

CS6375: Machine Learning

Gautam Kunapuli

Ensemble Methods: Boosting



THE UNIVERSITY OF TEXAS AT DALLAS

Erik Jonsson School of Engineering and Computer Science

Ensemble Methods

Idea: Train models on different data set samples to reduce the model variance and/or bias

Problem: Only one training set; where do multiple models come from?

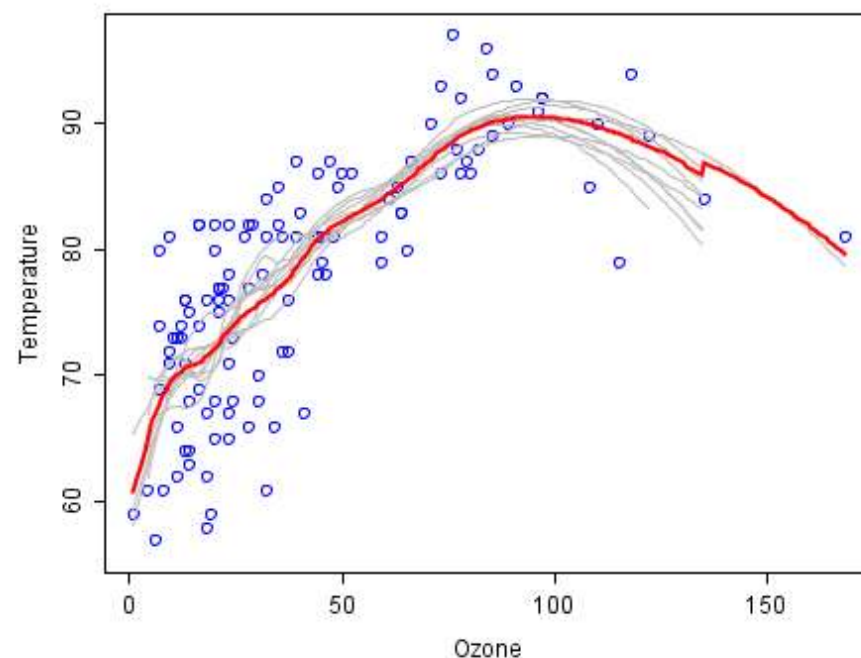
Solution: Take a **single** learning algorithm and generate multiple variations called **ensembles**

Bagging

- Uses **bootstrap sampling** to learn many (strong) models, whose predictions can be **aggregated**
- bagging reduces **variance** (and in practice, maybe increases bias slightly)
- bagging learns with **bootstrap samples** of the **same size as the original data set**
 - with decision trees, typically learns full trees
 - computational complexity is higher
- each model is learned independently of other models
 - insight from one model does not influence the learning of the next model

Boosting

- uses weak learners that have high bias
 - e.g., decision stumps (decision trees with depth 1)
- boosting reduces both **bias** and **variance**
- **iterative algorithm** that increases weights on hard examples
 - insight from previous iterations guides learning



AdaBoost: Adaptive Boosting

Basic idea behind Boosting: **examples** are given **weights**: at each iteration, a new hypothesis is learned and examples are **reweighted** to enable focus on examples that **most recently learned classifier** got wrong

Basic Algorithm for Boosting:

Initialize: set all examples to have equal weights

for each $t = 1, \dots, T$,

Learn a hypothesis h_t from weighted examples

Decrease weights of examples h_t classifies **correctly**

Calculate α_t , the weight of the current weak learner, h_t

return $h(x) = \sum_{t=1}^T \alpha_t h_t(x)$

Weighted examples: Base (weak) learner must focus on correctly classifying the most highly weighted examples while strongly avoiding over-fitting.

Weighted Hypotheses: During testing, each of the T hypotheses get a weighted vote proportional to their accuracy on the training data.

- Originally developed by computational learning theorists to guarantee performance improvements on fitting training data for a **weak learner**
 - a **weak learner** only needs to generate a hypothesis with a training accuracy **greater than 0.5** (Schapire, 1990);
 - **often, weak learners are only slightly better than random**
- practical algorithm, **AdaBoost**, for building ensembles that **empirically improves generalization** (Freund & Schapire, 1996).

