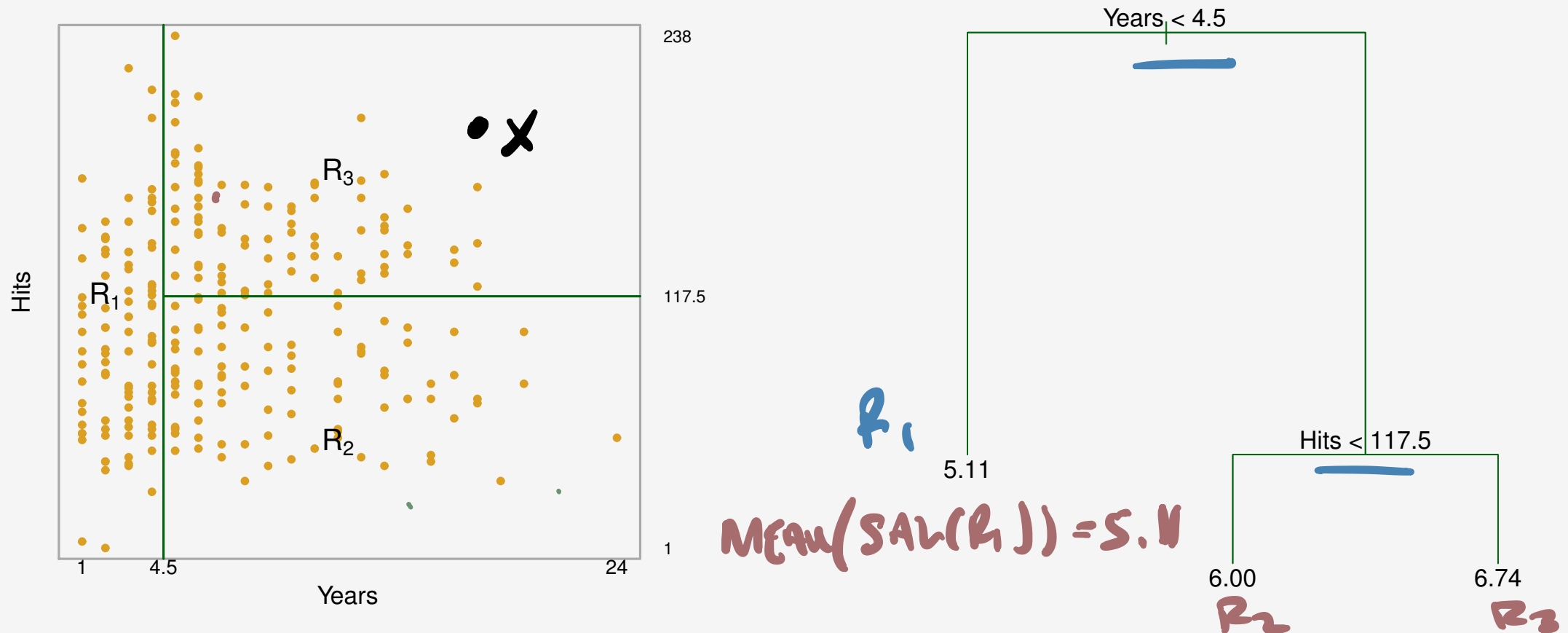# Bootstrapped Aggregation and Random Forests

# Bonus: Regression Trees

o Suppose you want to PREDICT the salary of a MLB player based on two features: how many hits they average a year and how many years they've been in the league

