

Ensemble learning



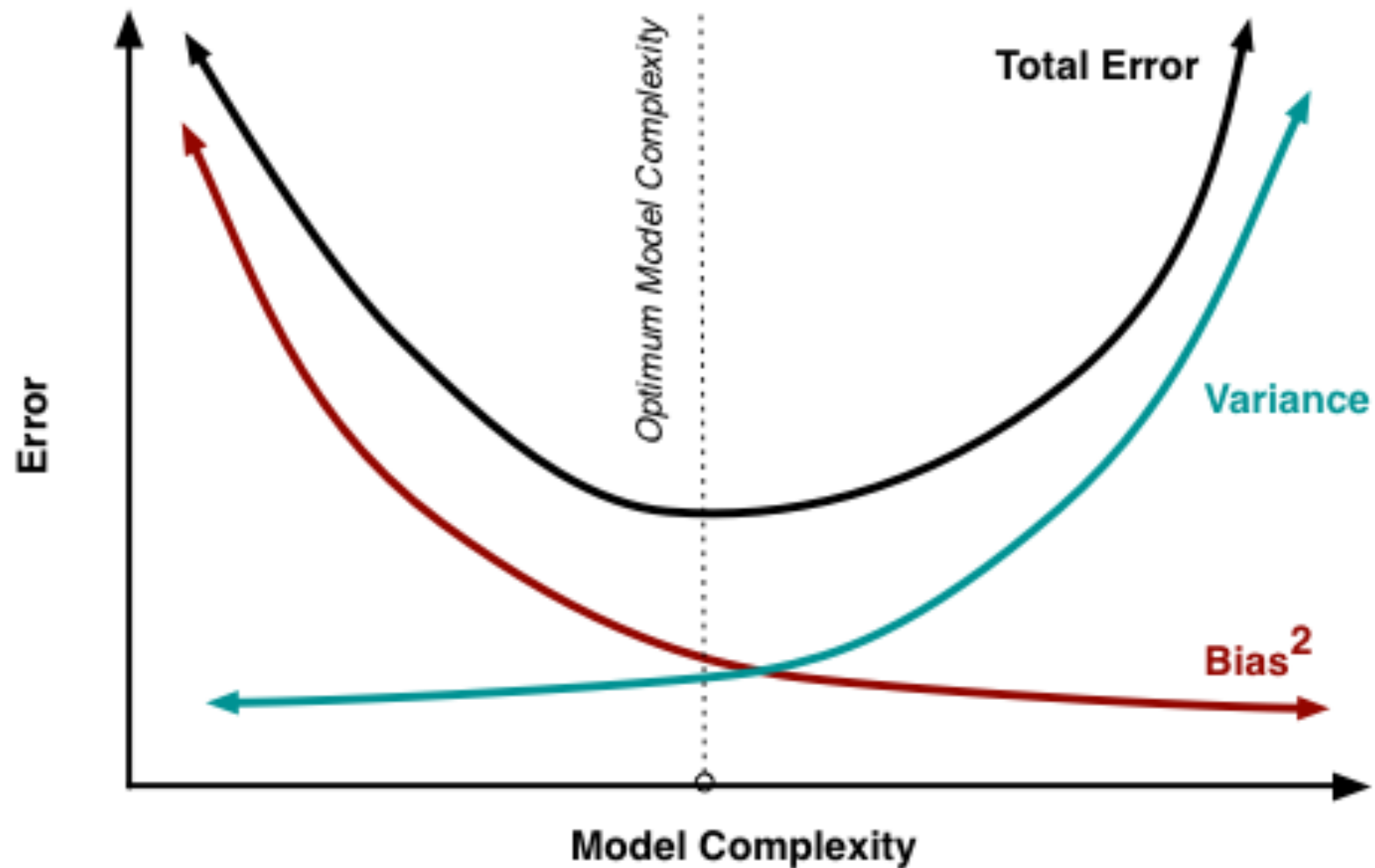
Reminders/Comments

- Midterm marks out, average about 65%
- For final, I will give slightly less weight to bonus questions
- The letter grades in eclass are not accurate (just show what they were for last year)
 - each year I adjust it based on performance in the course
 - this year, marking scheme has changed, I have no doubt the letter grade assignment to percentage will change
- Initial drafts due today

Reasons for lost marks

- The answer was wrong
- The answer was sort of for a different question, or misunderstood the question
- The answer was unclear and/or stated incorrect things
 - I could try to guess what you meant, and sometimes I do. But, at the same time, I have to evaluate what you actually wrote, because otherwise I add a lot of bias
- I did not give partial marks for just writing something, if the answer was wrong
 - Giving partial marks for wrong solutions just artificially inflates scores, and doesn't really fix penalties for having unclear solutions
 - In the end, I scale everything anyway, but I scale something that more accurately reflects what was given

How did everyone draw this picture so well?



- To clarify, bias-variance is about (a) generalization error or (b) about getting estimated parameters to be closer to true parameters, in terms of mean-squared error

Q3 was the most difficult

- $Y = w_0 + w_1 X_1 + w_2 X_2 + \text{epsilon}$, with $w_1 = 0$
- What do we know about $p(y \mid x_1, x_2)$ and $p(y \mid x_1)$
- I wanted you to think about independence and conditional independence

Bonus question

- Do not have a fixed batch of data, rather have a constant stream of data (e.g., imagine you are Google), but still want to normalize the data (either by min and max, or by mean and stddev)
- I was looking for you to state that you could keep a running mean or variance, possibly with an exponential weighting to prefer newer data more
- Why reasonable and why deficient?

Topic today: Collections of models

- Have mostly discussed learning one single “best” model
 - best linear regression model
 - best neural network model
- Can we take advantage of multiple learned models?

Rationale

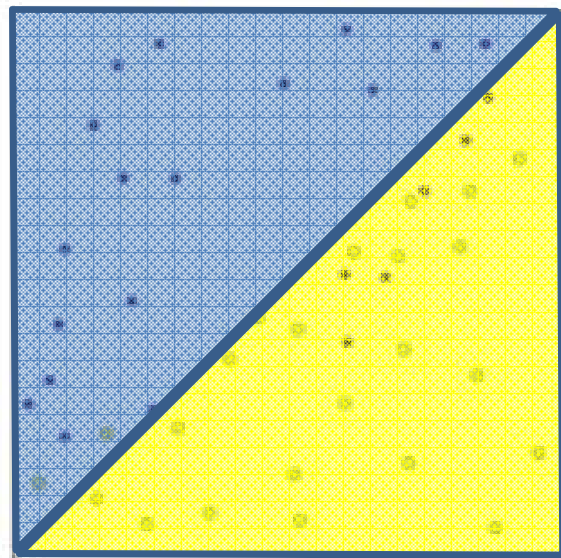
- There is no algorithm that is always the most accurate
- Different learners can use different
 - Algorithms (e.g., logistic regression or SVMs)
 - Parameters (e.g., regularization parameters)
 - Representations (e.g., polynomial basis or kernels)
 - Training sets (e.g., two different random subsamples of data)
- The problem: how to combine them

