{x})$



1/16



1/4



1/16



1/16



1/4



1/16



1/4

- Draw k samples, each \mathbf{x} with probability equal to $\mathbf{D}(\mathbf{x})$:



re-sampled examples

Ada Boost: Step 1

1. Find best weak classifier $h_t(\mathbf{x})$ using weights $D(\mathbf{x})$
 - Give to the classifier the re-sampled examples:



Ada Boost: Step 1

1. Find best weak classifier $h_t(x)$ using weights $D(x)$

- Give to the classifier the re-sampled examples:



- To find the best weak classifier, go through **all** weak classifiers, and find the one that gives the smallest error on the re-sampled examples

weak classifiers	$h_1(x)$	$h_2(x)$	$h_3(x)$	$h_m(x)$
errors:	0.46	0.36	0.16		0.43

the best classifier $h_t(x)$
to choose at iteration t

Ada Boost: Step 2

2. Compute ϵ_t the error rate as

$$\epsilon_t = \sum_{i=1}^N D(x^i) \cdot I[y^i \neq h_t(x^i)]$$



1/16



1/4



1/16



1/16



1/4



1/16



1/4



$$\epsilon_t = \frac{1}{4} + \frac{1}{16} = \frac{5}{16}$$

- ϵ_t is the weight of all misclassified examples added
 - the error rate is computed over original examples, not the re-sampled examples
- If a weak classifier is better than random, then $\epsilon_t < 1/2$

Ada Boost: Step 3

3. compute weight α_t of classifier h_t

$$\alpha_t = \log((1 - \epsilon_t) / \epsilon_t)$$

In example from previous slide:

$$\epsilon_t = \frac{5}{16} \Rightarrow \alpha_t = \log \frac{1 - \frac{5}{16}}{\frac{5}{16}} = \log \frac{11}{5} \approx 0.8$$

Ada Boost: Step 3

3. compute weight α_t of classifier h_t

$$\alpha_t = \log((1 - \epsilon_t) / \epsilon_t)$$

In example from previous slide:

$$\epsilon_t = \frac{5}{16} \Rightarrow \alpha_t = \log \frac{1 - \frac{5}{16}}{\frac{5}{16}} = \log \frac{11}{5} \approx 0.8$$

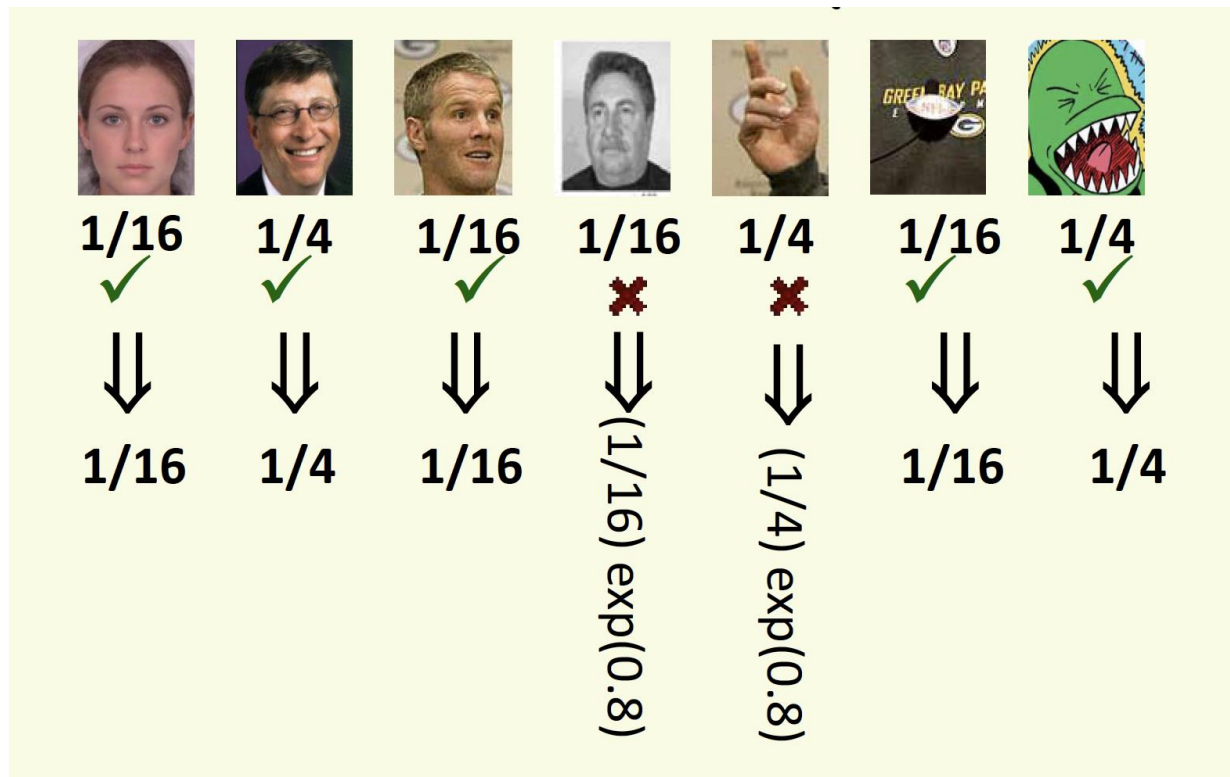
- Recall that $\epsilon_t < 1/2$
- Thus $(1 - \epsilon_t) / \epsilon_t > 1 \Rightarrow \alpha_t > 0$
- The smaller is ϵ_t , the larger is α_t , and thus the more importance (weight) classifier $h_t(x)$

