



University of Colorado **Boulder**

Department of Computer Science

CSCI 5622: Machine Learning

Chenhao Tan

Lecture 16: Ensemble Learning

Slides adapted from Chris Ketelsen, Jordan Boyd-Graber,  
and Noah Smith

# Administrivia

- Final project teamwork
- Peer feedback for a bonus point
- Be constructive

# Learning objectives

- Understand the general idea behind ensembling
- Learn about Adaboost
- Learn the math behind boosting (bonus)

# Outline

- Ensemble methods
- Bagging and random forest
- General idea of boosting
- Adaboost

# Outline

- Ensemble methods
- Bagging and random forest
- General idea of boosting
- Adaboost

# Outline for CSCI 5622

We've already covered stuff in **blue!**

- Problem formulations: classification, regression
- Supervised techniques: decision trees, nearest neighbors, perceptron, linear models, neural networks, support vector machine, kernel methods
- Unsupervised techniques: clustering, linear dimensionality reduction, topic modeling
- “Meta-techniques”: ensembles, expectation-maximization, variational inference
- Understanding ML: limits of learning, practical issues, bias & fairness
- Recurring themes: (stochastic) gradient descent

## Ensemble methods

We have learned

- Decision trees
- Perceptron
- KNN
- Naïve Bayes
- Logistic regression
- Neural networks
- Support vector machines

A meta question: why do we only use a single model?

## Ensemble methods

We have learned

- Decision trees
- Perceptron
- KNN
- Naïve Bayes
- Logistic regression
- Neural networks
- Support vector machines

A meta question: why do we only use a single model?

In practice, especially to win competitions, many models are used together.

## Ensemble methods

---

There are many techniques to use multiple models.

- Bagging
  - Train classifiers on subsets of data
  - Predict based on majority vote
- Boosting
  - Build a sequence of dumb models
  - Modify training data along the way to focus on difficult examples
  - Predict based on weighted majority vote of all the models
- Stacking
  - Take multiple classifiers' outputs as inputs and train another classifier to make final prediction

# Outline

- Ensemble methods
- Bagging and random forest
- General idea of boosting
- Adaboost

## Recap of decision tree

---

**Algorithm:** DTREETRAIN

**Data:** data  $D$ , feature set  $\Phi$

**Result:** decision tree

*if all examples in  $D$  have the same label  $y$ , or  $\Phi$  is empty and  $y$  is the best guess*

**then**

    | return LEAF( $y$ );

**else**

    | **for each feature  $\phi$  in  $\Phi$  do**

        | | partition  $D$  into  $D_0$  and  $D_1$  based on  $\phi$ -values;

        | | let mistakes( $\phi$ ) = (non-majority answers in  $D_0$ ) + (non-majority answers in  $D_1$ );

    | **end**

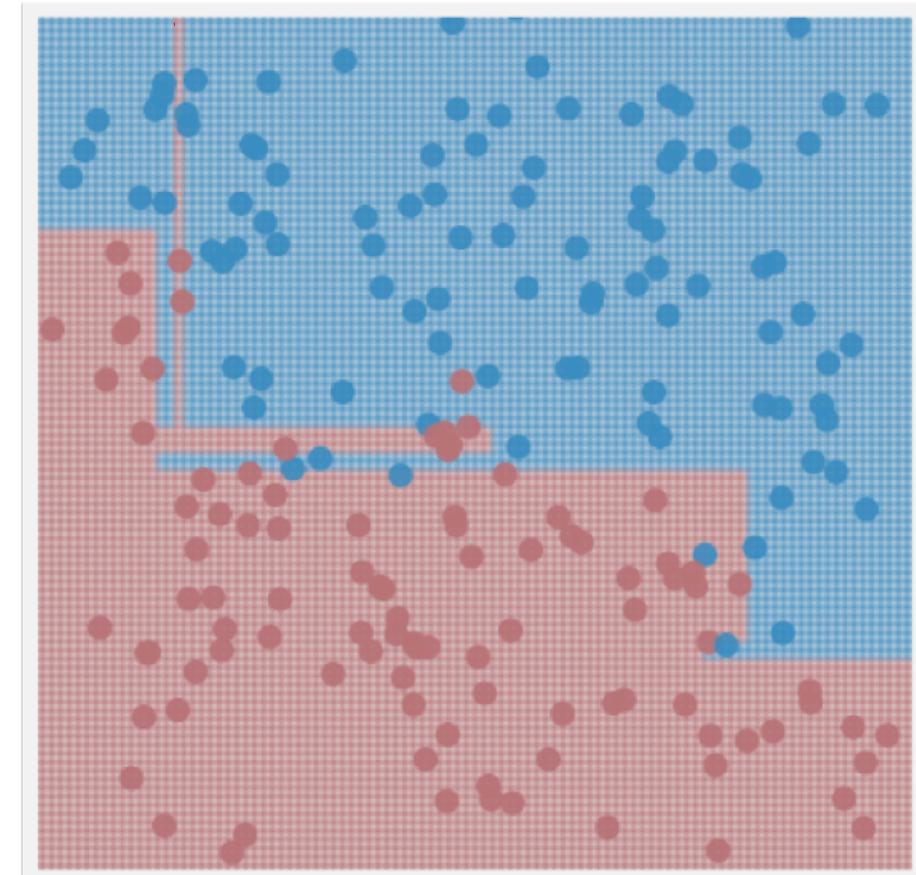
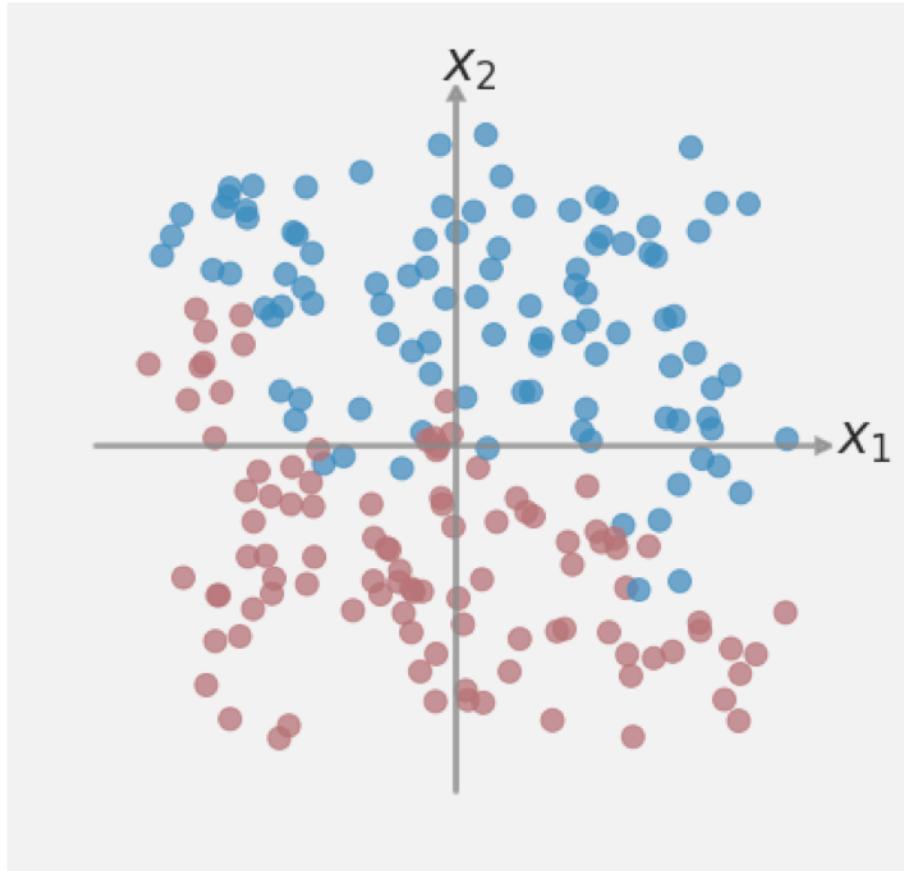
    | let  $\phi^*$  be the feature with the smallest number of mistakes;

    | return NODE( $\phi^*$ , {0 → DTREETRAIN( $D_0$ ,  $\Phi \setminus \{\phi^*\}$ )}, 1 → DTREETRAIN( $D_1$ ,  $\Phi \setminus \{\phi^*\}$ ));