# Bonus: Regression Trees

o Suppose you want to PREDICT the salary of a MLB player based on two features: how many hits they average a year and how many years they've been in the league
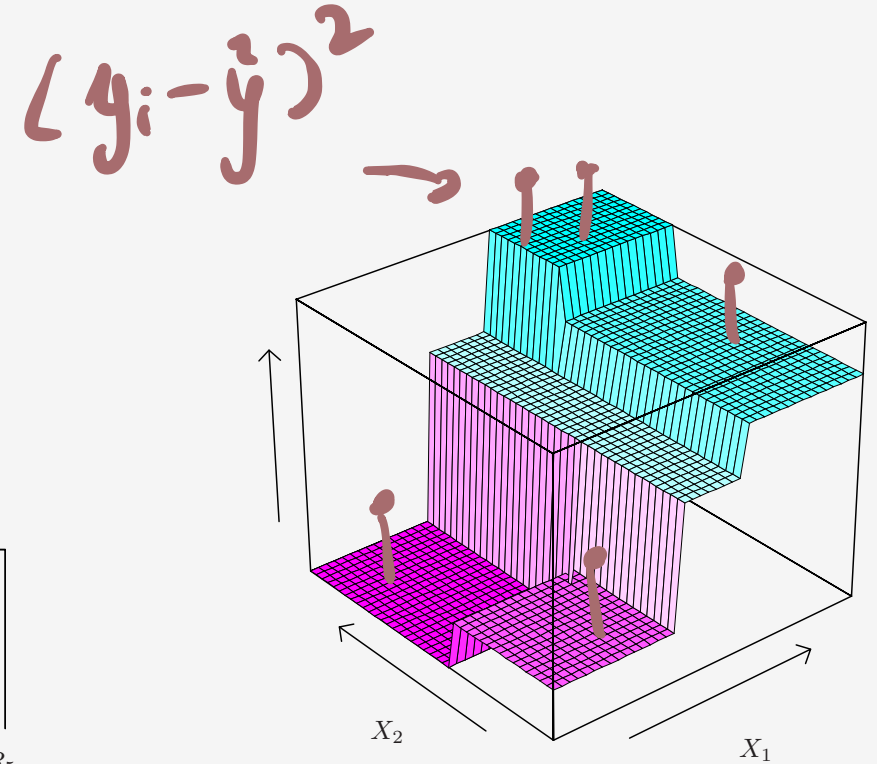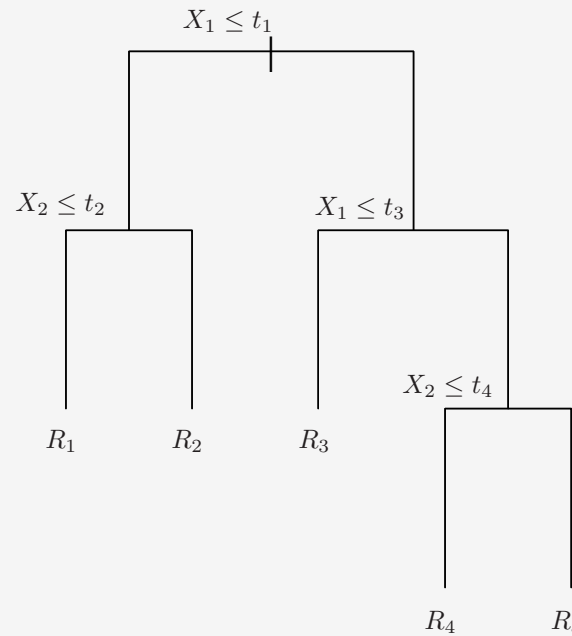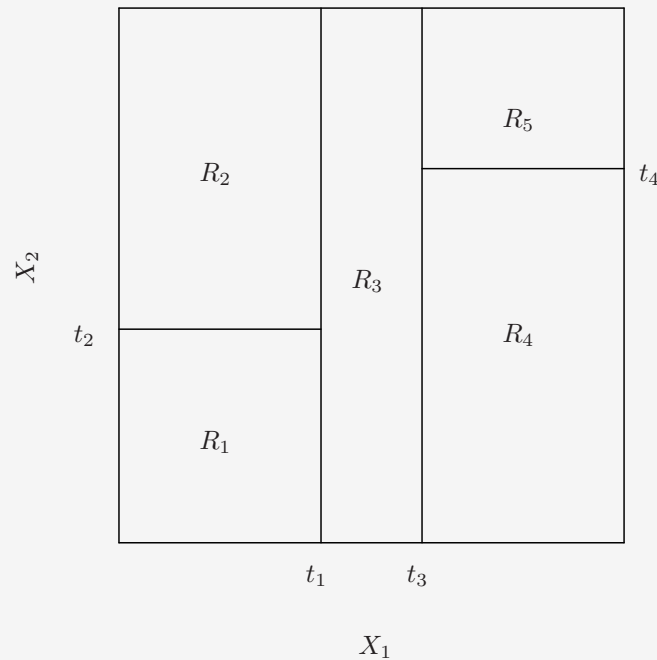
o Perform the usual binary splitting by feature

o Instead of classification, predict the response based on average response in leaf node

# Bonus: Regression Trees

o Instead of classification, predict the response based on average response in leaf node

$$\sum (y_i - \hat{y})^2$$

# Bonus: Regression Trees

o  What measure do we use to decide which feature to split on?