Ensembles

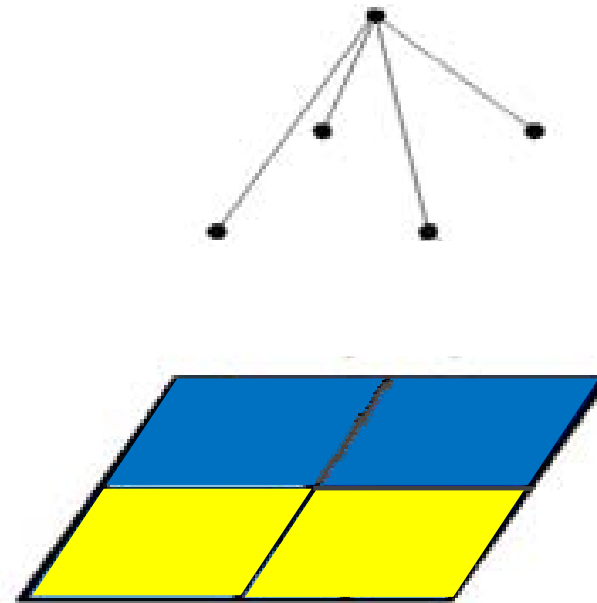
- Can a set of **weak learners** create a single strong learner?
- Answer: yes! See seminal paper: “The Strength of Weak Learnability” Schapire, 1990
- Why do we care?
 - can be easier to specify weak learners e.g., shallow decision trees, set of neural networks with smaller number of layers, etc.
 - fighting the bias-variance trade-off

Weak learners

- **Weak learners:** naive Bayes, logistic regression, decision stumps (or shallow decision trees)

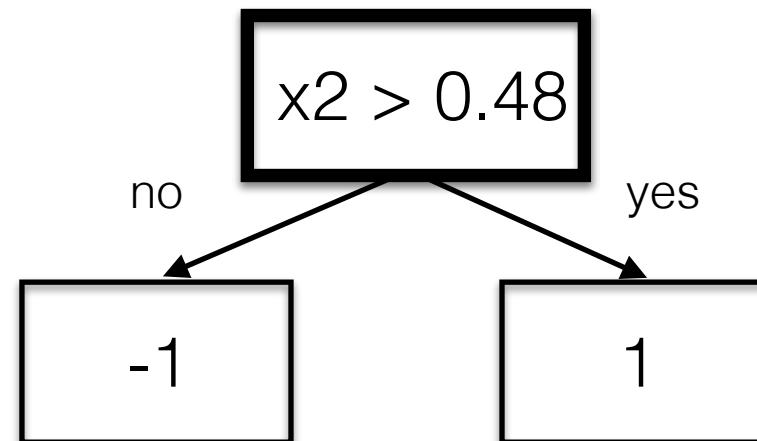


logistic regression



decision stump

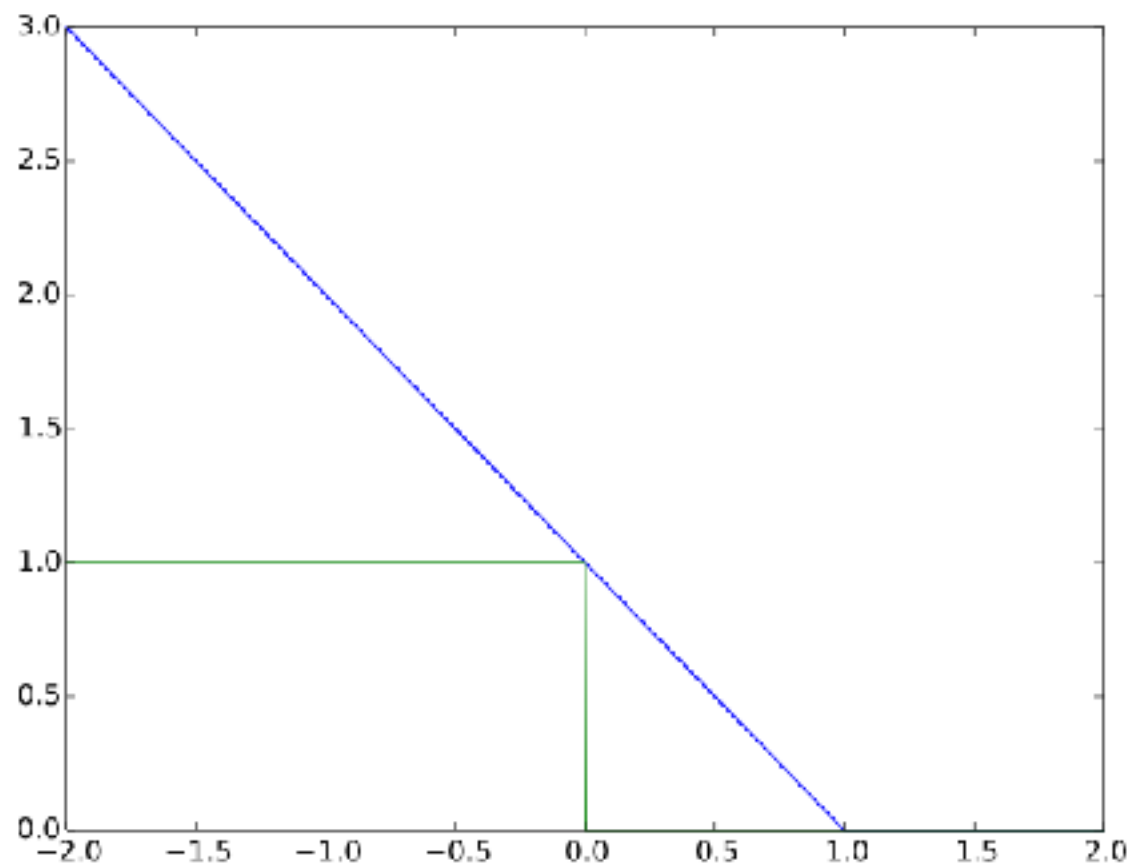
Example of a decision stump



Decision tree provides more splits;
decision stump is a one level decision tree

How learn signed prediction?