**end**

## Recap of decision tree

Full decision tree classifiers tend to overfit

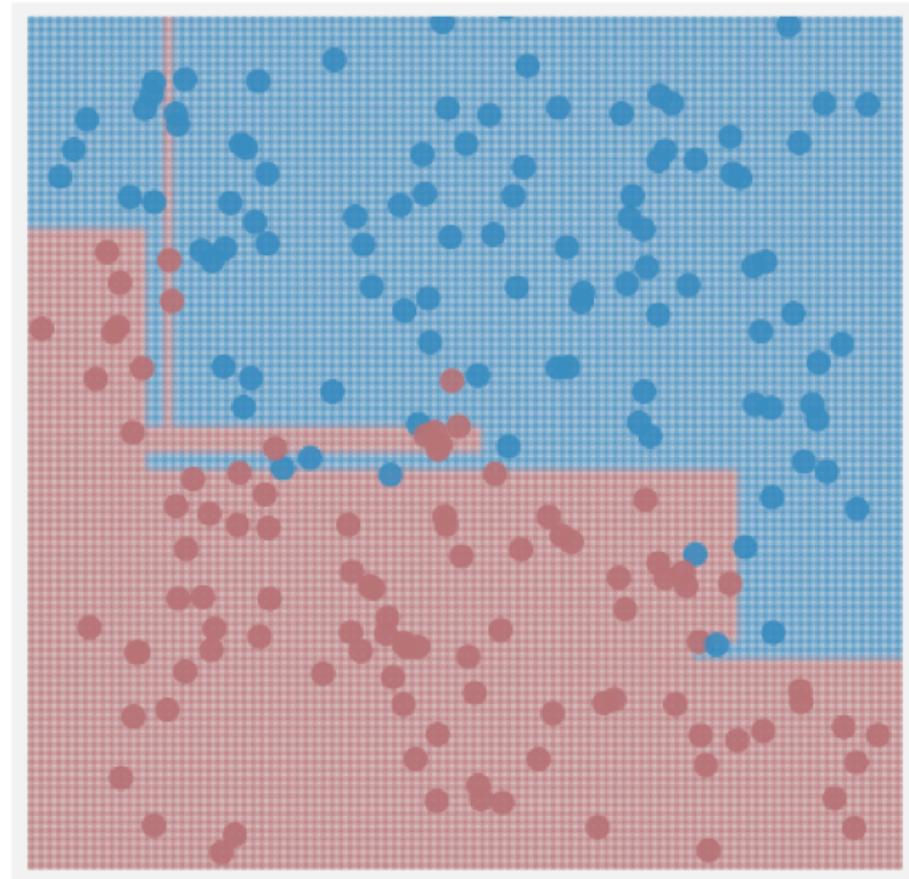


## Recap of decision tree

Full decision tree classifiers tend to overfit

We could alleviate this a bit with pruning.

- Prepruning
  - Don't let the tree have too many levels
  - Stopping conditions
- Postpruning
  - Build the complete tree
  - Go back and remove vertices that don't affect performance too much



## Bagging

We will see how we can use many decision trees.

General idea:

- Train numerous slightly different classifiers
- { Use each classifier to make a prediction on a test instance
- Predict by taking majority vote from individual classifiers

## Bagging

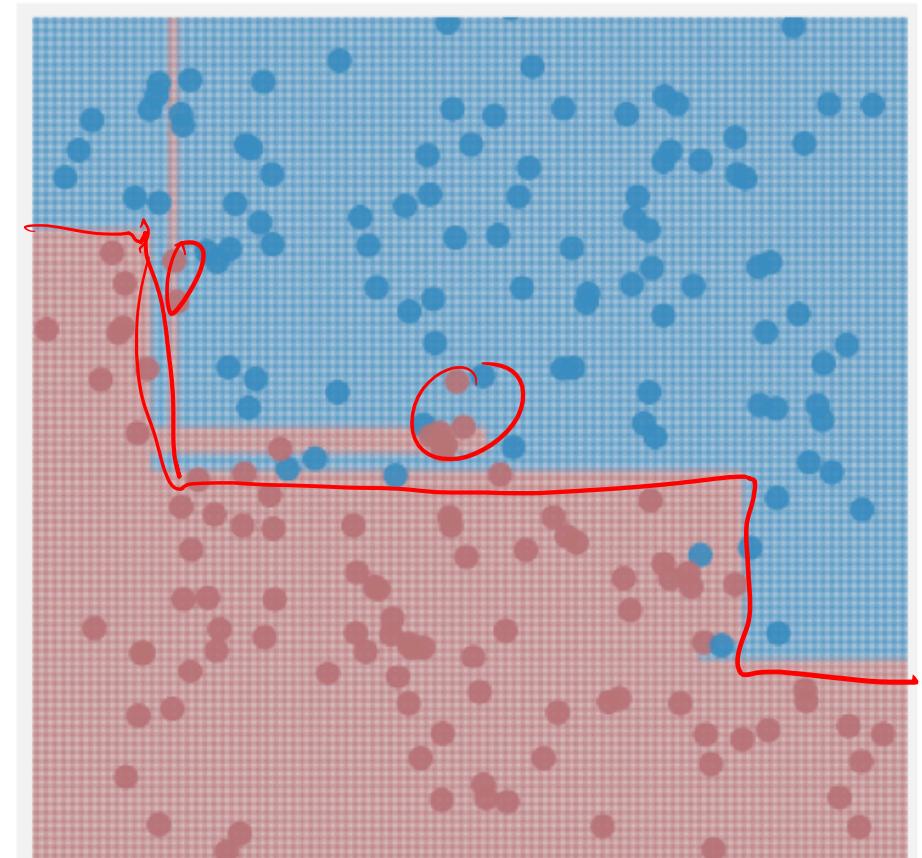
---

Today we will look at two different types of ensembled decision tree classifiers. The simplest is called *bootstrapped aggregation* (or bagging).

## Bagging

Today we will look at two different types of ensembled decision tree classifiers. The simplest is called *bootstrapped aggregation* (or bagging).

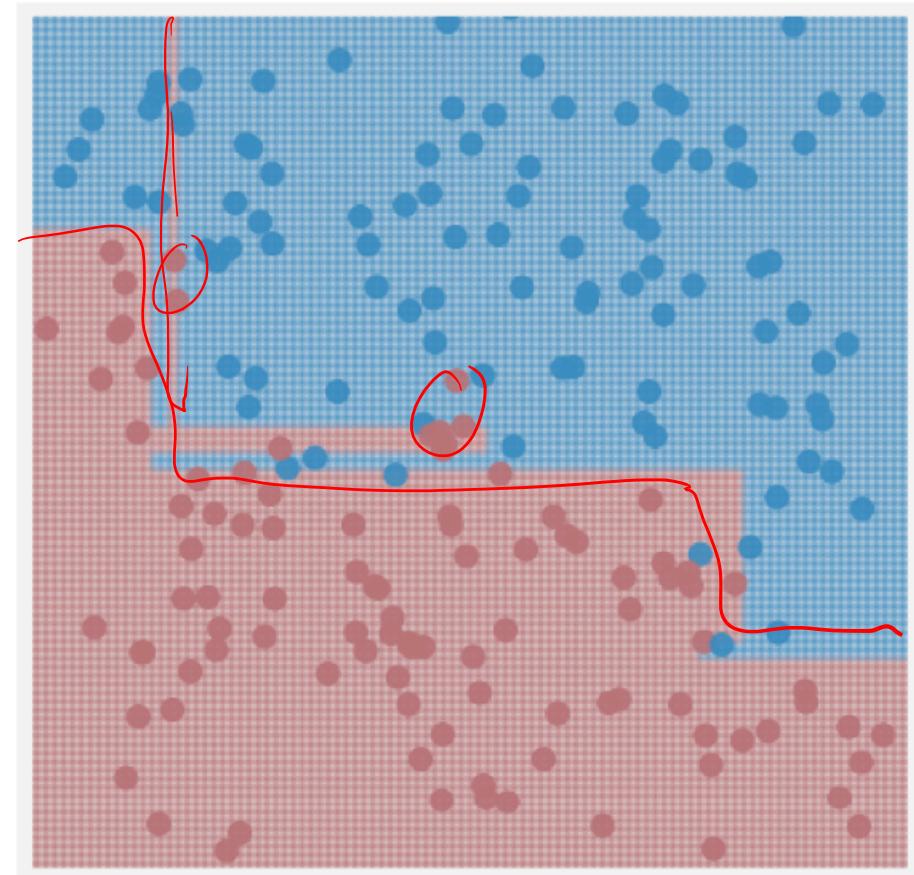
Idea: What would have happened if those training examples that we clearly overfit to weren't there?