o  The goal is to find boxes that minimize the $RSS = \sum\limits_{j=1}^{J} \sum\limits_{i \in R_j} \left(y_i - \hat{y}_{R_j}\right)^2$
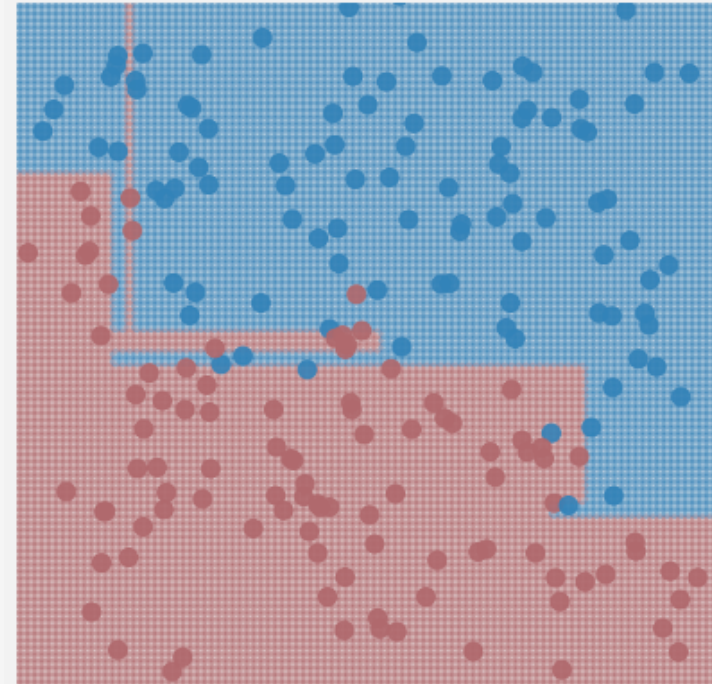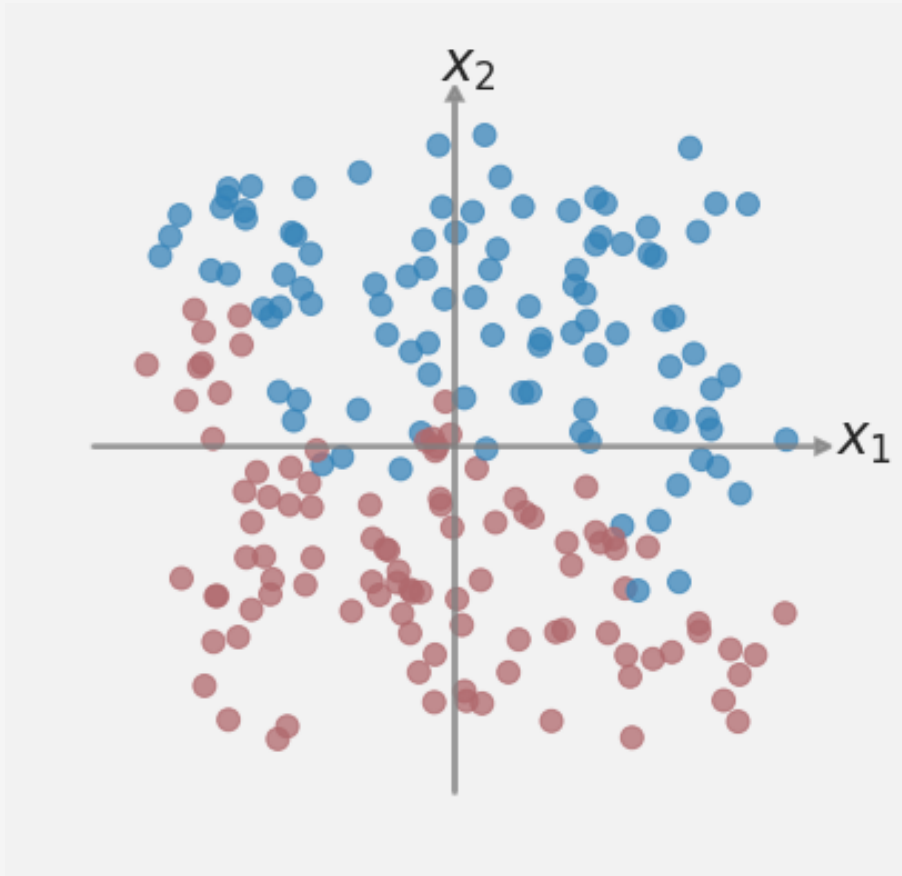
o  Consider splitting the data on a node into two boxes. Choose feature and split to minimize

$$\sum\limits_{i: \ x_i \in R_1(j,s)} \left(y_i - \hat{y}_{R_1}\right)^2 + \sum\limits_{i: \ x_i \in R_2(j,s)} \left(y_i - \hat{y}_{R_2}\right)^2$$