AdaBoost: Weak Learners

Basic Algorithm for Boosting:

Initialize: set all examples to have equal weights

for each $t = 1, \dots, T$,

Learn a hypothesis h_t from weighted examples

Decrease weights of examples h_t classifies **correctly**

Calculate α_t , the weight of the current weak learner, h_t

return $h(x) = \sum_{t=1}^T \alpha_t h_t(x)$

a **weak learner** is typically easy to train and is simple, that is, of low complexity

- **high bias**, low variance
- **boosting** takes a weak learner and converts it to a strong learner
- just has to achieve an **accuracy slightly better than random guessing**, that is, error $\epsilon < 0.5$
- a weak learner achieves **accuracy-to-error ratio**:

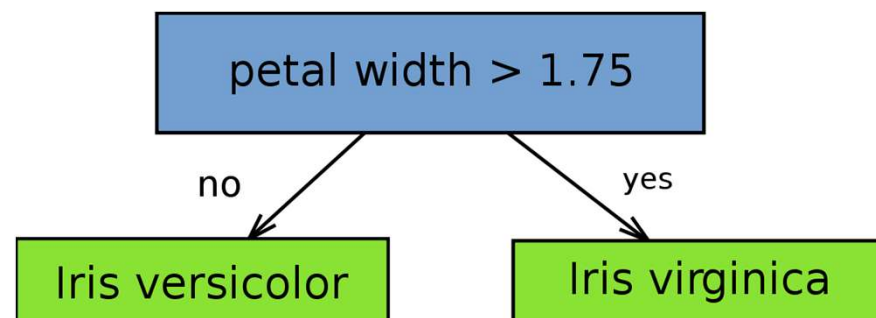
$$\frac{1-\epsilon}{\epsilon} > 1$$

- we can make the weight of a weak learner during boosting depend on its accuracy

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

- stronger learners have higher weights

Decision stumps are classical and often-used weak learners; Naïve Bayes, Logistic Regression also return probability of classification



Weak learners commonly used in practice:

- Decision stumps (axis parallel splits)
- Shallow decision trees
- Multi-layer neural networks
- Radial basis function networks

Note: There is nothing inherently weak about weak learners – we just think of them this way. In fact, **any learning algorithm** can be used as a weak learner.

AdaBoost: Training Set Distributions

Basic idea behind AdaBoost: maintain a **distribution** over **examples** that reflects their “hardness” of classification; a new hypothesis is learned and the distribution is **updated** to enable focus on examples that **most recently learned classifier** got wrong

Basic Algorithm for Boosting:

Initialize: set all examples to have equal weights

for each $t = 1, \dots, T$,

Learn a hypothesis h_t from weighted examples

Decrease weights of examples h_t classifies **correctly**

Calculate α_t , the weight of the current weak learner, h_t

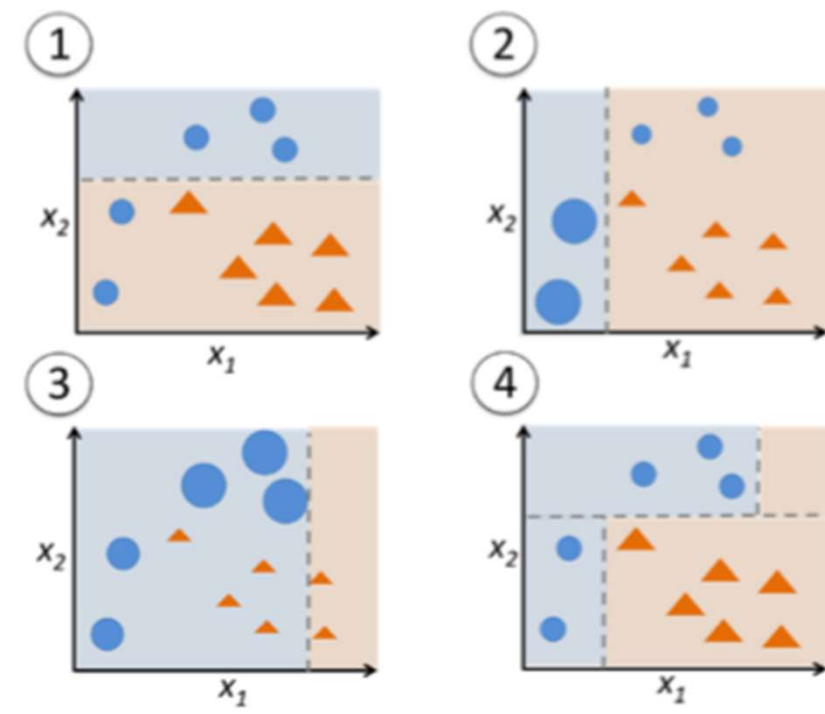
return $h(x) = \sum_{t=1}^T \alpha_t h_t(x)$

weights on examples can be converted to a **distribution** that reflects their “**hardness of classification**”

thus, each training example (x_i, y_i) has weights $D(i)$, with

$$\sum_{i=1}^n D(i) = 1$$

- most misclassified points get highest weights
- this ensures that the algorithm can focus on training examples with higher weights