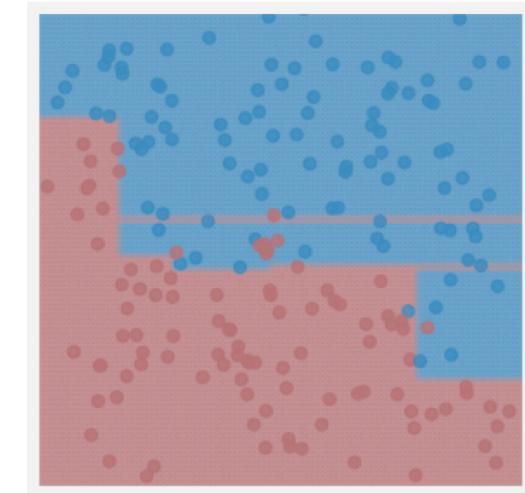
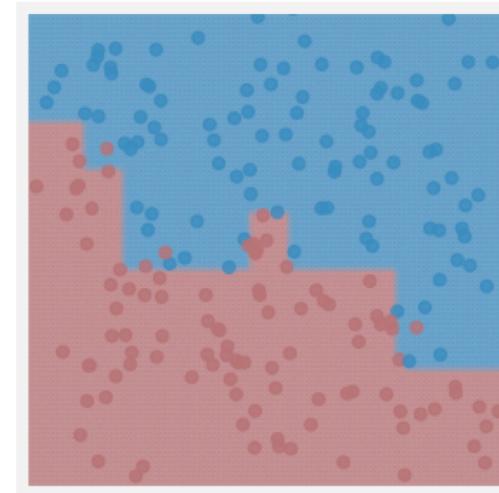
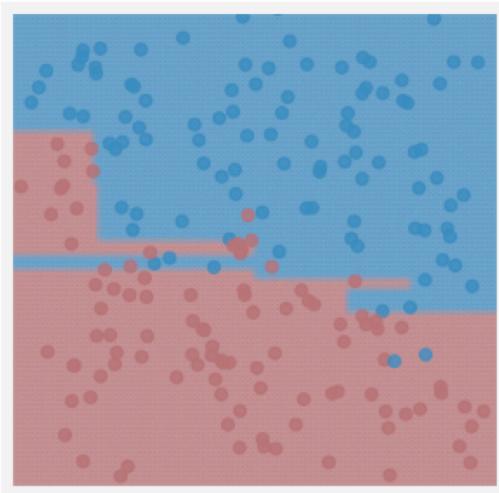
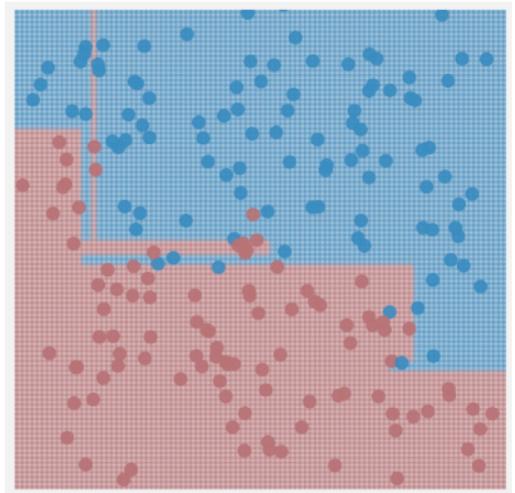
## Bagging

The simplest is called *bootstrapped aggregation* (or bagging).

- Sample  $n$  training examples with replacement
- Fit decision tree classifier to each sample
- Repeat many times
- Aggregate results by majority vote



## Bagging



## Bagging

$$i \in \{(x, y)\} \quad n = |\{(x, y)\}|$$

How many unique instances show up in  $n$  samples with replacement?

not showing up  $(1 - \frac{1}{n})^n$   $\lim_{n \rightarrow \infty} (1 - \frac{1}{n})^n = \frac{1}{e}$

showing up  $1 - \frac{1}{e}$

## Bagging

How many unique instances show up in  $n$  samples with replacement?

The likelihood that an instance is never chosen in  $n$  draws is

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = \frac{1}{e}$$

## Bagging

---

In general, assume that  $y \in \{-1, 1\}$  and each individual DCT has  $h_k$   
Then we can define the bagged classifier as

$$H(x) = \text{sign} \sum_{k=1}^K h_k(x)$$

## Bagging Pros and Cons

---

Pros:

- Results in a much lower variance classifier than a single decision tree

Cons:

- Results in a much less interpretable model

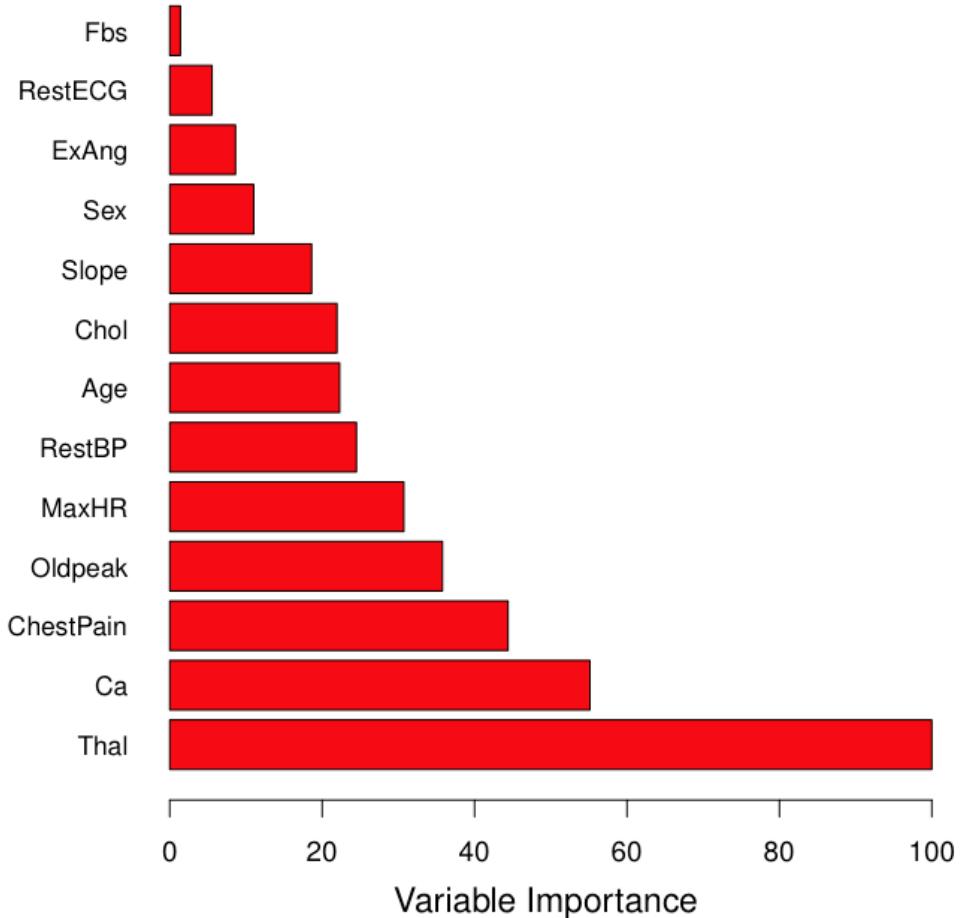
We essentially give up some interpretability in favor of better prediction accuracy.  
But we can still get insight into what our model is doing using bagging.

## Bagging and feature importance

---

Although we cannot follow the tree structure any more, we can estimate average feature importance of single features.

Feature importance in a single tree can be defined as information gain achieved by splitting on this feature.



## An even better approach: random forests

It turns out, we can take the bagging idea and make it even better.

Suppose you have a particular feature that is a very strong predictor for the response.

So strong that in almost all Bagged Decision Trees, it's the first feature that gets split on.

## An even better approach: random forests

It turns out, we can take the bagging idea and make it even better.

Suppose you have a particular feature that is a very strong predictor for the response.

So strong that in almost all Bagged Decision Trees, it's the first feature that gets split on.

When this happens, the trees can all look very similar. They're too correlated.

Maybe always splitting on the same feature isn't the best idea.

## An even better approach: random forests

---

It turns out, we can take the bagging idea and make it even better.

Suppose you have a particular feature that is a very strong predictor for the response.

So strong that in almost all Bagged Decision Trees, it's the first feature that gets split on.

When this happens, the trees can all look very similar. They're too correlated.

Maybe always splitting on the same feature isn't the best idea.

Maybe, like with bootstrapping training examples, we split on a random subset of features too.