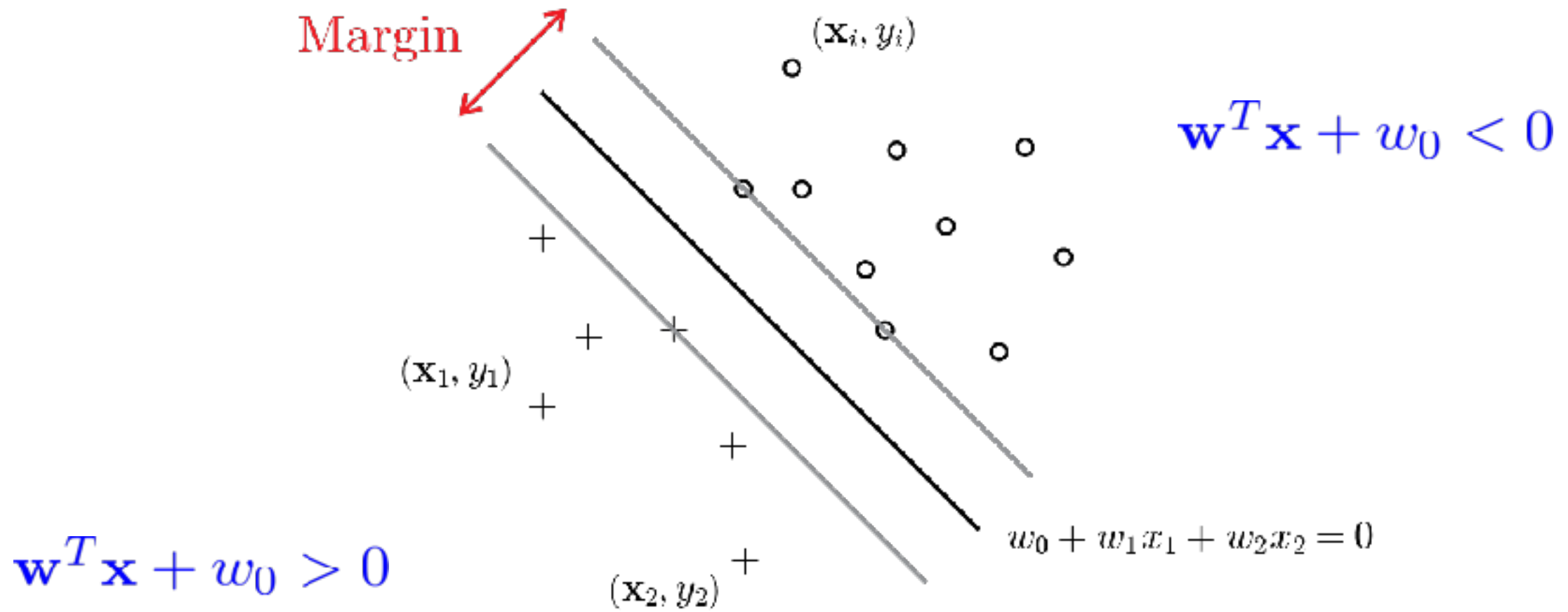
- Decision-stump outputs $\text{sign}(\langle x, w \rangle)$
- Logistic regression and linear regression: take learned w , and prediction is set to $\text{sign}(\langle x, w \rangle)$
- Support vector machines: minimize hinge loss



Green is zero-one loss
Blue is hinge loss

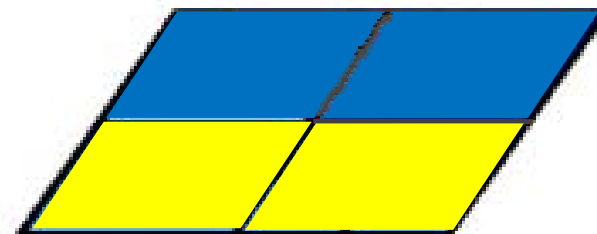
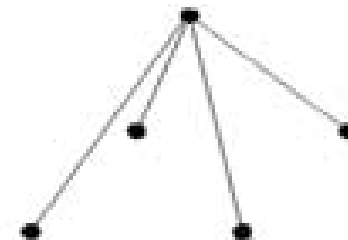
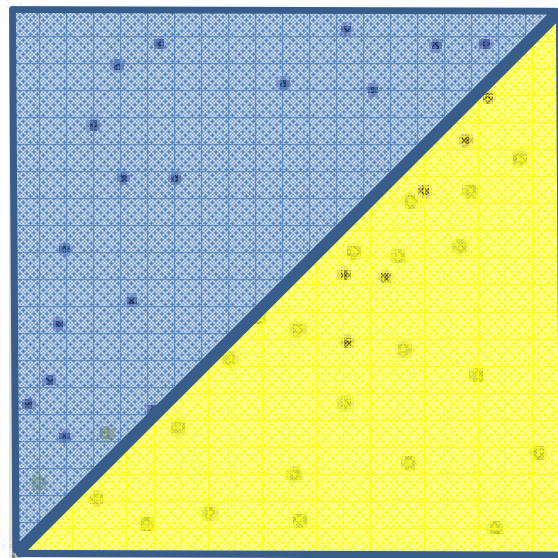
A side point about SVMs

- Support vector machine: minimize hinge loss while also adding goal to maximize the margin



Weak learners

- **Weak learners:** naive Bayes, logistic regression, decision stumps (or shallow decision trees)



Are good 😊 - Low variance, don't usually overfit

Are bad ☹️ - High bias, can't solve hard learning problems

Bias-variance tradeoff

- We encountered this trade-off for weights in linear regression
- Regularizing introduced bias, but reduced variance

$$\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \left(\text{Bias}(\hat{\theta}, \theta) \right)^2.$$

- More generally, when picking functions

$$\text{Var}(\hat{f}) + \text{Bias}(\hat{f}, f)^2$$

How might you specify bias between functions?

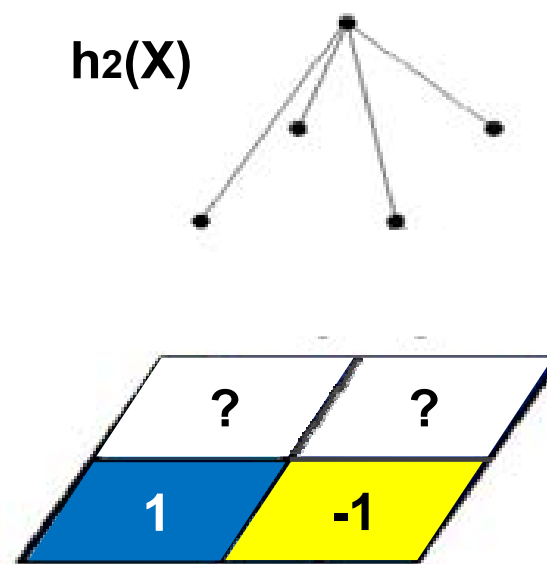
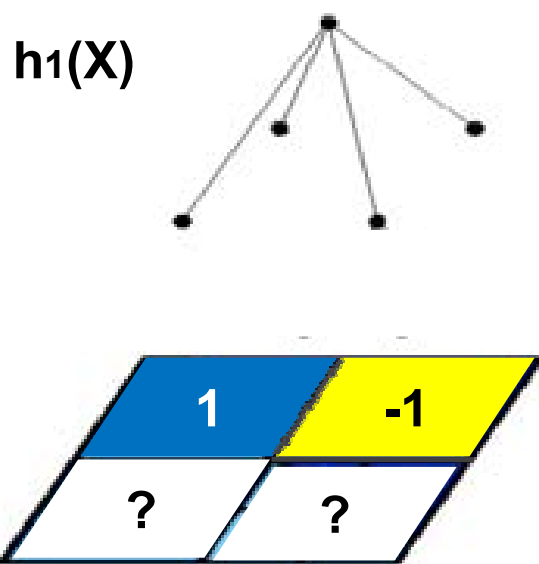
Voting (Ensemble methods)

- Instead of learning a single (weak) classifier, learn **many weak classifiers** that are **good at different parts of the input space**

- **Output class:** (Weighted) vote of each classifier

Classifiers that are most “sure” will vote with more conviction

Classifiers will be most “sure” about a particular part of the space. On average, do better than single classifier!



$$H: X \rightarrow Y (-1,1)$$

$$H(X) = h_1(X) + h_2(X)$$

$$H(X) = \text{sign}\left(\sum_t \alpha_t h_t(X)\right)$$

↓
weights

Voting (Ensemble methods)

- Instead of learning a single (weak) classifier, learn **many weak classifiers** that are **good at different parts of the input space**
- **Output class:** (Weighted) vote of each classifier
 - Classifiers that are most “sure” will vote with more conviction
 - Classifiers will be most “sure” about a particular part of the space
 - On average, do better than single classifier!
- **But how do you**
 - force classifiers h_t to learn about different parts of the input space?
 - weight the votes of different classifiers? α_t