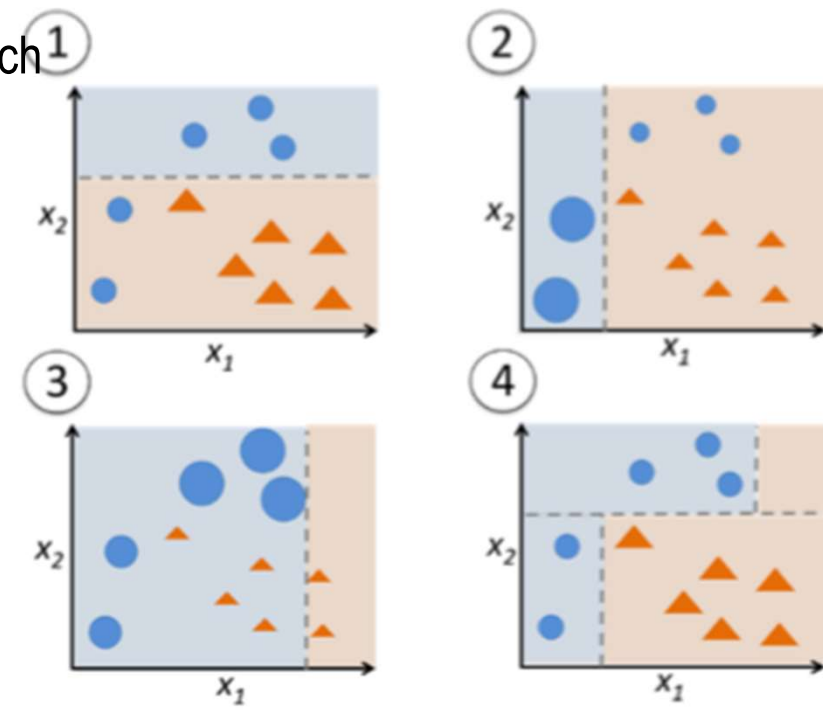
For boosting, **we need a weak learner** that can

- handle weighted **examples/distributions**
- alternately, **sample** training examples **according to the distribution** (more on next slide)
 - contrast this with bagging!

Learning with Weighted Training Examples

In a **weighted dataset** we have a weight associated with each training example:

- $D(i)$ is the weight of i -th training example (x_i, y_i)
- i -th training example counts as $D(i)$ training examples; if we “resampled” data, we would get more samples of “heavier” data points
- Now, in all calculations, the i -th **training example counts as $D(i)$ “examples”**



Example 1: in Maximum Likelihood Estimation

Unweighted data: $\#(y = c) = \sum_i \mathbf{1}(y = c)$

Weighted data: $\#(y = c) = \sum_i D(i) \mathbf{1}(y = c)$

Example 2: in Decision Stumps:

- first, when computing $H(Y) = -\sum_y P(Y = c) \log_2 P(Y = c)$ for a class c , use the weights; for instance, in the binary classification case:

$$P(y = 0) = \frac{\#(y=0)}{\#(y=0) + \#(y=1)} \text{ (unweighted)} \quad P(y = 0) = \frac{\sum w(y=0)}{\sum w(y=0) + \sum w(y=1)} \text{ (weighted)}$$

- second, when computing $H(Y|X) = -\sum_x P(X = x) \sum_y P(Y = c | X = x) \log_2 (Y = c | X = x)$
- alternately, since decision stumps are easy to compute, simply compute **all possible decision stumps and select the one with the smallest weighted error** as the best weak learners

AdaBoost: Full Algorithm

Given: $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in X, y_i \in Y = \{-1, +1\}$
Initialize $D_1(i) = 1/m$.
For $t = 1, \dots, T$:

best decision stump is the one that **minimizes** the **weighted training error**