where  $R_1(j,s) = \{X \mid X_j < s\}$ and $R_2(j,s) = \{X \mid X_j \geq s\}$

*(handwritten annotations: #REGIONS, Dprure)*

# Previously on CSCI 4622

Last time we saw that full Decision Tree Classifiers tended to overfit

# Previously on CSCI 4622

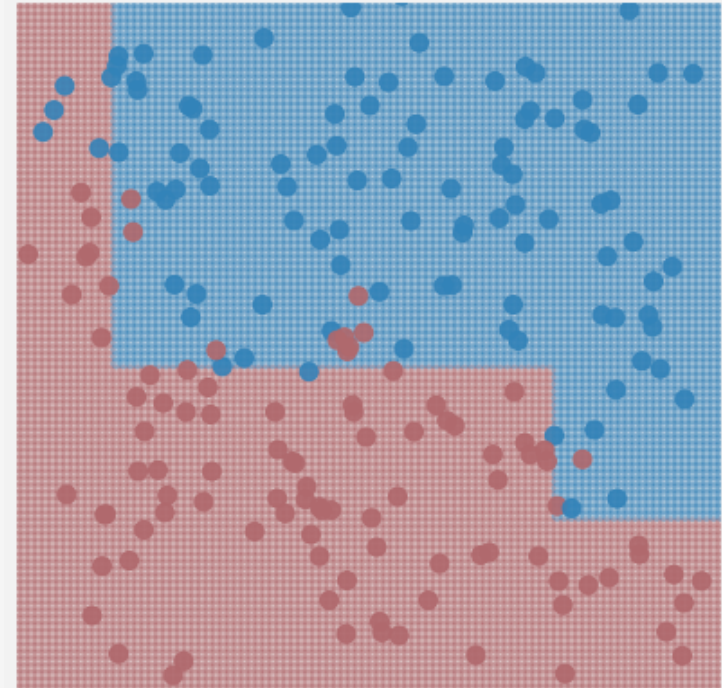Last time we saw that full Decision Tree Classifiers tended to overfit

We could alleviate this a bit with pruning

**Prepruning**:

o Don't let the tree have too many levels

o Stopping conditions

**Postpruning**:

o Build the complete tree

o Go back and remove vertices that don't affect performance too much

# Ensemble Methods

Today we'll see how we can use many decision trees to do even better

**Ensemble methods** are very popular in Machine Learning of late

**General Idea**:

o  Train numerous slightly different classifiers

o  Use each classifier to make a prediction on a query point

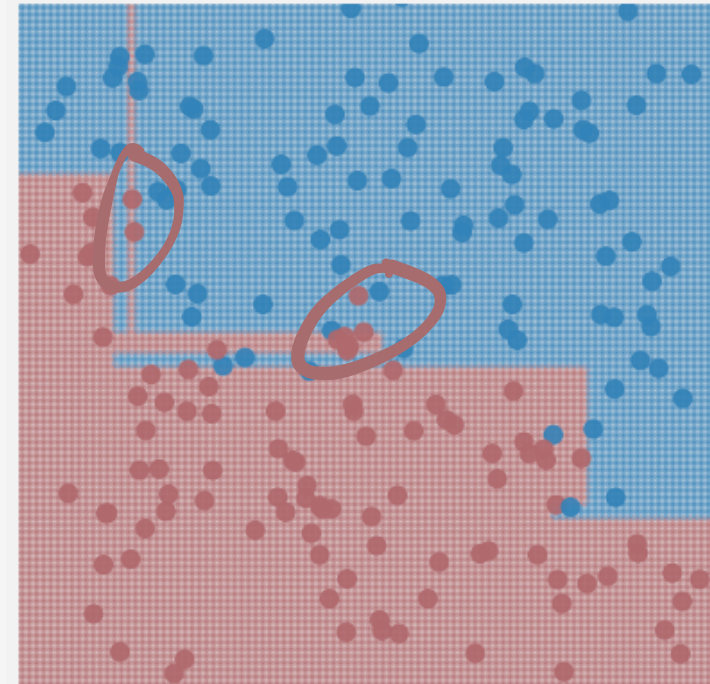o  Make the ensembled prediction by taking majority vote from individual classifiers