Boosting [Schapire 89]

- **Idea:** given a weak learning algorithm, run it multiple times on (reweighted) training data, then let learned classifiers vote
- On each iteration t :
 - weight each training example by how incorrectly it was classified
 - Learn a weak hypothesis – h_t
 - Obtain a strength for this hypothesis – α_t
- Final classifier:
$$H(X) = \text{sign}(\sum \alpha_t h_t(X))$$

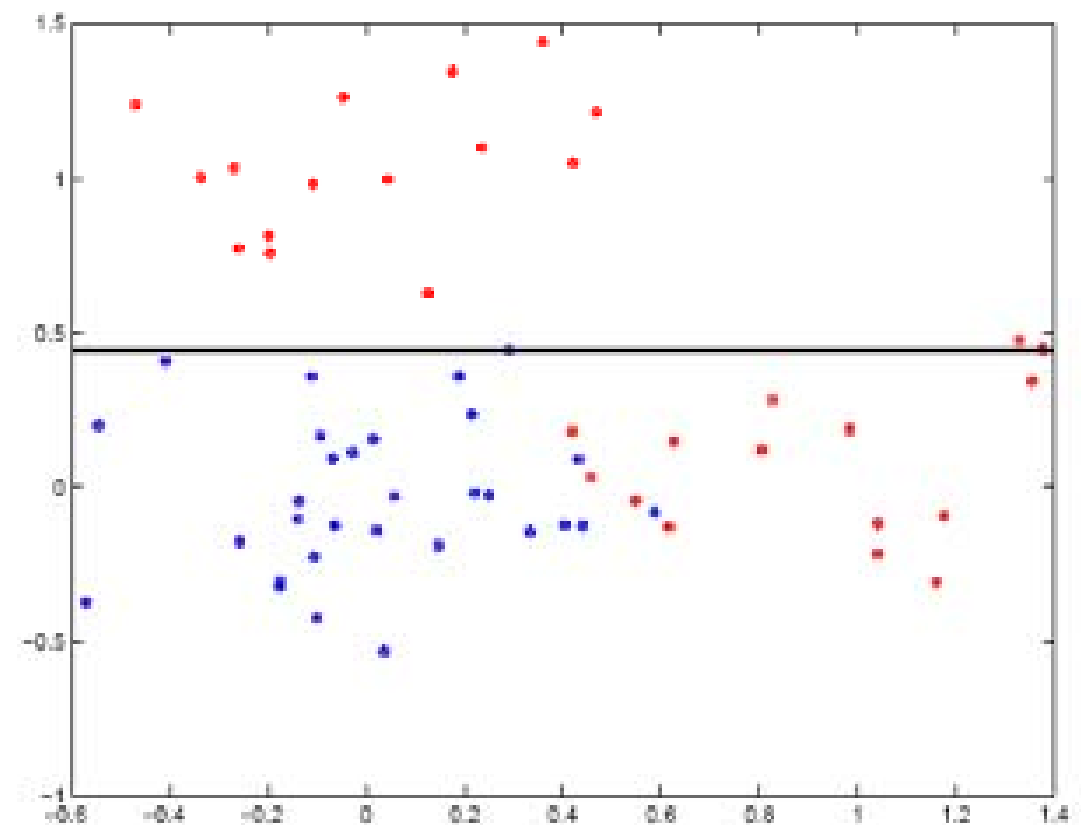
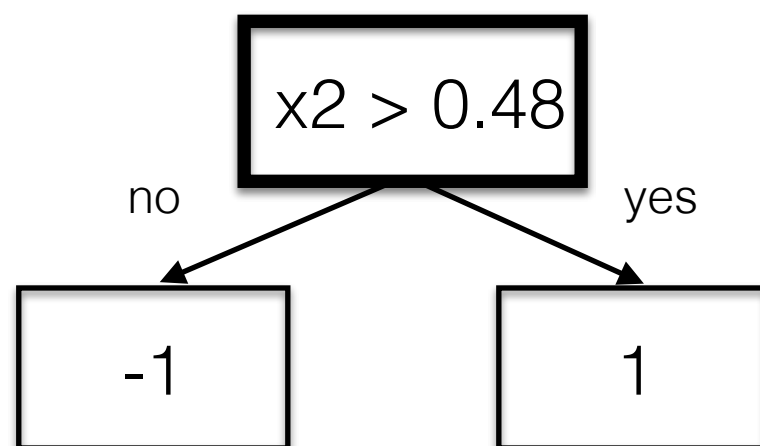
Combination of classifiers

- Suppose we have a family of component classifiers (generating ± 1 labels) such as decision stumps:

$$h(x; \theta) = \text{sign}(wx_k + b)$$

where $\theta = \{k, w, b\}$

- Each decision stump pays attention to only a single component of the input vector



$$w = 1, k = 2, b = 0.48$$

Combination of classifiers

- We'd like to combine the simple classifiers additively so that the final classifier is the sign of

$$\hat{h}(\mathbf{x}) = \alpha_1 h(\mathbf{x}; \theta_1) + \dots + \alpha_m h(\mathbf{x}; \theta_m)$$