- Train weak learner using distribution D_t .
- Get weak hypothesis $h_t : X \rightarrow \{-1, +1\}$ with error

$$\epsilon_t = \Pr_{i \sim D_t} [h_t(x_i) \neq y_i]. \quad \epsilon_t = \frac{1}{\sum_{i=1}^n D_t(i)} \sum_{i=1}^n D_t(i) \cdot \mathbb{I}[h_t(\mathbf{x}_i) \neq y_i]$$

- Choose $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$.
- Update:

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(x_i) = y_i \\ e^{\alpha_t} & \text{if } h_t(x_i) \neq y_i \end{cases} \\ &= \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t} \end{aligned}$$

re-normalize the weights into a distribution

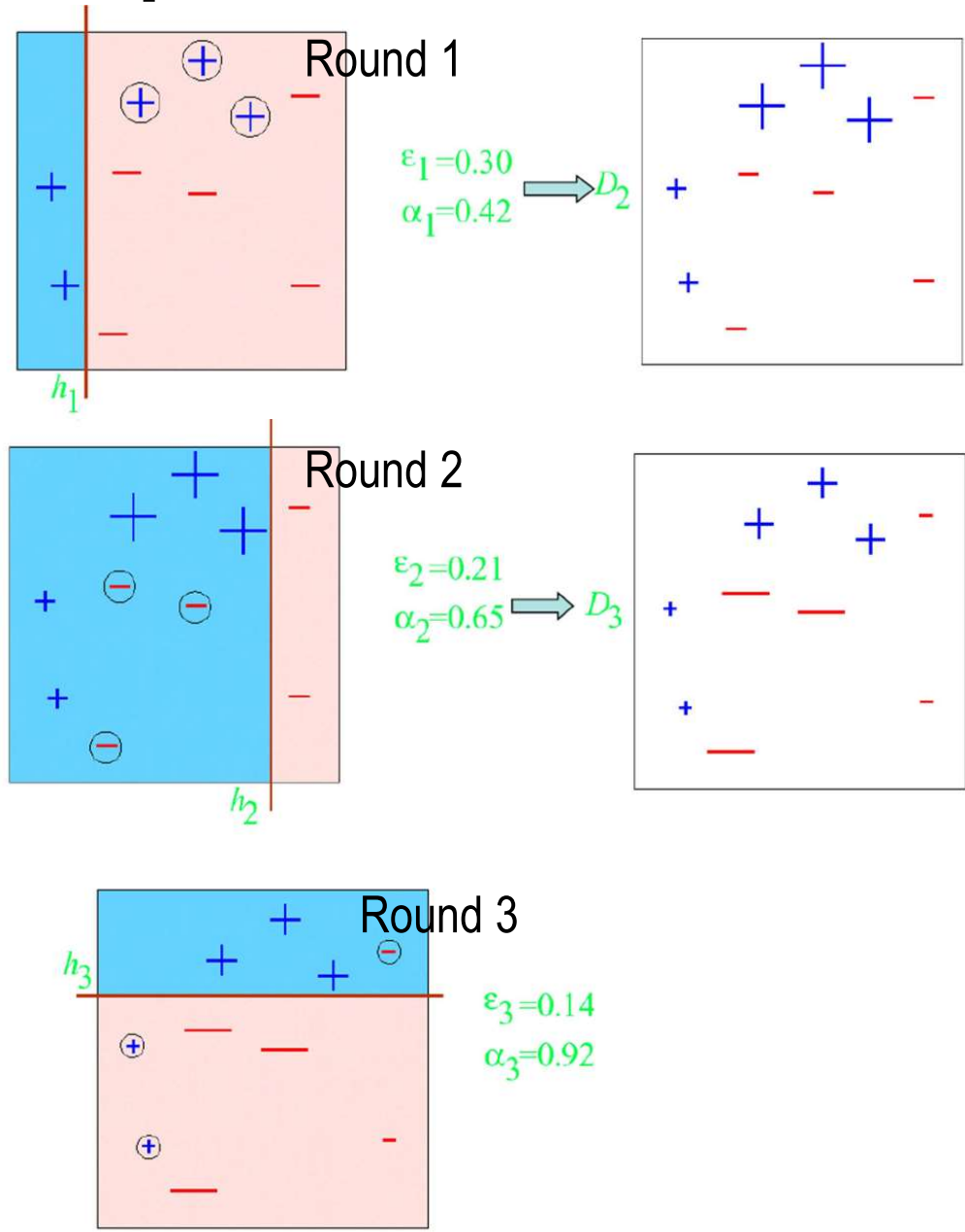
$$Z_t = \sum_{i=1}^n D_t(i) \cdot \exp(-\alpha_t y_i h_t(\mathbf{x}_i))$$

where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution).