$$\text{final}(\mathbf{x}) = \text{sign} \left[\sum \alpha_t h_t(\mathbf{x}) \right]$$

Ada Boost: Step 4

4. For each \mathbf{x}^i , $\mathbf{D}(\mathbf{x}^i) = \mathbf{D}(\mathbf{x}^i) \cdot \exp(\alpha_t \cdot \mathbf{I}[y^i \neq h_t(\mathbf{x}^i)])$

from previous slide $\alpha_t = 0.8$



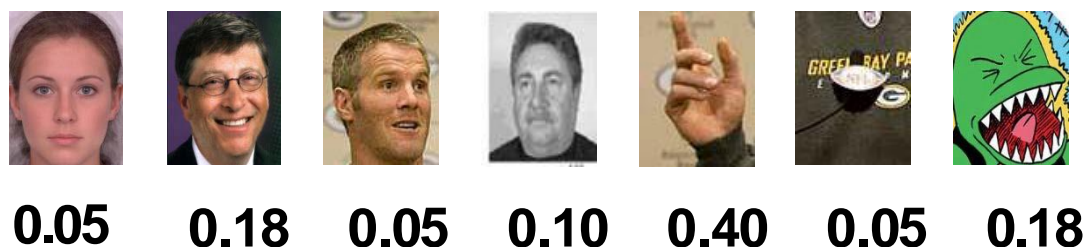
- weight of misclassified examples is increased