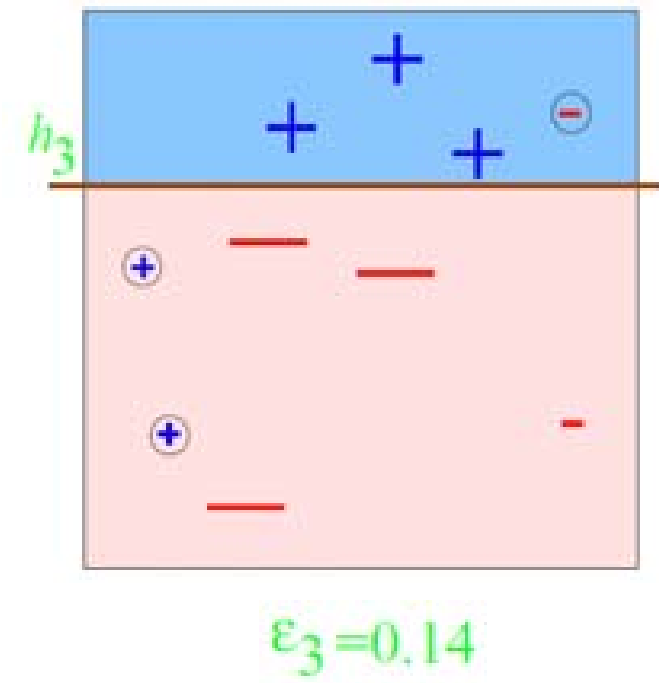
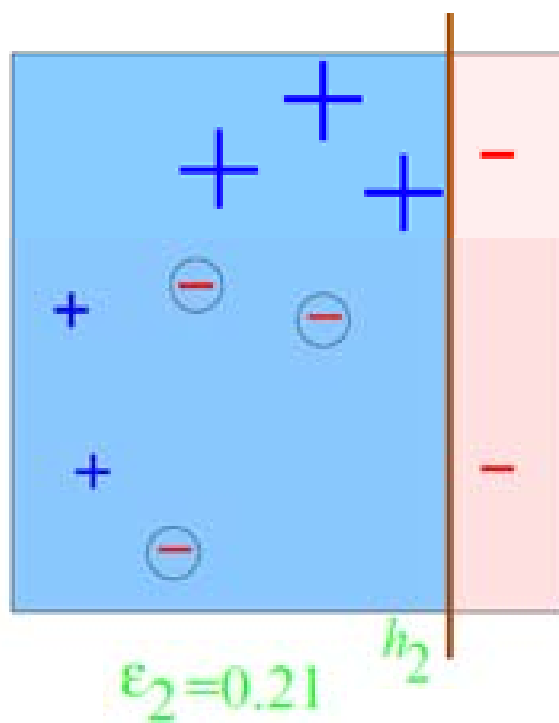
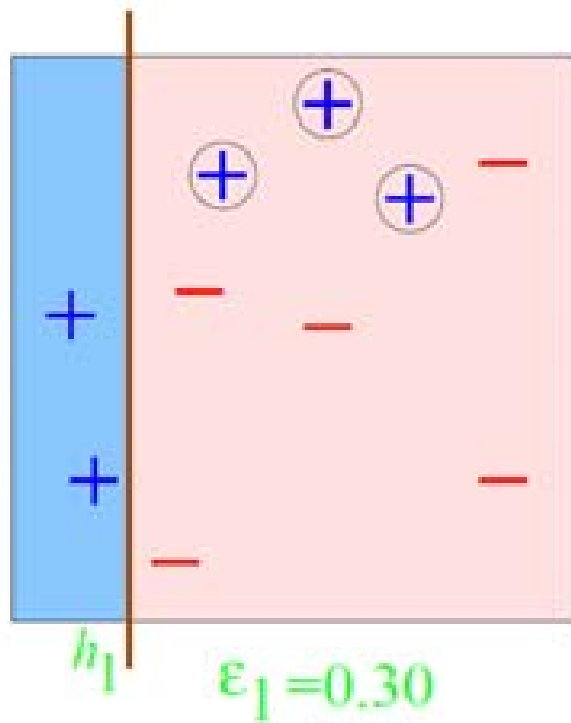
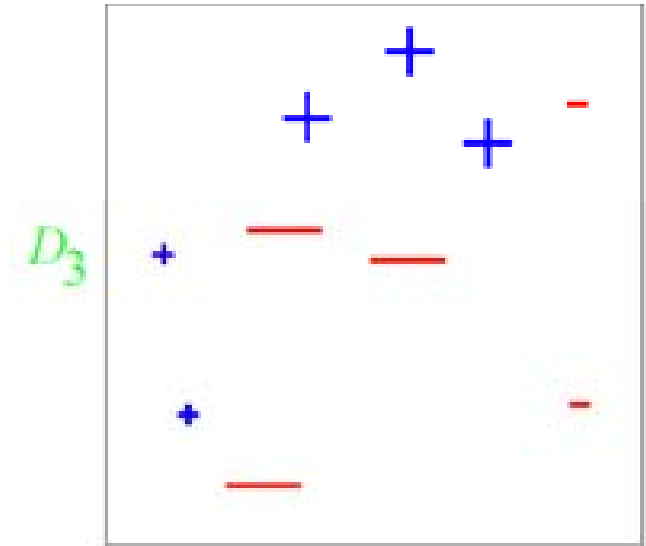
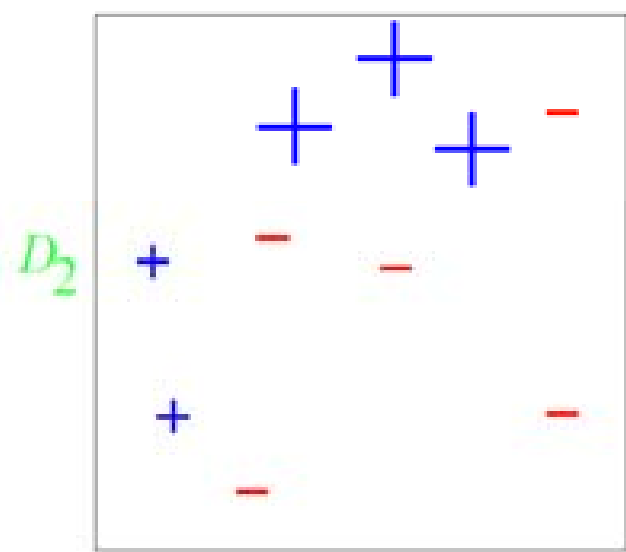
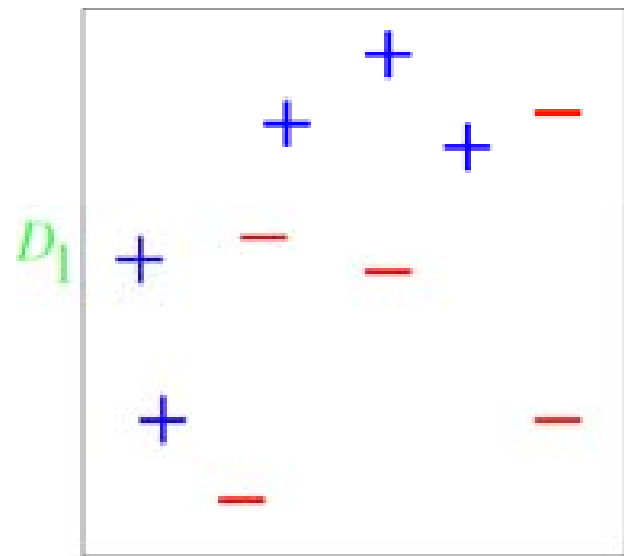
where the “votes” $\{\alpha_i\}$ emphasize component classifiers that make more reliable predictions than others

Recall

- On each iteration t :
 - weight each training example by how incorrectly it was classified
 - Learn a weak hypothesis – h_t
 - Obtain a strength for this hypothesis – α_t