# Bagging

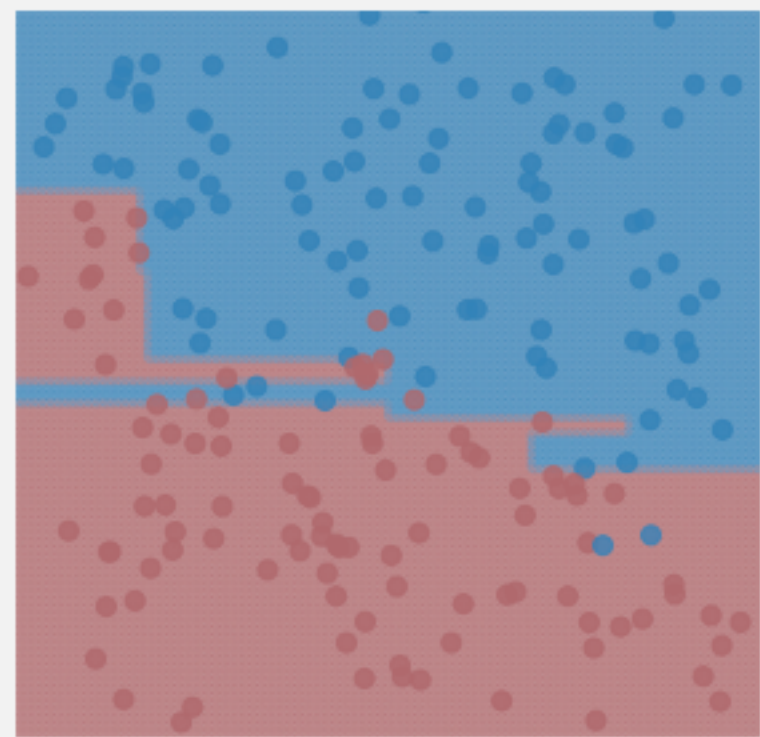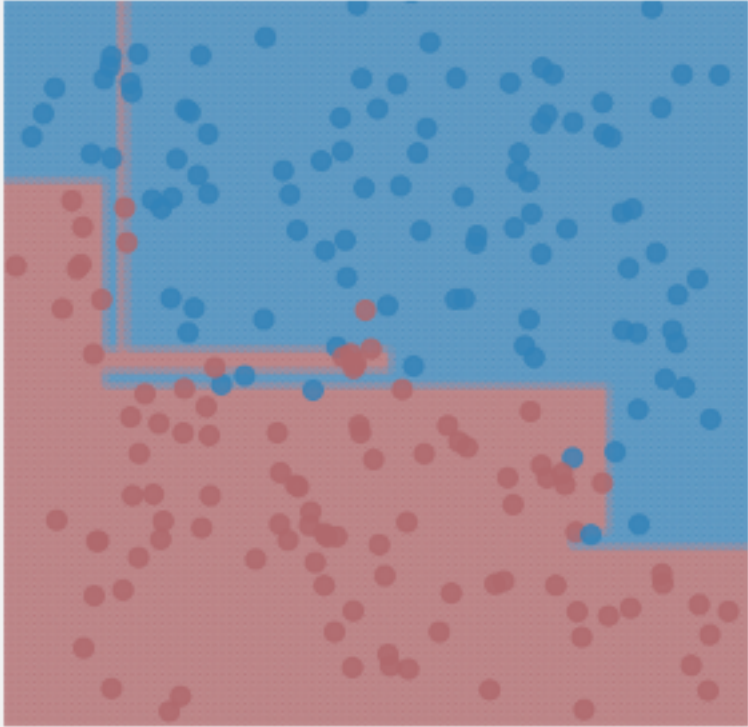This week we'll look at three different types of Ensembled Decision Tree Classifiers

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

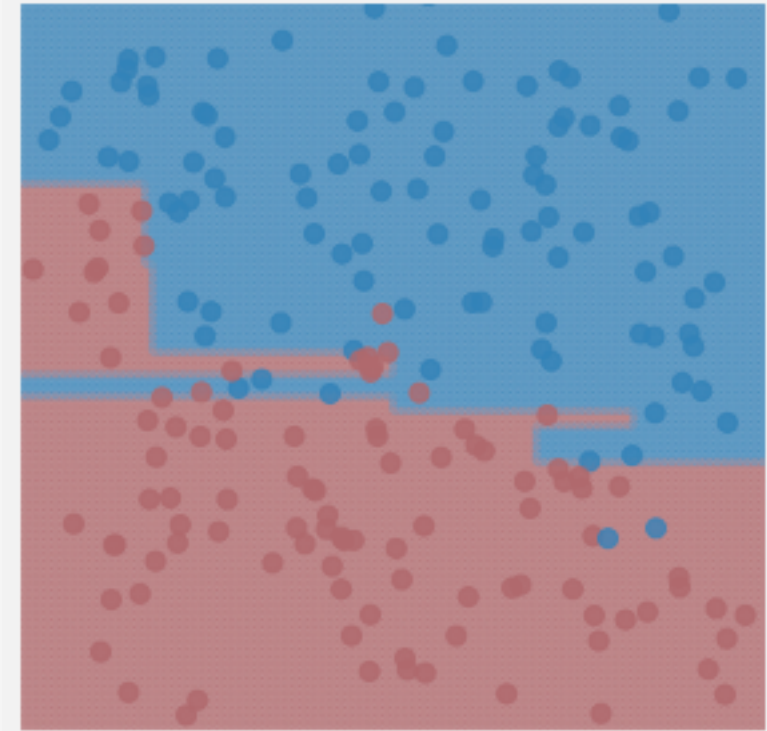The simplest is called **Bootstrapped Aggregation** (or Bagging)