## Random forests

---

- Still doing bagging, in the sense that we fit bootstrapped resamples of training data
- But every time we have to choose a feature to split on, don't consider all  $p$  features
- Instead, consider only  $\sqrt{p}$  features to split on

## Random forests

---

- Still doing bagging, in the sense that we fit bootstrapped resamples of training data
- But every time we have to choose a feature to split on, don't consider all  $p$  features
- Instead, consider only  $\sqrt{p}$  features to split on

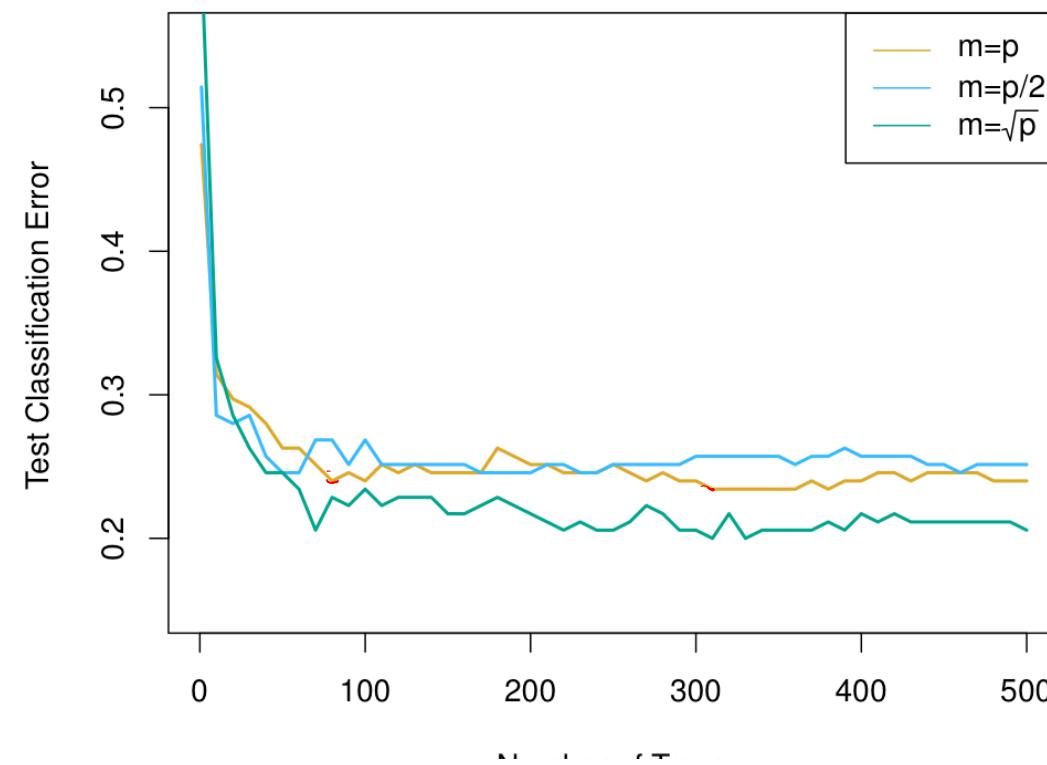
Advantages:

- Since at each split, considering less than half of features, this gives very different trees
- Very different trees in your ensemble decreases correlation, and gives more robust results
- Also, slightly cheaper because you don't have to evaluate each feature to choose best split

## Random forests

---

Look at expression of 500 different genes in tissue samples. Predict 15 cancer states



# Outline

- Ensemble methods
- Bagging and random forest
- General idea of boosting
- Adaboost

## Boosting intuition

---

Boosting is an ensemble method, but with a different twist.

Idea:

- Build a sequence of dumb models
- Modify training data along the way to focus on difficult to classify examples
- Predict based on weighted majority vote of all the models

Challenges:

- What do we mean by dumb?
- How do we promote difficult examples?
- Which models get more say in vote?

## Boosting intuition

What do we mean by dumb?

Each model in our sequence will be a weak learner

$$\text{err} = \frac{1}{m} \sum_{i=1}^m I(y_i \neq h(\mathbf{x}_i)) = \frac{1}{2} - \gamma, \gamma > 0$$

Most common weak learner in Boosting is a decision stump - a decision tree with a single split

## Boosting intuition

How do we promote difficult examples?

After each iteration, we'll increase the importance of training examples that we got wrong on the previous iteration and decrease the importance of examples that we got right on the previous iteration

Each example will carry around a weight  $w_i$  that will play into the decision stump and the error estimation

Weights are normalized so they act like a probability distribution

$$\sum_{i=1}^m w_i = 1$$

## Boosting intuition

Which models get more say in vote?

The models that performed better on training data get more say in the vote

For our sequence of weak learners:  $h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_K(\mathbf{x})$

Boosted classifier defined by

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

$$H(\mathbf{x}) = \text{sign} \left[ \sum_{k=1}^K \alpha_k h_k(\mathbf{x}) \right]$$

Weight  $\alpha_k$  is measure of accuracy of  $h_k$  on training data

# Outline

- Ensemble methods
- Bagging and random forest
- General idea of boosting
- Adaboost

## Adaboost

---

- Training set  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$
- Feature vector  $\mathbf{x} \in \mathbb{R}^d$
- Labels are  $y_i \in \{-1, 1\}$

## Adaboost

---

- 1 . Initialize data weights to  $w_i = \frac{1}{m}$ ,  $i = 1, \dots, m$
- 2 . For  $k = \underline{1}$  to  $K$ :
  - (a) Fit classifier  $h_k(\mathbf{x})$  to training data with weights  $w_i$
  - (b) Compute weighted error  $\text{err}_k = \frac{\sum_{i=1}^m w_i I(y_i \neq h_k(\mathbf{x}_i))}{\sum_{i=1}^m w_i}$
  - (c) Compute  $\alpha_k = \frac{1}{2} \log((1 - \text{err}_k)/\text{err}_k)$
  - (d) Set  $w_i \leftarrow \frac{w_i}{Z_k} \cdot \exp[-\alpha_k y_i h_k(\mathbf{x}_i)]$ ,  $i = 1, \dots, m$
- 3 . Output  $H(\mathbf{x}) = \text{sign} \left[ \sum_{k=1}^K \alpha_k h_k(\mathbf{x}) \right]$

## Adaboost

---

1 . Initialize data weights to  $w_i = \frac{1}{m}$ ,  $i = 1, \dots, m$

Weights are initialized to uniform distribution. Every training example counts equally on first iteration.

## Adaboost

---

2a. Fit classifier  $h_k(\mathbf{x})$  to training data with weights  $w_i$

Decide split based on info gain with weighted entropy:

$$\underline{H(D) = -p \log_2 p - (1-p) \log_2(1-p)}$$

$$I(D, a) = \underline{H(D)} - \sum_{v \in vals(a)} \frac{|\{x \in D | x_a = v\}|}{|D|} H(\{x \in D | x_a = v\})$$

$$p = \frac{\sum_{i=1}^m \underline{w_i} \cdot I(y_i = 1)}{\sum_{i=1}^m w_i \cdot I(y_i = \pm 1)}$$

## Adaboost

---

2b. Compute weighted error

$$\text{err}_k = \frac{\sum_{i=1}^m w_i I(y_i \neq h_k(\mathbf{x}_i))}{\sum_{i=1}^m w_i}$$

## Adaboost

---

2b. Compute weighted error

$$\text{err}_k = \frac{\sum_{i=1}^m w_i I(y_i \neq h_k(\mathbf{x}_i))}{\sum_{i=1}^m w_i}$$

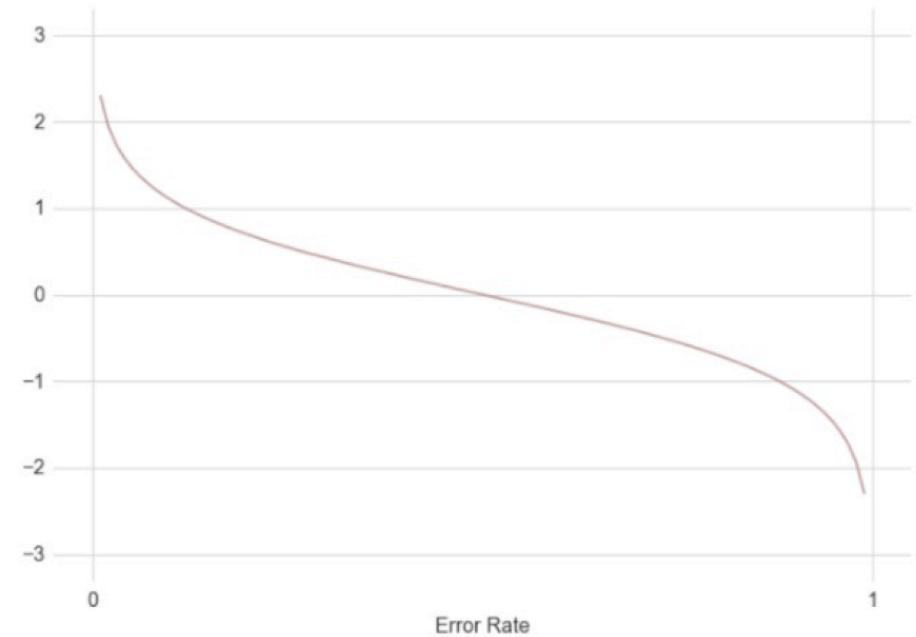
Still gives error rate in [0, 1]