Boosting example with decision stumps

$d = 2$
 $n = 10$

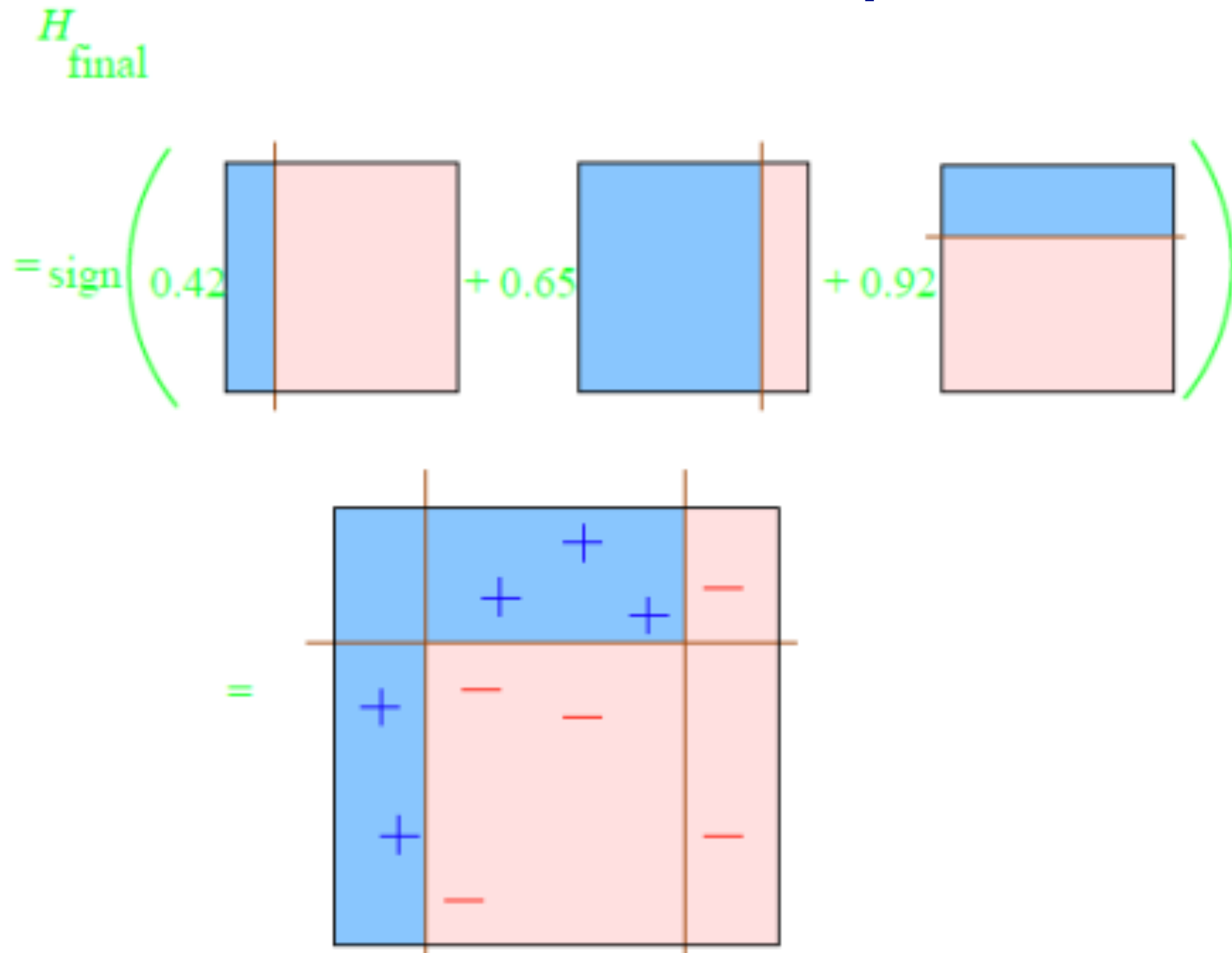


$$\alpha_1 = 0.42$$

$$\alpha_2 = 0.65$$

$$\alpha_3 = 0.92$$

Boosting example with decision stumps



AdaBoost

- **Input:**
 - N examples $\mathcal{S}_N = \{(x_1, y_1), \dots, (x_N, y_N)\}$
 - a weak base learner $h = h(x, \theta)$
- **Initialize:** equal example weights $w_i = 1/N$ for all $i = 1..N$
- **Iterate for $t = 1..T$:**
 1. train base learner according to **weighted example** set (w_t, x) and obtain hypothesis $h_t = h(x, \theta_t)$
 2. compute hypothesis error ε_t
 3. compute **hypothesis weight** α_t
 4. update **example weights** for next iteration w_{t+1}
- **Output:** final hypothesis as a linear combination of h_t

At k -th iteration, we have

$$f(x) = \sum_{j=1}^k \alpha_j h(x ; \theta_j)$$

Adaboost

- At the k th iteration we find (**any**) classifier $h(\mathbf{x}; \theta_k)$ for which the weighted classification error:

$$\varepsilon_k = \frac{\sum_{i=1}^n W_i^{k-1} I(y_i \neq h(\mathbf{x}_i; \theta_k))}{\sum_{i=1}^n W_i^{k-1}}$$

is better than chance.

- This is meant to be "easy" --- weak classifier
- Determine how many "votes" to assign to the new component classifier:

$$\alpha_k = 0.5 \log\left(\frac{1 - \varepsilon_k}{\varepsilon_k}\right)$$

epsilon small,
(1-epsilon)/epsilon is big
epsilon = 0.5 (random),
alpha = 0

 - stronger classifier gets more votes
- Update the weights on the training examples:

$$W_i^k = W_i^{k-1} \exp\{-y_i \alpha_k h(\mathbf{x}_i; \theta_k)\}$$

or equivalently

$$W_i^k = \exp(-y_i f(x_i)) \quad f(x) = \sum_{j=1}^k \alpha_j h(x; \theta_j)$$

Base learners


- Weak learners used in practice:
 - Decision stumps
 - Decision trees (e.g. C4.5 by Quinlan 1996)
 - Multi-layer neural networks
 - Radial basis function networks
- Can base learners operate on weighted examples?
 - In many cases they can be modified to accept weights along with the examples
 - In general, we can sample the examples (with replacement) according to the distribution defined by the weights

Exercise and break

- How can we modify logistic regression to use different weights for each example?
- Can we modify naive Bayes to use different weights for each example?

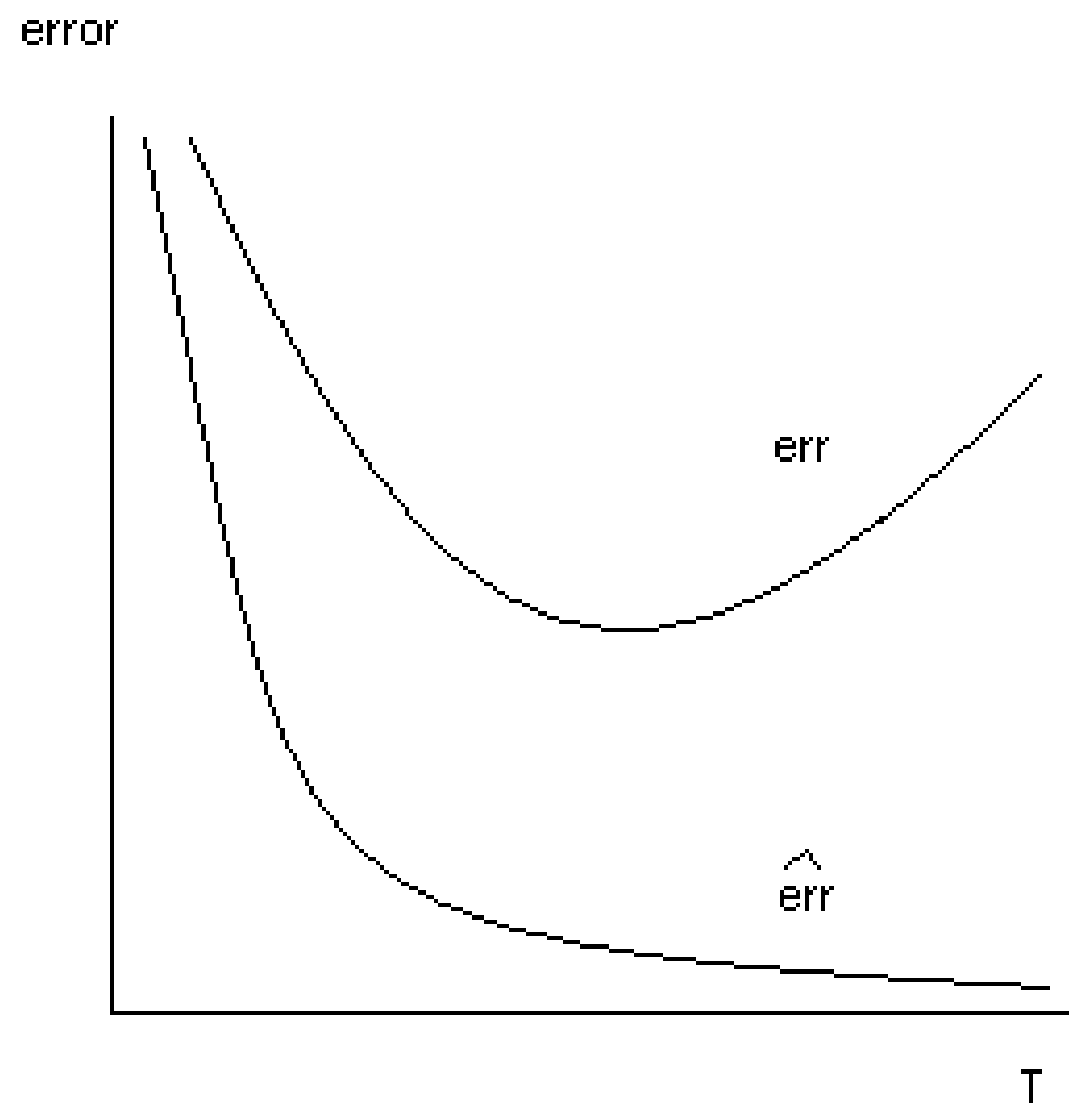
Generalization error bounds for Adaboost