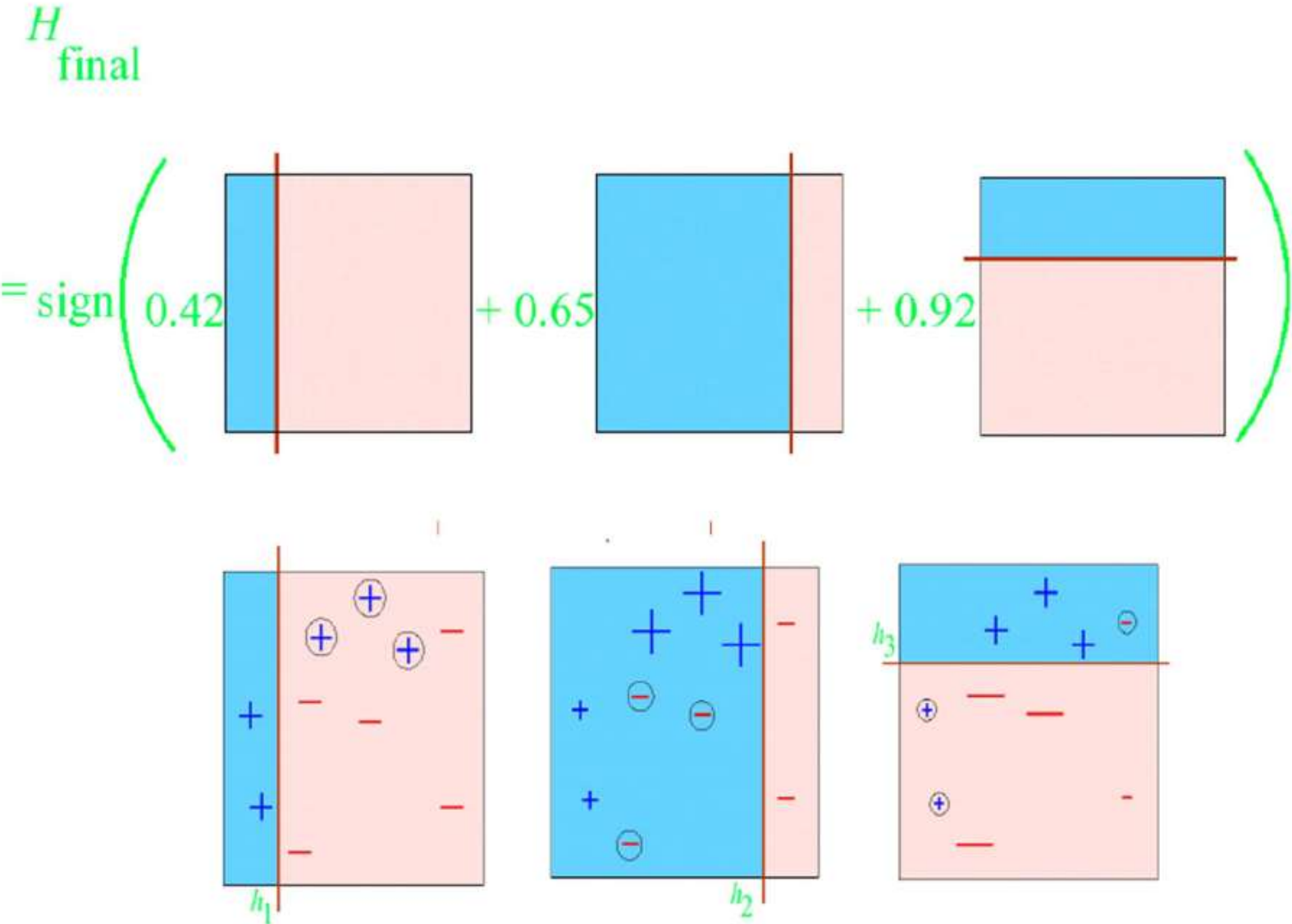
Output the final hypothesis:

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

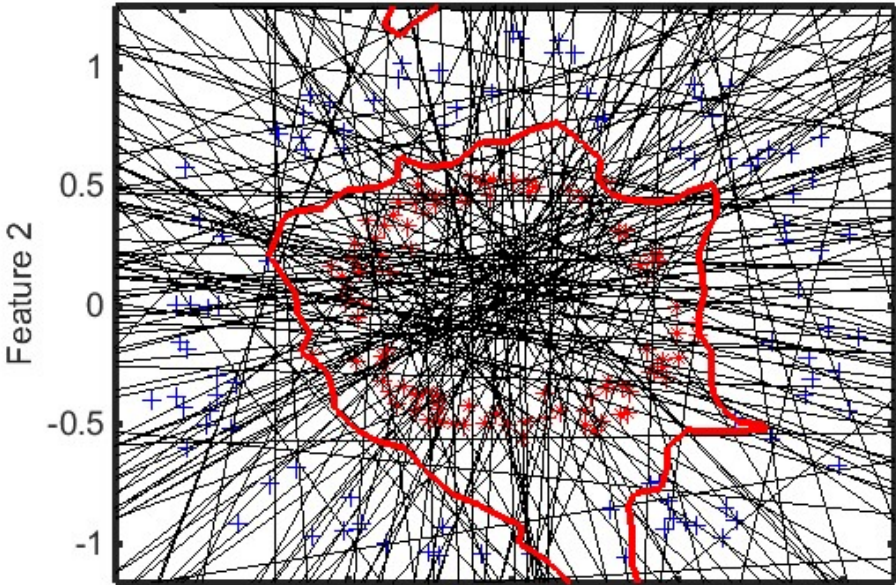
AdaBoost: Example



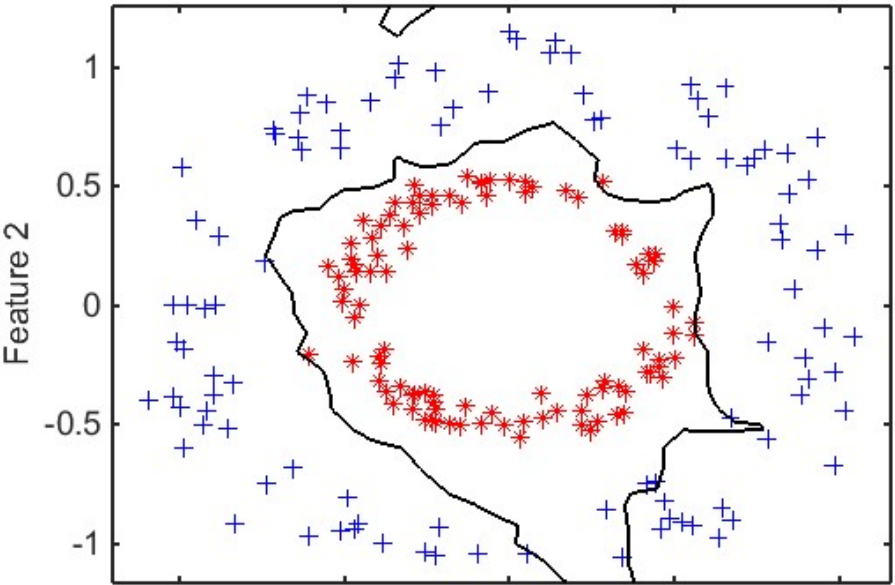
AdaBoost: Example



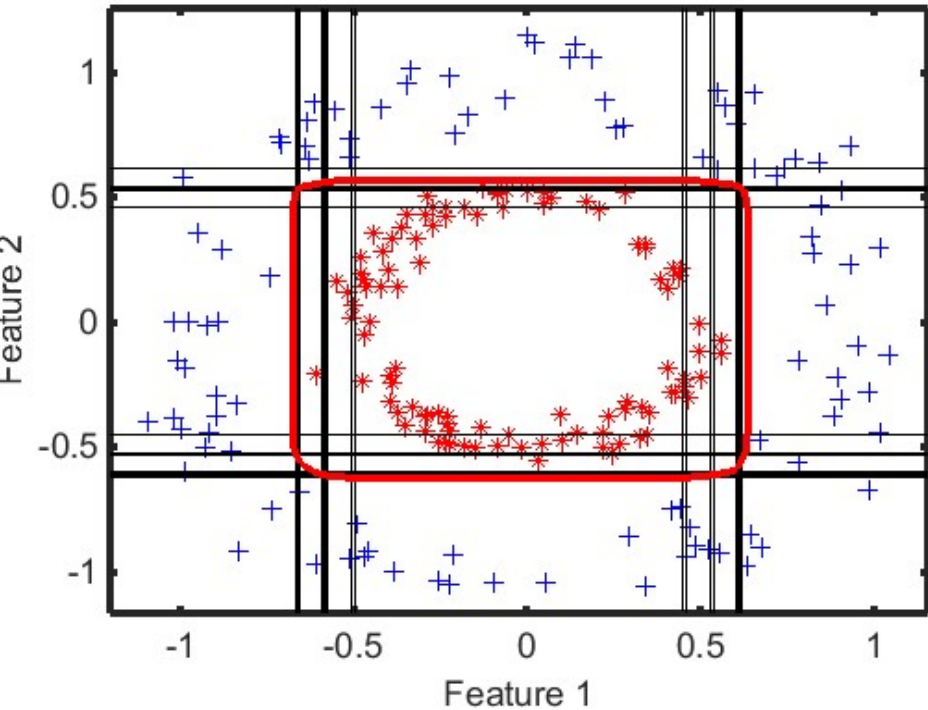
200 base classifiers: single epoch linear perceptron