Ada Boost: Step 5

5. Normalize $D(x^i)$ so that $\sum D(x^i) = 1$
from previous slide:

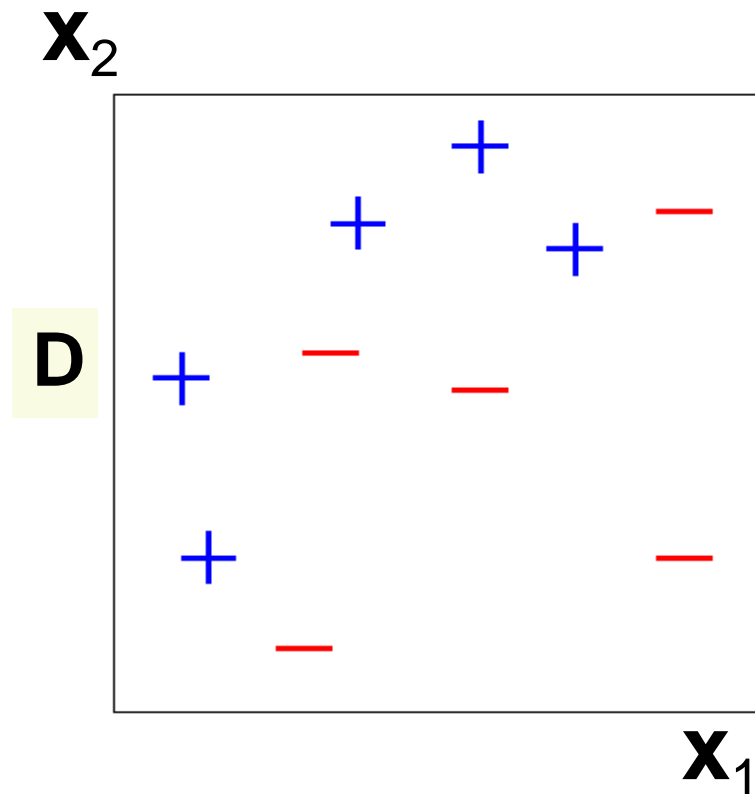


- after normalization



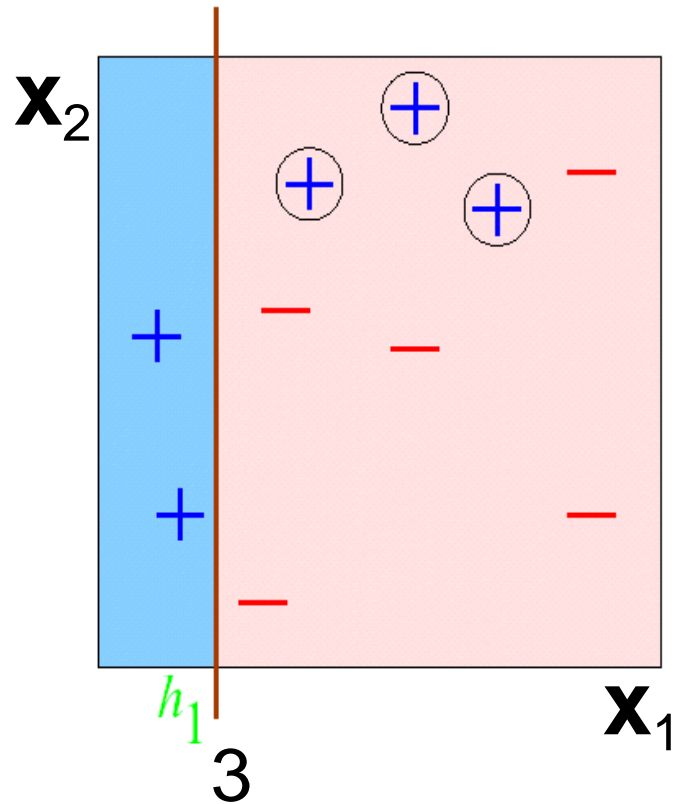
AdaBoost Example

- Initialization: all examples have equal weights



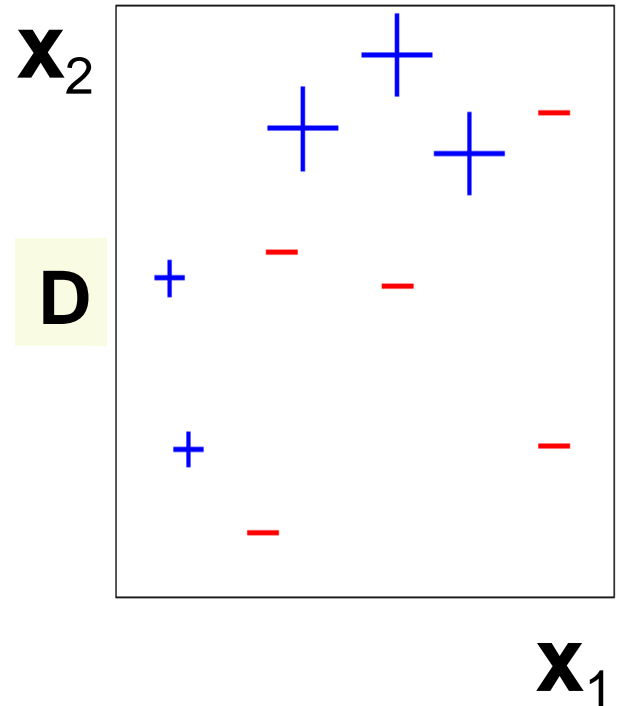
from “A Tutorial on Boosting” by Yoav Freund and Rob Schapire