## Adaboost

2c. Compute  $\alpha_k = \frac{1}{2} \log((1 - \text{err}_k)/\text{err}_k)$

Models with small err get promoted  
(exponentially)

Models with large err get demoted  
(exponentially)



## Adaboost

---

$$\check{A}_k > 0$$

2d. Set  $w_i \leftarrow \frac{w_i}{Z_k} \cdot \exp[-\alpha_k y_i h_k(\mathbf{x}_i)], i = 1, \dots, m$

If example was misclassified weight goes up

If example was classified correctly weight goes down

How big of a jump depends on accuracy of model

$Z_k$  is just a normalizing constant to ensure that the  $w$ 's are a distribution

$$y_i h_k(x_i) < 0$$

$$y_i h_k(x_i) > 0$$

$$\exp(-2y_i h_k(x_i)) > 0$$

$$\exp(-2y_i h_k(x_i)) < 1$$

## Adaboost

---

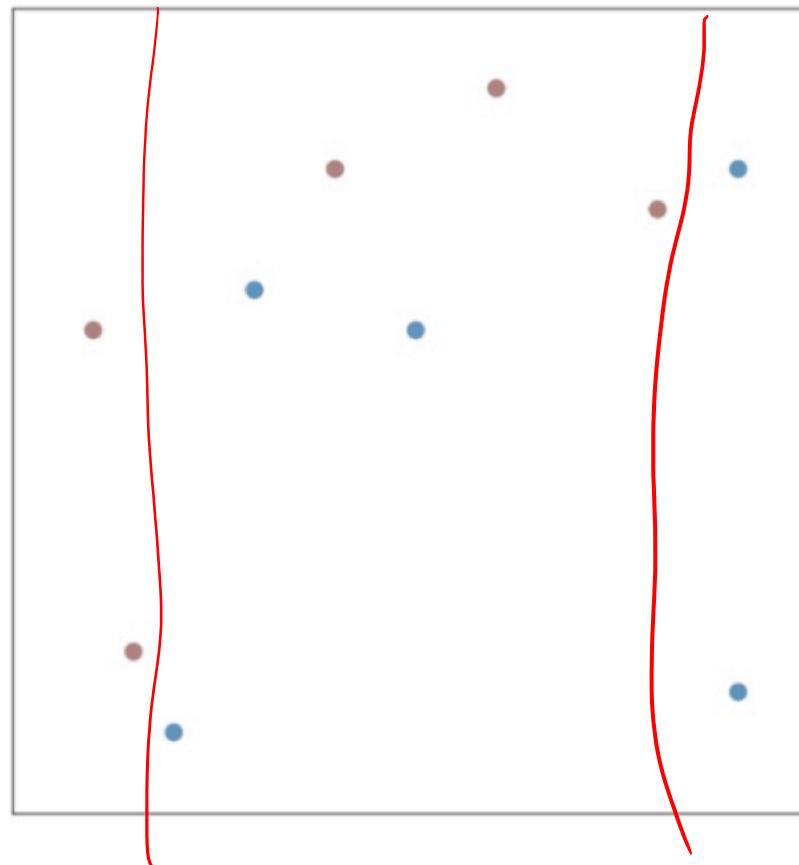
$$3. \text{ Output } H(\mathbf{x}) = \text{sign} \left[ \sum_{k=1}^K \alpha_k h_k(\mathbf{x}) \right]$$

Sum up weighted votes from each model

Classify  $y = 1$  if positive and  $y = -1$  if negative

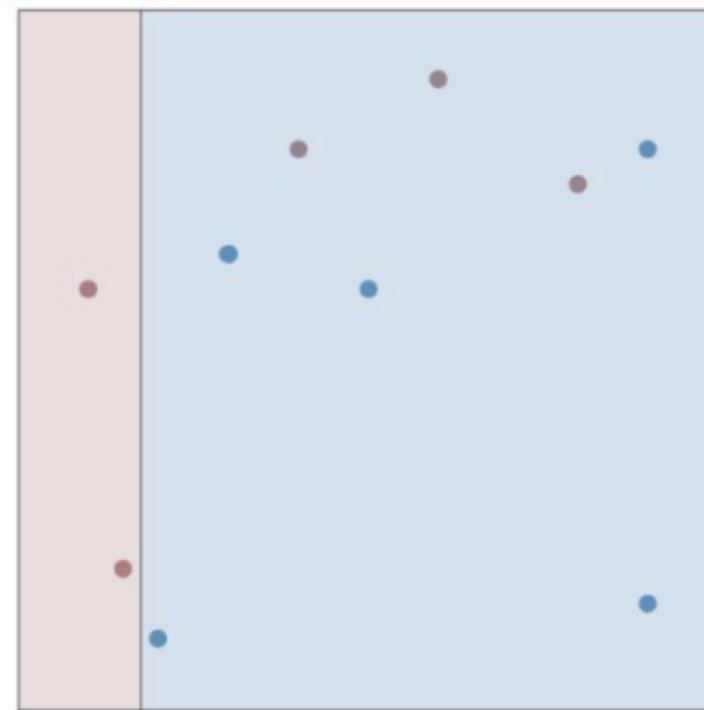
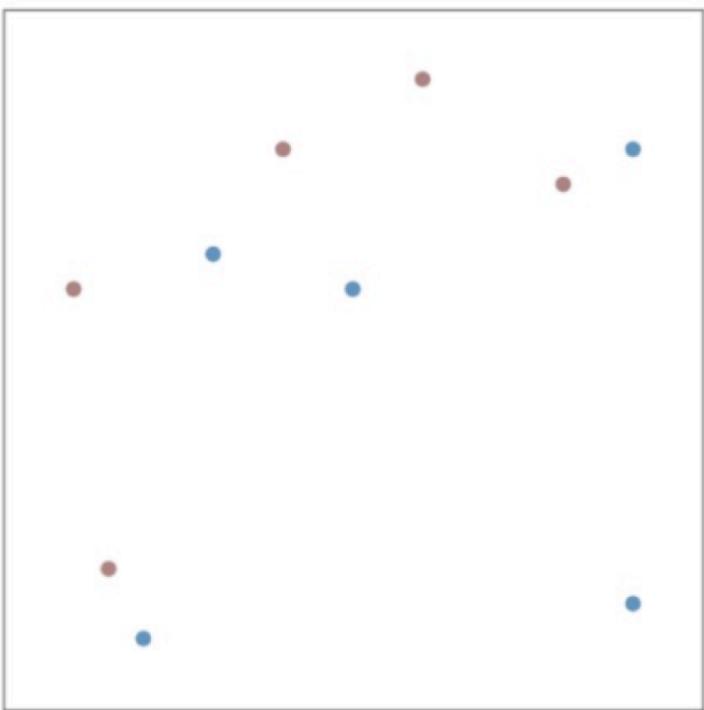
## Adaboost example