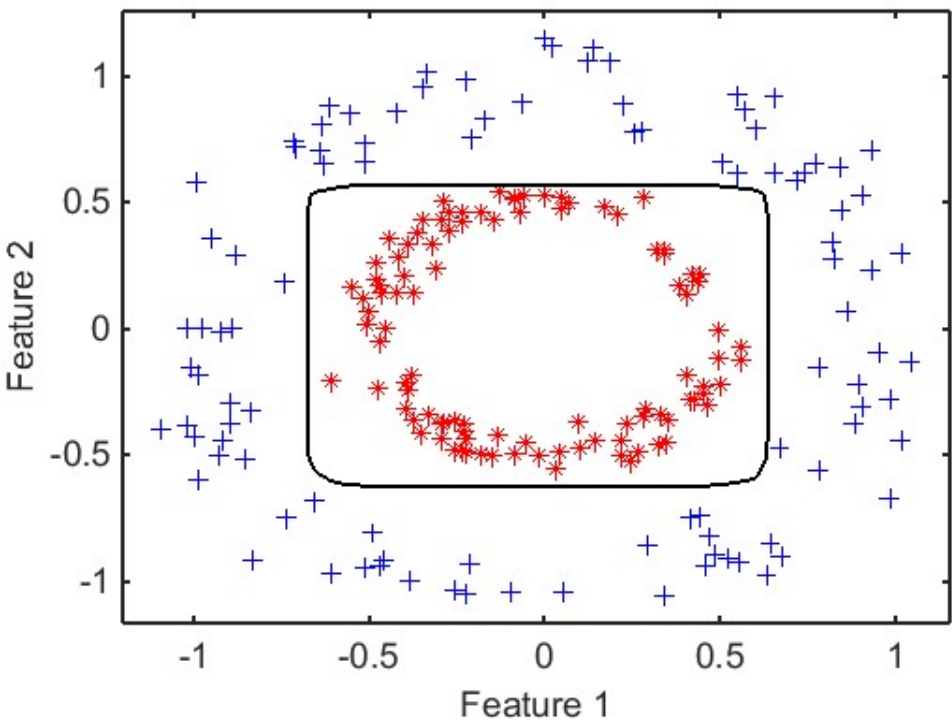
Result adaboost, 200 weak classifiers combined



200 base classifiers: single epoch decision stump



Result adaboost, 200 weak classifiers combined



Boosting Optimizes Exponential Loss

In 2000, Friedman et al. interpreted AdaBoost as stage-wise forward additive model that actually minimizes the exponential loss function, $L(y, f(x)) = E[e^{-y f(x)}]$

Split the exponential loss into positive and negative components
 $E[e^{-y f(x)}] = e^{f(x)} P(y = -1) + e^{-f(x)} P(y = 1)$

Take the gradient and set to zero
 $\frac{d}{df(x)} E[e^{-y f(x)}] = e^{f(x)} P(y = -1) - e^{-f(x)} P(y = 1)$

Boosting learns $f(x)$ as below (compare with α_t on Slide 7)
$$f(x) = \frac{1}{2} \log \frac{P(y = 1)}{P(y = -1)}$$

Logistic regression:

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

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

where x_j predefined features
(linear classifier)

- Jointly optimize over all weights $w_0, w_1, w_2 \dots$

Boosting:

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

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

where $h_t(x)$ defined dynamically to fit data
(not a linear classifier)