$$error_{true}(H) \leq error_{train}(H) + \tilde{O}\left(\sqrt{\frac{Td}{m}}\right)$$

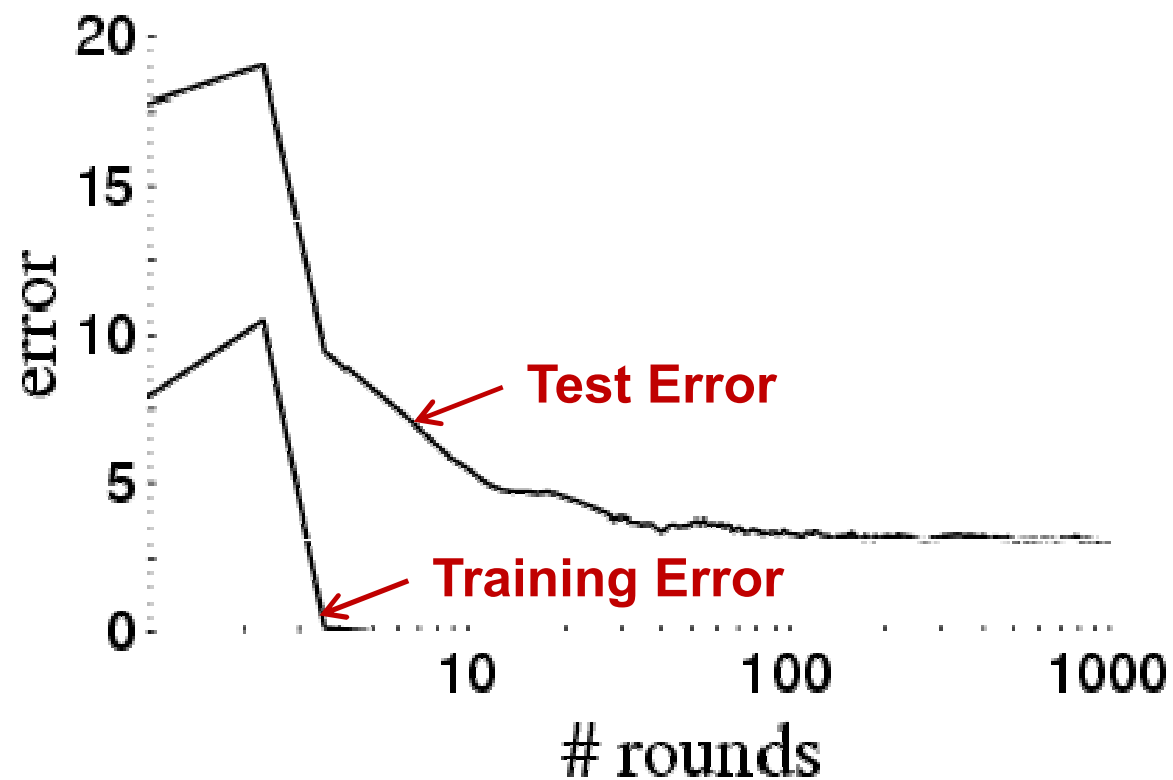
	bias	variance	
 tradeoff	large	small	T small
	small	large	T large

- T – number of boosting rounds
- d – VC dimension of weak learner, measures complexity of classifier
- m – number of training examples

Expected Adaboost behavior due to overfitting



Adaboost in practice



- Boosting often, **but not always**
 - Robust to overfitting
 - Test set error decreases even after training error is zero

Why does this seem to contradict the generalization bound?

Intuition

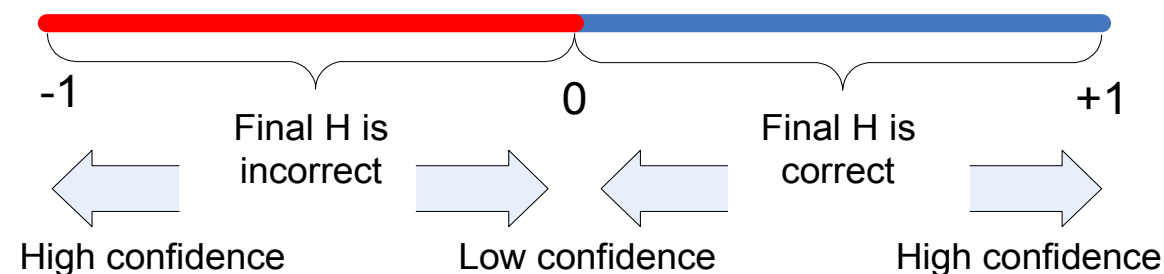
- Even when training error becomes zero, the confidence in the hypotheses continues to increase
- Large margin in training (increase in confidence) reduces the generalization error (rather than causing overfitting)
- Quantify with margin bound, to measure confidence of a hypothesis: when a vote is taken, the **more predictors agreeing**, the **more confident** you are in your prediction

Margin

$$\begin{aligned}\text{margin}(x, y) &= y f(x) \\ &= y \sum_t a_t h_t(x) \\ &= \sum_t a_t y h_t(x) \\ &= \sum_{t:h_t(x)=y} a_t - \sum_{t:h_t(x)\neq y} a_t\end{aligned}$$

h returns a -1 or a $+1$

where y is the correct label of instance x , and a_t is a normalized version of α_t such that $\alpha_t \geq 0$ and $\sum_t a_t = 1$. The expression $\sum_{t:h_t(x)=y} a_t$ stands for the weighted fraction of correct votes, and $\sum_{t:h_t(x)\neq y} a_t$ stands for the weighted fraction of incorrect votes. Margin is a number between -1 and 1 as shown in Figure 4.