# Bagging

*n TRAINING EXS*
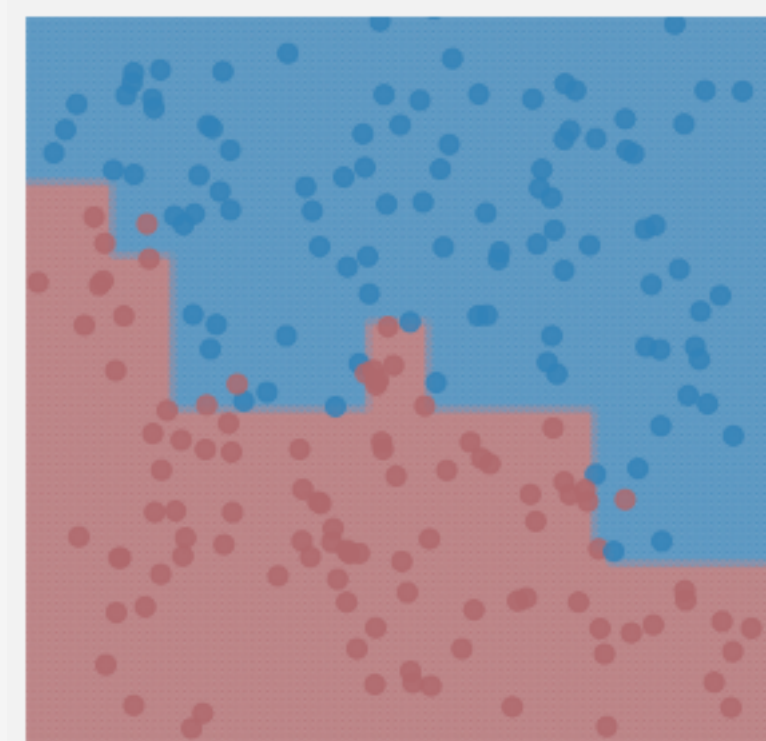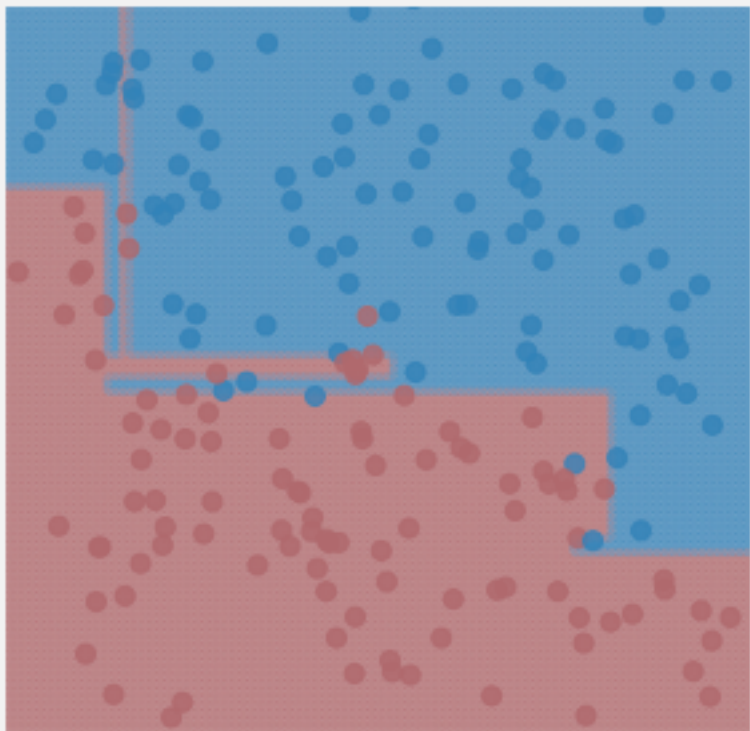
The simplest is called **Bootstrapped Aggregation** (or Bagging)