- Weights α_t learned per iteration incrementally



Boosting Increases The Margin

We can write the **combined classifier** in a more useful form by dividing the predictions by the “total number of votes”:

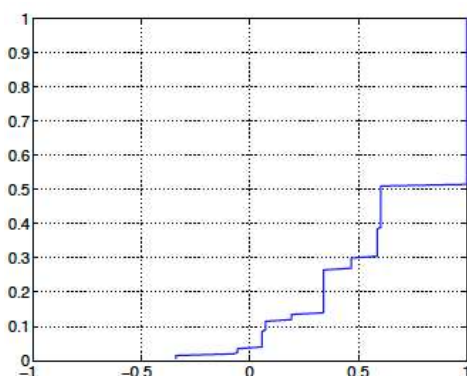
$$h_{t+1}(\mathbf{x}_i) = \frac{\alpha_1 h_1(\mathbf{x}_i) + \cdots + \alpha_t h_t(\mathbf{x}_i)}{\alpha_1 + \cdots + \alpha_t}$$

- This allows us to define a clear notion of “**voting margin**” that the combined classifier achieves for each training example:

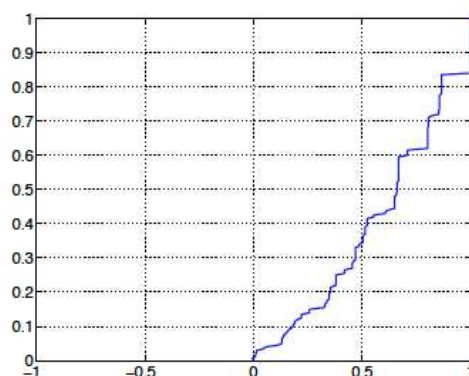
$$\text{margin}(\mathbf{x}_i) = y_i h_{t+1}(\mathbf{x}_i)$$

- The margin lies in $[-1, 1]$ and is negative for all misclassified examples.
- Successive boosting iterations **improve the majority vote** or **margin** for the training examples

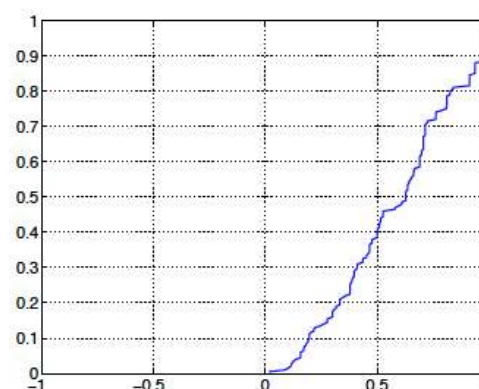
the **margin of a single data point** is defined to be the **distance** from the data point **to the decision boundary**



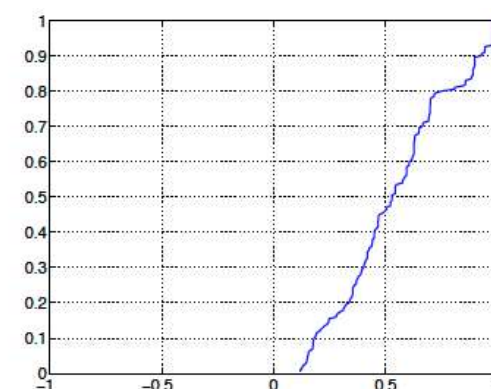
4 iterations



10 iterations



20 iterations



50 iterations

Cumulative distributions of margin values

Boosting: Pros and Cons

Pros

- **fast**
- **simple** and easy to program
- **no parameters** to tune (except **T**)
- **flexible** — can combine with **any** learning algorithm
- **no prior knowledge** needed about weak learner
- **provably effective**, provided can consistently find rough rules of thumb
 - shift in mind set — goal now is merely to find classifiers barely better than random guessing
- **versatile**
- can use with data that is textual, numeric, discrete, etc.
- has been extended to learning problems well beyond binary classification

Cons

- performance of AdaBoost depends on **data** and **weak learner**
- consistent with theory, AdaBoost can **fail** if
- weak classifiers too **complex** (! overfitting)
- weak classifiers too **weak** (! underfitting)
- empirically, AdaBoost seems especially susceptible to uniform noise

Good 😊 : Can identify outliers since focuses on examples that are hard to categorize

Bad ☹️ : Too many outliers can degrade classification performance dramatically increase time to convergence