see “Boosting the margin: A new explanation for the effectiveness of voting methods”, Schapire et al. 1997

General Boosting

- “Boosting algorithms as gradient descent”, Mason et al, 2000
- Adaboost is only one of many choices, with exponential loss
- Other examples and comparison: see “Cost-sensitive boosting algorithms: Do we really need them?” Nikolaou et al., 2016
- Main idea: given some loss L , (implicit) set of hypotheses and a weak learning algorithm,
 - generate hypothesis h_t that point in a descent direction
 - assign weight relative to how much pointing in descent direction

Boosting and logistic regression

Logistic regression equivalent to minimizing log loss

$$\sum_{i=1}^m \ln(1 + \exp(-y_i f(x_i)))$$

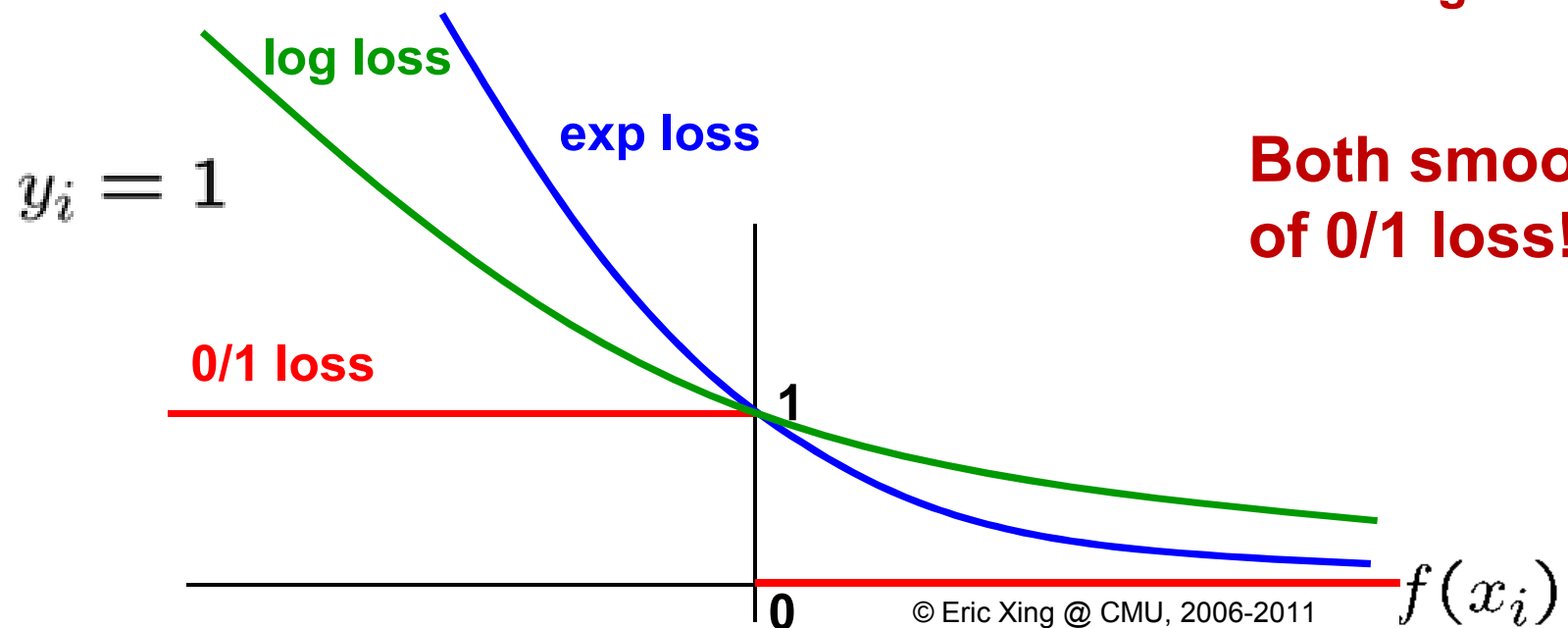
$$f(x) = w_0 + \sum_j w_j x_j$$

Boosting minimizes similar loss function!!

$$\frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i))$$

$$f(x) = \sum_t \alpha_t h_t(x)$$

Weighted average of weak learners



Any Boost

Algorithm 1 : AnyBoost

Require :

- An inner product space $(\mathcal{X}, \langle \cdot, \cdot \rangle)$ containing functions mapping from X to some set Y .
- A class of base classifiers $\mathcal{F} \subseteq \mathcal{X}$.
- A differentiable cost functional $C: \text{lin}(\mathcal{F}) \rightarrow \mathbb{R}$.
- A weak learner $\mathcal{L}(F)$ that accepts $F \in \text{lin}(\mathcal{F})$ and returns $f \in \mathcal{F}$ with a large value of $-\langle \nabla C(F), f \rangle$.

Let $F_0(x) := 0$.

for $t := 0$ to T **do**

 Let $f_{t+1} := \mathcal{L}(F_t)$.

if $-\langle \nabla C(F_t), f_{t+1} \rangle \leq 0$ **then**

 return F_t .

end if

 Choose w_{t+1} .

 Let $F_{t+1} := F_t + w_{t+1} f_{t+1}$

end for

return F_{T+1} .

Can be thought of as a stepsize



How does AdaBoost fit into this?

(see “Boosting algorithms as gradient descent”, Mason et al, 2000)

AdaBoost as AnyBoost

- Loss function is the exponential loss

$$C(F) := \frac{1}{m} \sum_{i=1}^m c(y_i F(x_i))$$

$$-\langle \nabla C(F), f \rangle = -\frac{1}{m^2} \sum_{i=1}^m y_i f(x_i) c'(y_i F(x_i)).$$

- Such an f corresponds to minimizing a weighted error with weights

$$\frac{c'(y_i F(x_i))}{\sum_{i=1}^m c'(y_i F(x_i))}$$

Diversity of the ensemble

- An important property appears to be diversity of the ensemble
- We get to define the hypothesis space: does not have to be homogenous (e.g., the set of linear classifiers)
- Strategies to promote this include:
 - using different types of learners (e.g., naive Bayes, logistic regression and decision trees)
 - pruning learners that are similar
 - random learners, which are more likely to be different than strong/deliberate algorithms which might learn similar predictions

Exercise: Can boosting be used for regression?

Algorithm 1 : AnyBoost

Require :

- An inner product space $(\mathcal{X}, \langle \cdot, \cdot \rangle)$ containing functions mapping from X to some set Y .
- A class of base classifiers $\mathcal{F} \subseteq \mathcal{X}$.
- A differentiable cost functional $C: \text{lin}(\mathcal{F}) \rightarrow \mathbb{R}$.
- A weak learner $\mathcal{L}(F)$ that accepts $F \in \text{lin}(\mathcal{F})$ and returns $f \in \mathcal{F}$ with a large value of $-\langle \nabla C(F), f \rangle$.

Let $F_0(x) := 0$.

for $t := 0$ to T **do**

 Let $f_{t+1} := \mathcal{L}(F_t)$.

if $-\langle \nabla C(F_t), f_{t+1} \rangle \leq 0$ **then**

 return F_t .

end if

 Choose w_{t+1} .

 Let $F_{t+1} := F_t + w_{t+1} f_{t+1}$

end for

return F_{T+1} .

If so, whats the loss?

How might this pseudocode change?