AdaBoost Example: Round 1

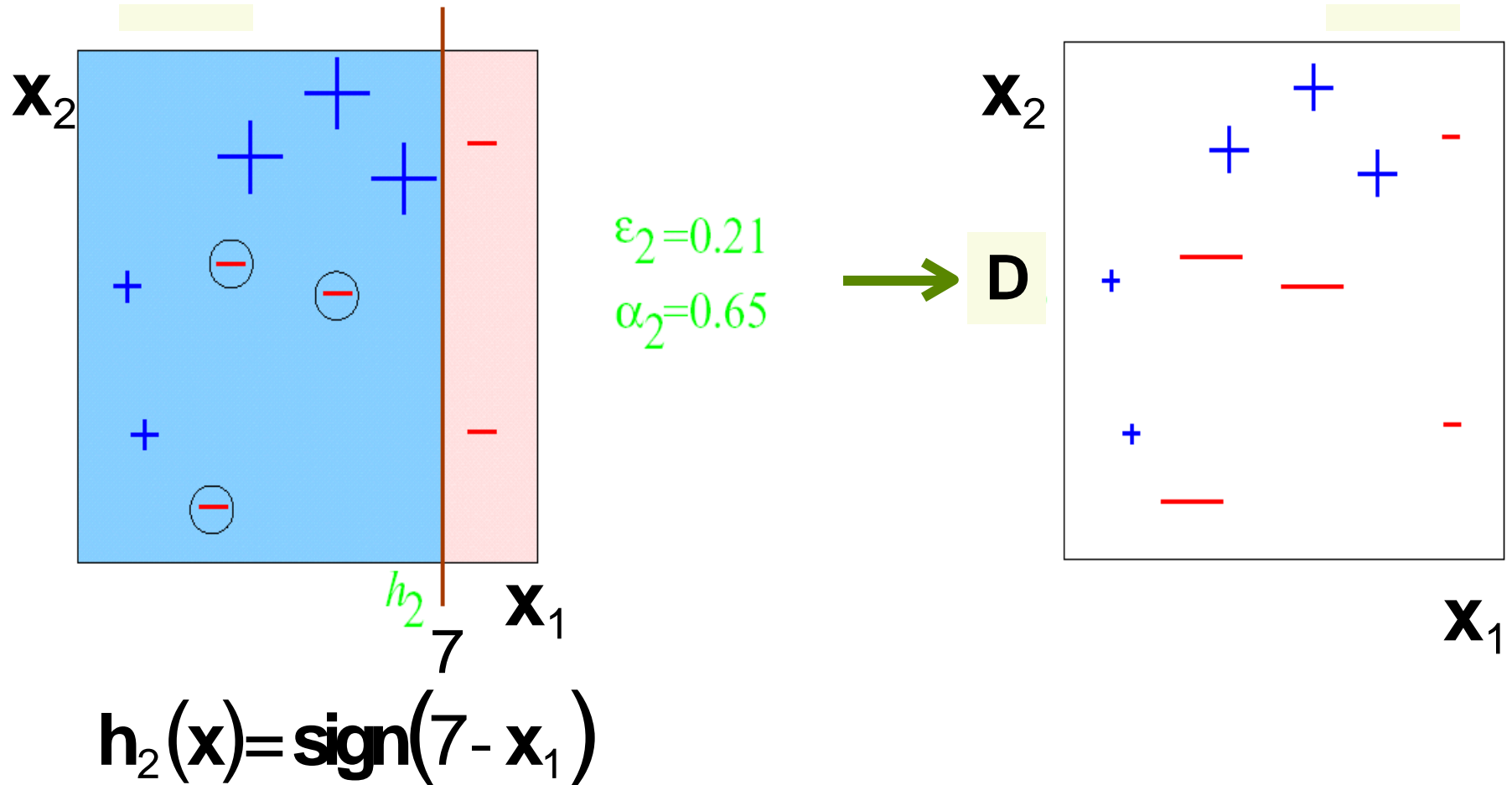


$$h_1(\mathbf{x}) = \text{sign}(3 - x_1)$$

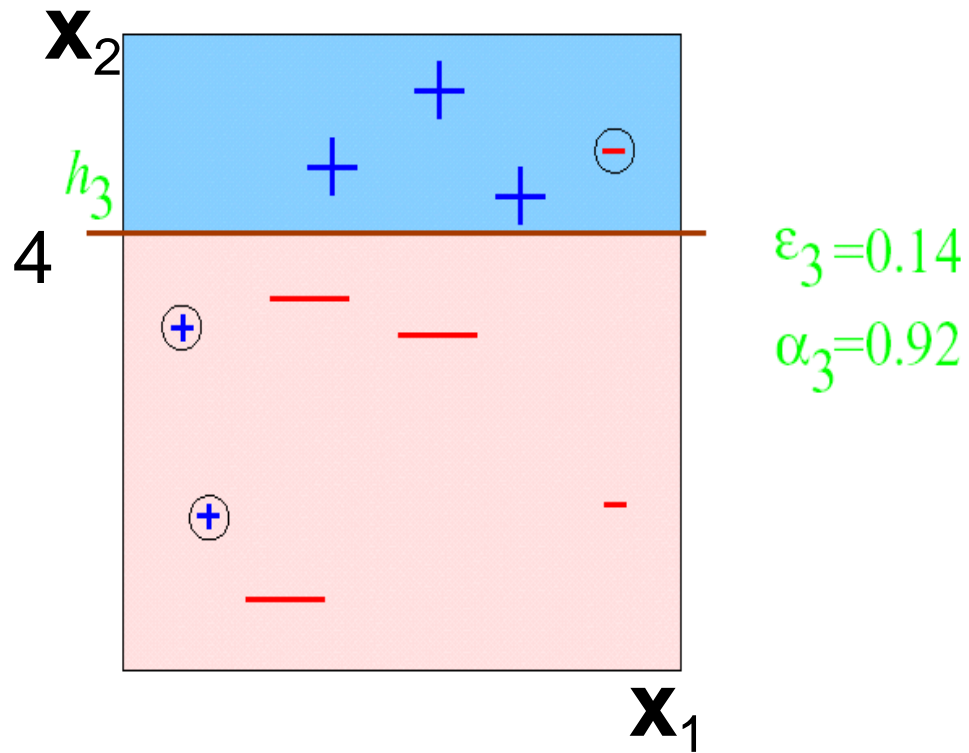
$$\begin{aligned}\epsilon_1 &= 0.30 \\ \alpha_1 &= 0.42\end{aligned}$$



AdaBoost Example: Round 2



AdaBoost Example: Round 3

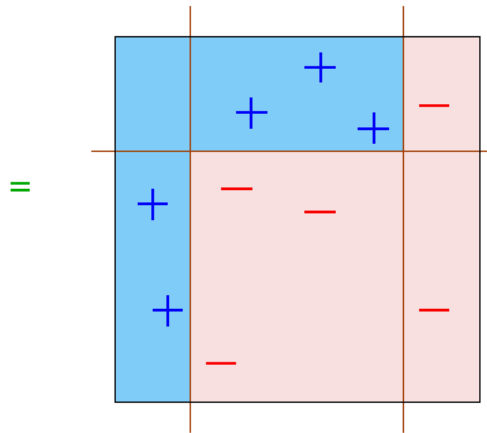


$$h_3(\mathbf{x}) = \text{sign}(x_2 - 4)$$

AdaBoost Example

$$\mathbf{f}_{\text{final}}(\mathbf{x}) =$$

$$\text{sign} \left(0.42 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} + 0.65 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} + 0.92 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} \right)$$



$$\mathbf{f}_{\text{final}}(\mathbf{x}) = \text{sign} \left(0.42 \text{sign}(3 - \mathbf{x}_1) + 0.65 \text{sign}(7 - \mathbf{x}_1) + 0.92 \text{sign}(\mathbf{x}_2 - 4) \right)$$

- Decision boundary non-linear

Practical Advantages of AdaBoost

- Can construct arbitrarily complex decision regions
- Fast and Simple
- Has only one parameter to tune, T
- Flexible: can be combined with any classifier
- provably effective (assuming weak learner)
 - shift in mind set: goal now is merely to find hypotheses that are better than random guessing