o   Sample n training examples **with replacement**

o   Fit full-depth Decision Tree Classifier

o   Repeat many times

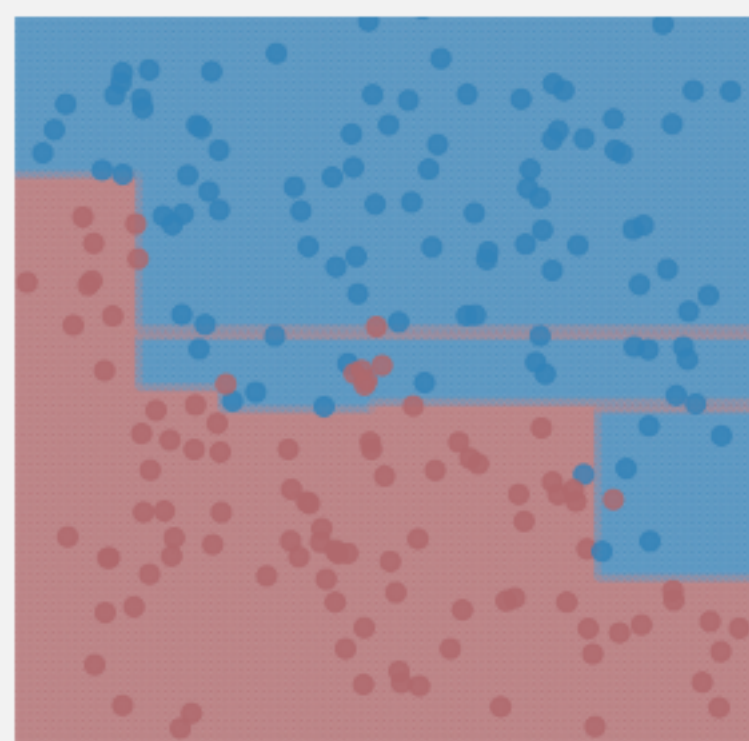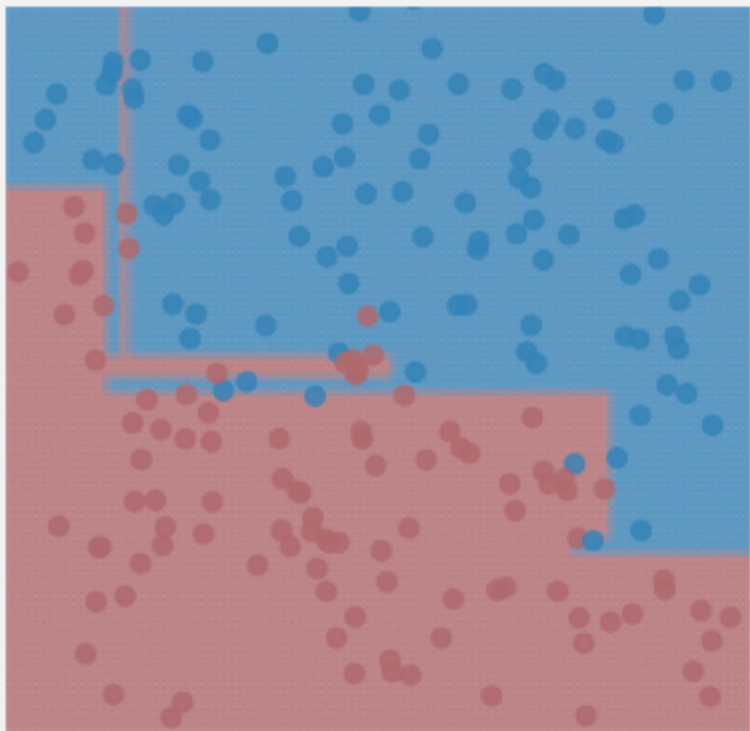o   Aggregate results by majority vote

# Bagging

The simplest is called **Bootstrapped Aggregation** (or Bagging)

# Bagging

The simplest is called **Bootstrapped Aggregation** (or Bagging)
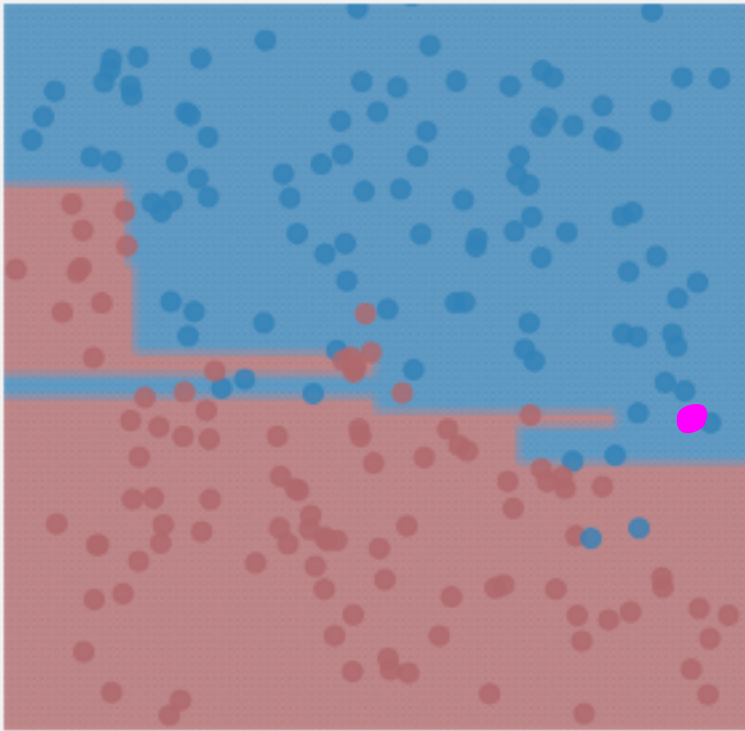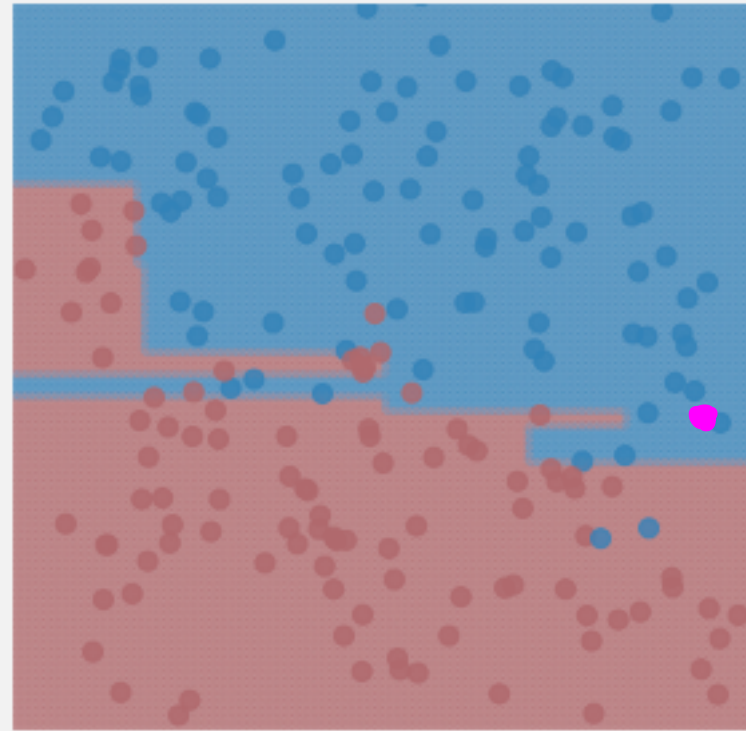
# Bagging