Suppose we have the following training data



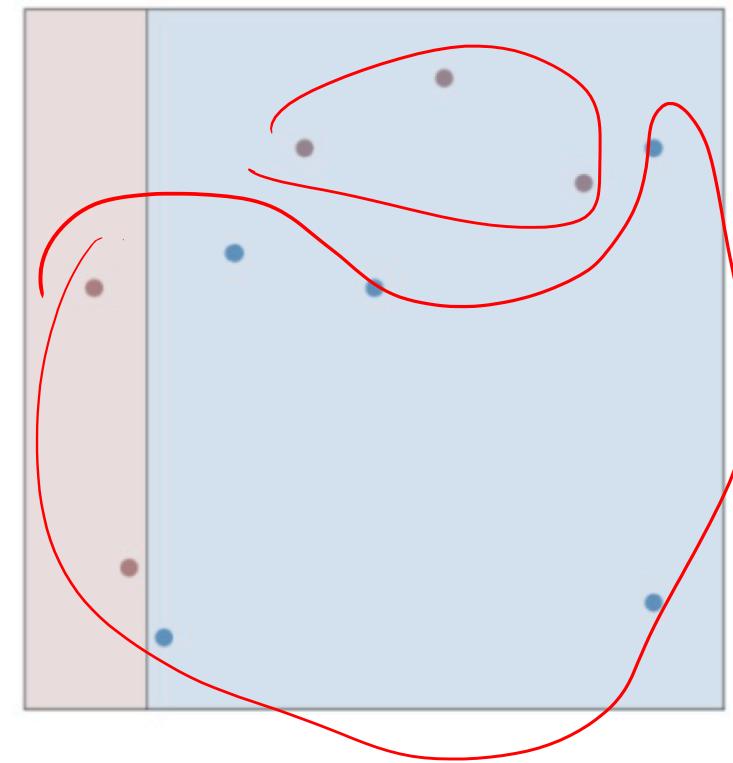
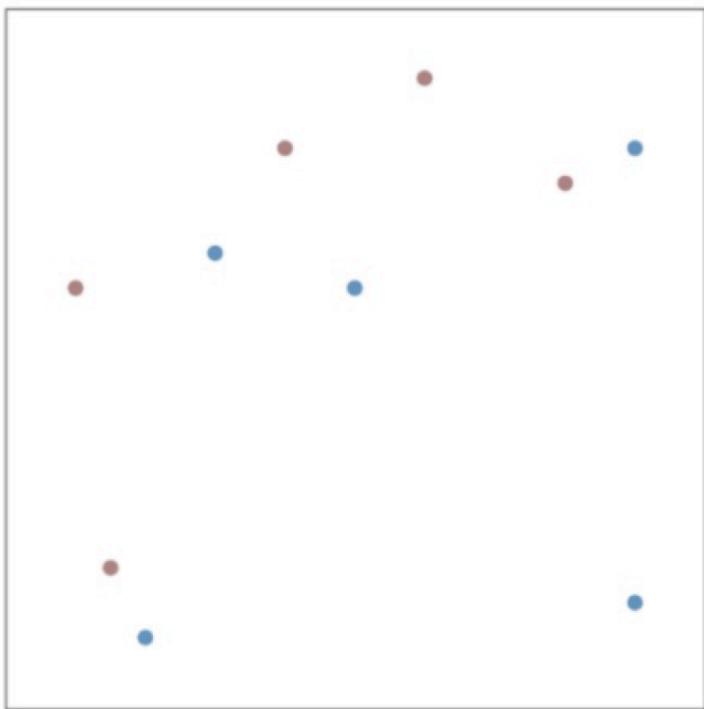
## Adaboost example

First decision stump



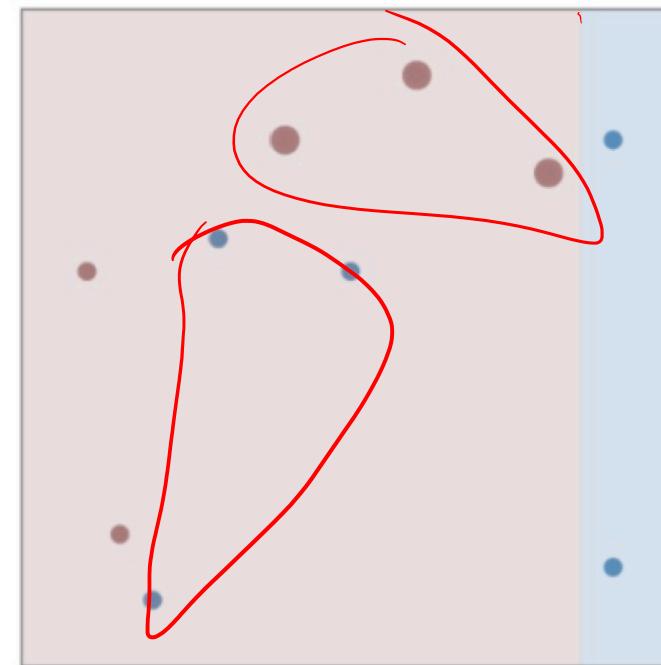
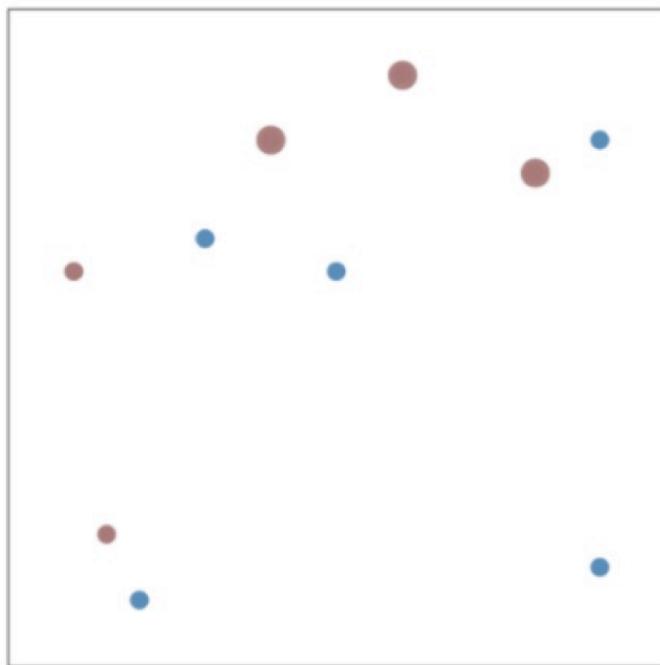
## Adaboost example

First decision stump ,  $err_1 = 0.3, \alpha_1 = 0.42$



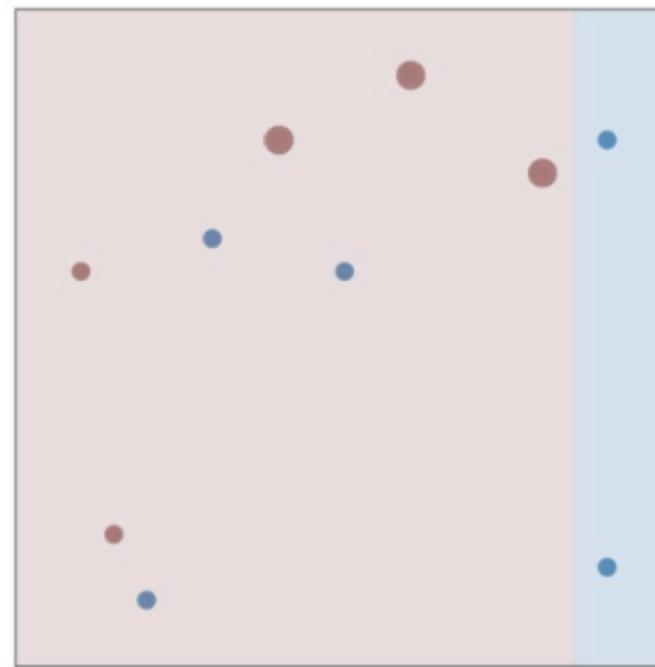
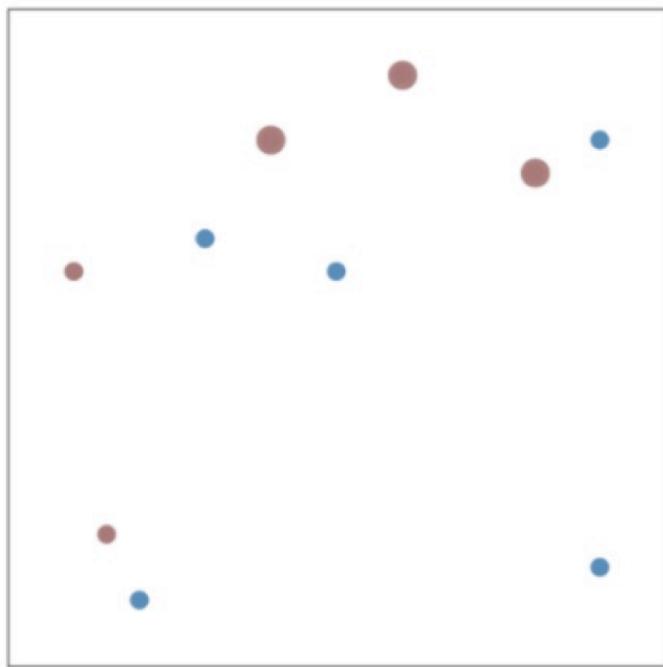
## Adaboost example

Second decision stump



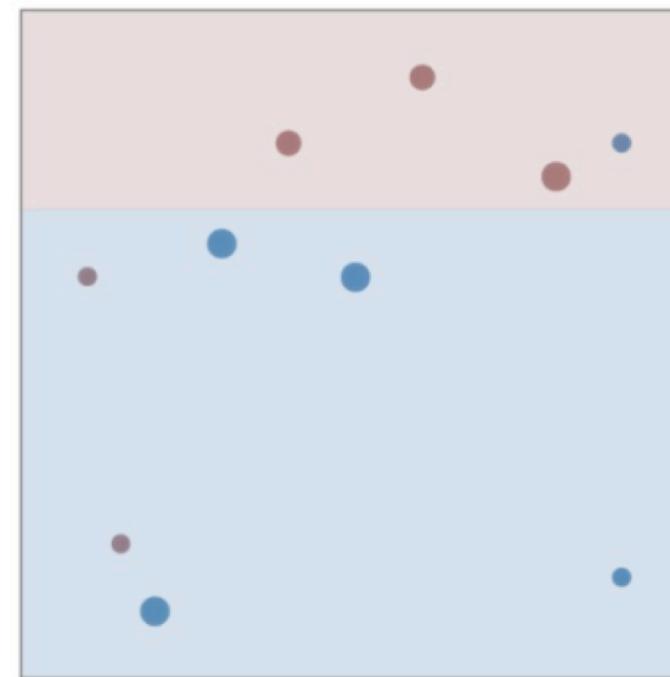
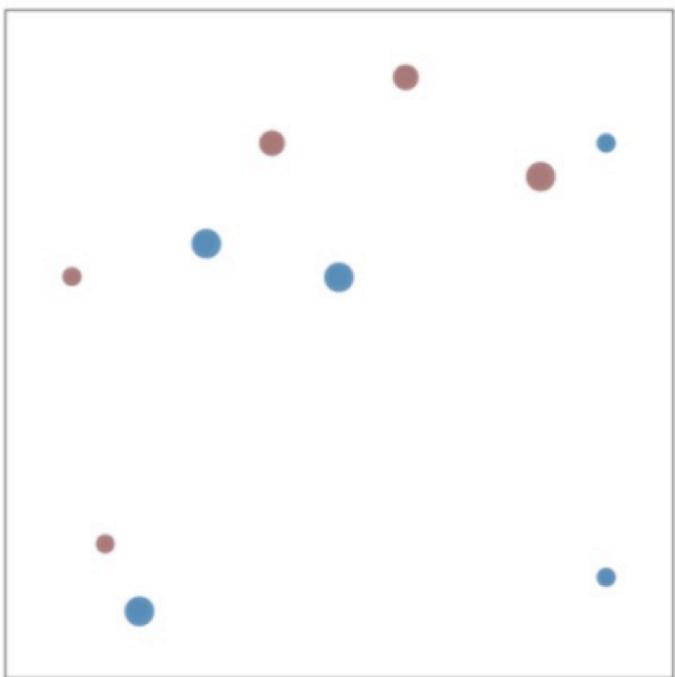
## Adaboost example

Second decision stump ,  $err_2 = 0.21, \alpha_1 = 0.65$



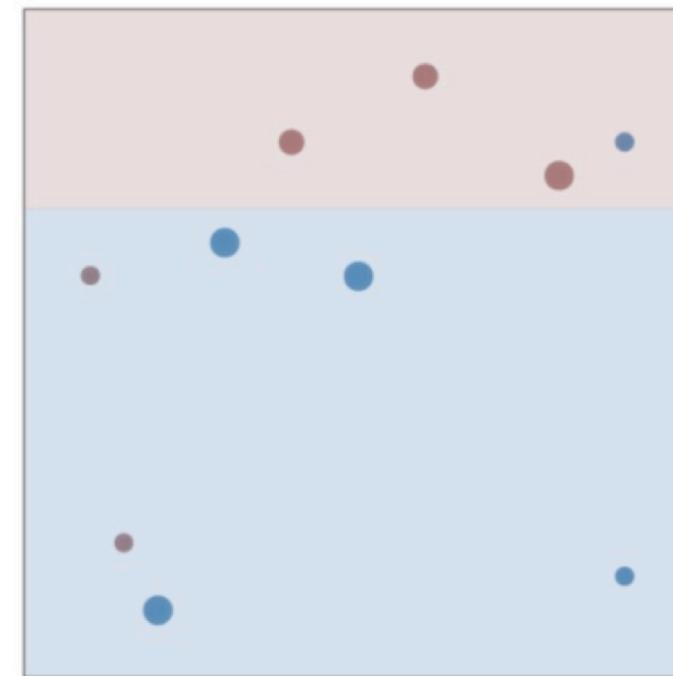
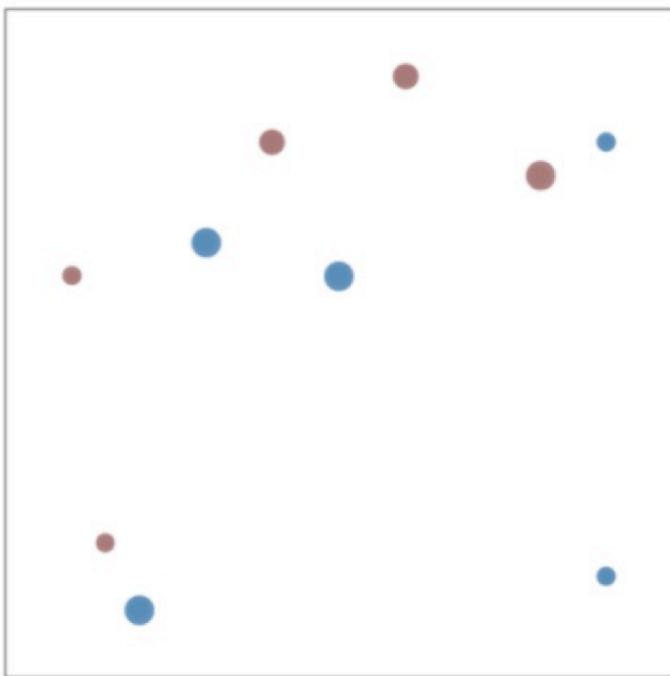
## Adaboost example

Third decision stump

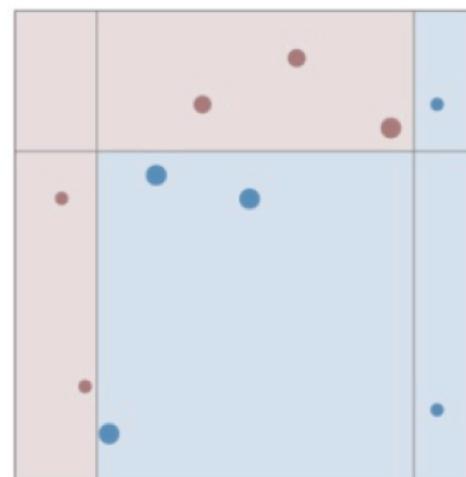
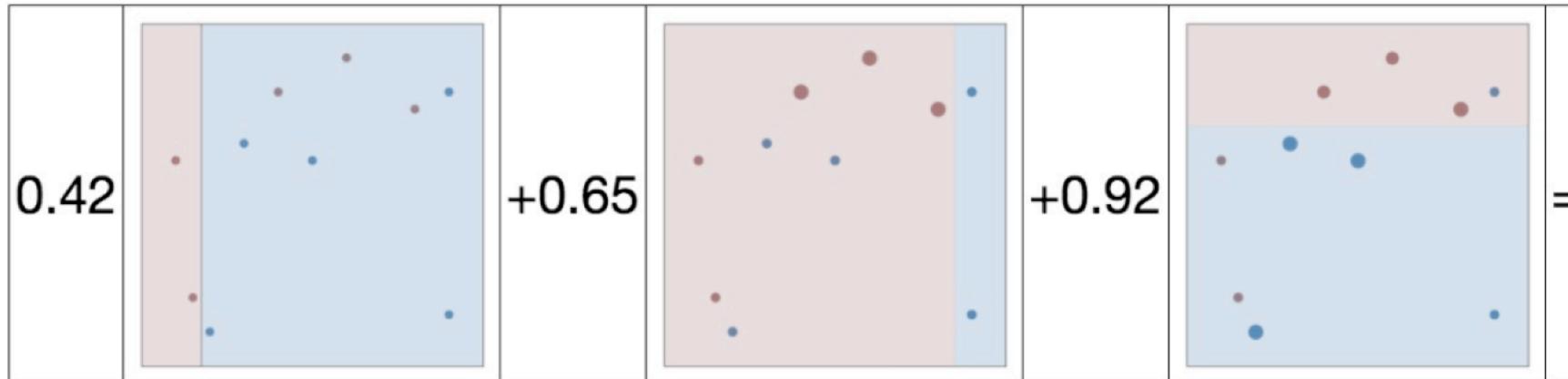


## Adaboost example

Third decision stump ,  $err_3 = 0.14, \alpha_3 = 0.92$



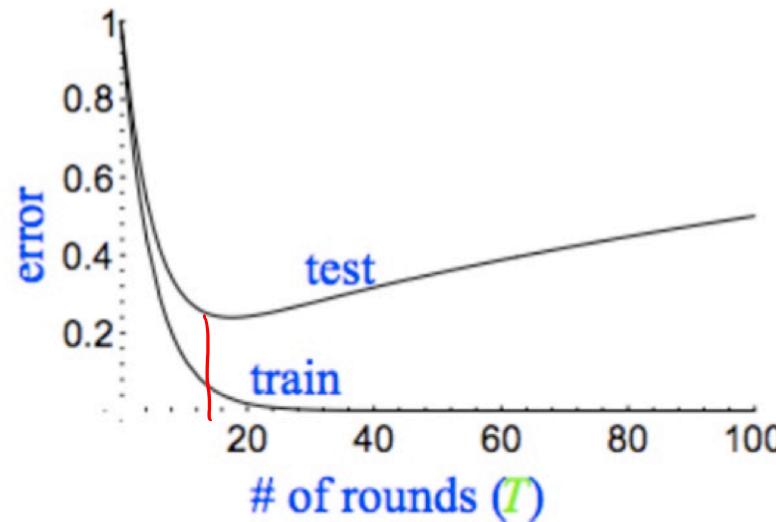
## Adaboost example



## Generalization performance

---

Recall the standard experiment of measuring test and training error vs. model complexity

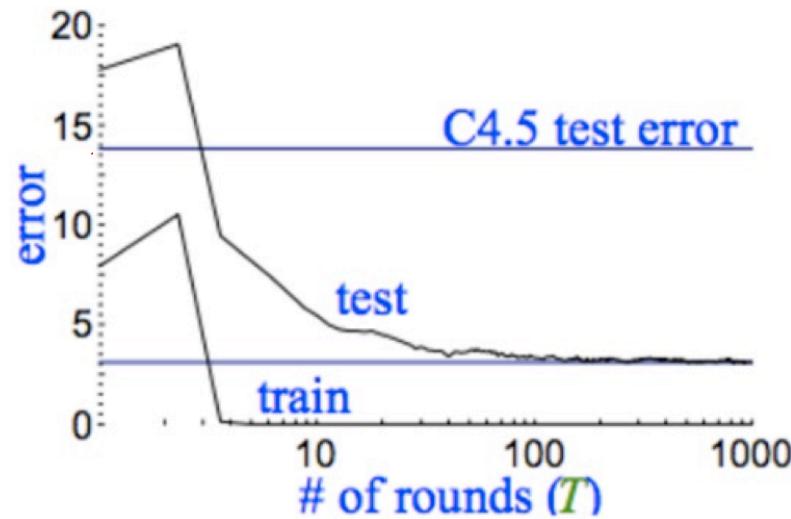


Once overfitting begins, test error goes up

## Generalization performance

---

Boosting has remarkably uncommon effect



Happens much slower with boosting

## The underlying math

So far this looks like a reasonable thing that just worked out  
But is there math behind it?

## The underlying math

---

Yep! It is minimization of a loss function, like always

Formulate the problem as finding a classifier  $H(\mathbf{x})$  such that

$$H^* = \arg \min_H \sum_{i=1}^m L(y_i, H(\mathbf{x}_i))$$

where here we take the loss function  $L$  to be the following exponential function

$$L = \exp[\underline{-y H(\mathbf{x})}]$$

Notice if  $y \neq H(\mathbf{x})$  we get a positive exponent and a large loss.

## The underlying math

---

Since we're doing this in an iterative way, we're going to assume a form of  $H(\mathbf{x})$  that is amenable to iterative improvement. Specifically

$$H(\mathbf{x}) = \sum_{k=1}^K \alpha_k \phi_k(\mathbf{x})$$

So the problem becomes to choose the optimal weights,  $\alpha$ , and optimal basis functions  $\phi_k$ .

For everything I will assume that we've already computed a good  $H_{k-1}$  and attempt to build a better  $H_k$ .

## The underlying math

At step  $k$  we have

$$L_k = \sum_{i=1}^m \exp [-y_i (\underline{H_{k-1}(\mathbf{x}_i)} + \underline{\alpha \phi(\mathbf{x}_i)})]$$

Taking the things we know out of the exponent gives

$$L_k = \sum_{i=1}^m \underline{w_{i,k}} \exp [-\underline{\alpha y_i \phi(\mathbf{x}_i)}], \quad w_{i,k} = \exp [\underline{-y_i H_{k-1}(\mathbf{x}_i)}]$$

Want to choose good  $\underline{\alpha}$  and  $\underline{\phi}$  to reduce loss

## The underlying math

Can rewrite as

$$L_k = e^\alpha \sum_{y_i \neq \phi(\mathbf{x}_i)} w_{i,k} + \underline{e^{-\alpha}} \sum_{y_i = \phi(\mathbf{x}_i)} w_{i,k}$$

.

Add zero in a fancy way

$\cancel{\cancel{0}}$

$$L_k = (\underline{e^\alpha - e^{-\alpha}}) \sum_{i=1}^m w_{i,k} I(y_i \neq \cancel{\cancel{\phi(\mathbf{x}_i)}}) + e^{-\alpha} \sum_{i=1}^m w_{i,k}$$

$$\lambda = \frac{1}{2} \log \left( \frac{1-\text{err}}{\text{err}} \right)$$

$$\begin{aligned} \frac{\partial L}{\partial \lambda} &= (e^\lambda + e^{-\lambda}) \sum_{i=1}^m w_{i,k} I(y_i \neq \phi(\mathbf{x}_i)) \\ &\quad - e^{-\lambda} \sum_{i=1}^m w_{i,k} = 0 \end{aligned}$$

$$\begin{aligned} (e^{2\lambda} + 1) \sum_{i=1}^m w_{i,k} I(y_i \neq \phi(\mathbf{x}_i)) &= \sum_{i=1}^m w_{i,k} \\ e^{2\lambda} + 1 &= \frac{1}{\text{err}} \quad \lambda = \frac{1}{2} \log \left( \frac{1}{\text{err}} - 1 \right) \end{aligned}$$

## The underlying math

Choose  $\phi$  and  $\alpha$  separately.

A good  $\phi$  would be one that minimizes weighted misclassifications

$$h_k = \arg \min_{\phi} w_{i,k} I(y_i \neq \phi(\mathbf{x}_i))$$

That's what our weak learner is for

## The underlying math

Plugging that into our Loss function gives

$$L_k = (e^\alpha - e^{-\alpha}) \sum_{i=1}^m w_{i,k} I(y_i \neq h_k(\mathbf{x}_i)) + e^{-\alpha} \sum_{i=1}^m w_{i,k}$$

Now we want to minimize w.r.t.  $\alpha$ . Take derivative, set equal to zero

$$0 = \frac{dL_k}{d\alpha} = (e^\alpha + e^{-\alpha}) \sum_{i=1}^m w_{i,k} I(y_i \neq h_k(\mathbf{x}_i)) - e^{-\alpha} \sum_{i=1}^m w_{i,k}$$

which gives  $e^{2\alpha} = \frac{\sum_i w_{i,k}}{\sum_i w_{i,k} I(y_i \neq h_k(\mathbf{x}_i))} - 1 = \frac{1}{\text{err}_k} - 1 = \frac{1 - \text{err}_k}{\text{err}_k}$

And finally  $\alpha_k = \frac{1}{2} \log \left( \frac{1 - \text{err}_k}{\text{err}_k} \right)$

## The underlying math

What about the weight update?

Remember we got the weights by pulling the already computed function out of  $L$ .  
For the new weights, we have

$$w_{i,k+1} = \exp[-y_i H_k(\mathbf{x}_i)] = \exp[-y_i \underline{H_{k-1}(\mathbf{x}_i)} - \underline{\alpha_k y_i h_k(\mathbf{x}_i)}] = \underline{w_{i,k}} \exp[-\underline{\alpha_k y_i h_k(\mathbf{x}_i)}]$$

which gives the update

$$w_i \leftarrow w_i \exp[-\alpha_k y_i h_k(\mathbf{x}_i)]$$

And finally,  $H_K(\mathbf{x}) = \text{sign} \left[ \sum_{k=1}^K \alpha_k h_k(\mathbf{x}) \right]$

This is exactly Adaboost

# Recap

- Random forests are easy to use and efficient
- Adaboost is supported by nice mathematical foundations
- Using many models usually outperforms a single model if we only care about performance