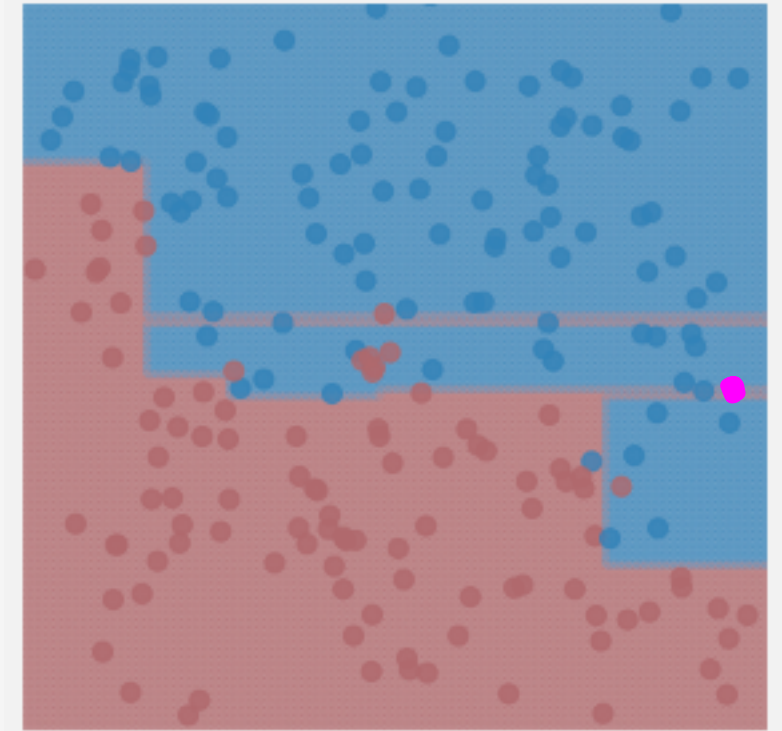The simplest is called **Bootstrapped Aggregation** (or Bagging)



B

B

R

# Bootstrap Refresher

INCOMES

10k
10k
20k
20k
20k
30k

40k



$$\frac{2}{7} \quad \frac{3}{7} \quad \frac{1}{7} \quad \frac{1}{7} \Rightarrow \sum = 1$$

10  20  30  40

EMPIRICAL DIST.

$R_1 = 10, 20, 20, 20, 20, 20$
$R_2 = 10, 10, 30, 30, 10, 20$

SAMPLE

POP

$[ \quad \bar{x} \quad ]$

# Bagging

In general, assume that we define $y \in \{-1, 1\}$ and each individual DCT has $h_k(\mathbf{x})$

Then we can define the Bagged Classifier as

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

How do we Validate the Bagged model?

o Could use cross-validation or a held-out validation set, as usual

o But there is a better way.

$1, 1, -1, -1, 1, 1$

$\Sigma = 2$

$\text{Sign}(2) = 1$

# Out-of-Bag Error Estimation

o Trees are repeatedly fit on bootstrapped re-samples of the training set

o Can be shown that around 2/3 of training examples present in each bootstrapped re-sample

o OOB Error Estimation makes use of the training examples NOT in the re-sample

**Here's How:**

$$x_i \Rightarrow sign\left[\sum_{\substack{k=1 \\ k \notin R_i}}^{500} h_k(x_i)\right] \Rightarrow +1 \text{ or } -1$$

# Bagging Pros and Cons

**Pros**:

o   Results in a much lower variance classifier than a single Decision Tree

**Cons**:

o   Results in a much less interpretable model

We essentially give up some interpretability in favor of better prediction accuracy

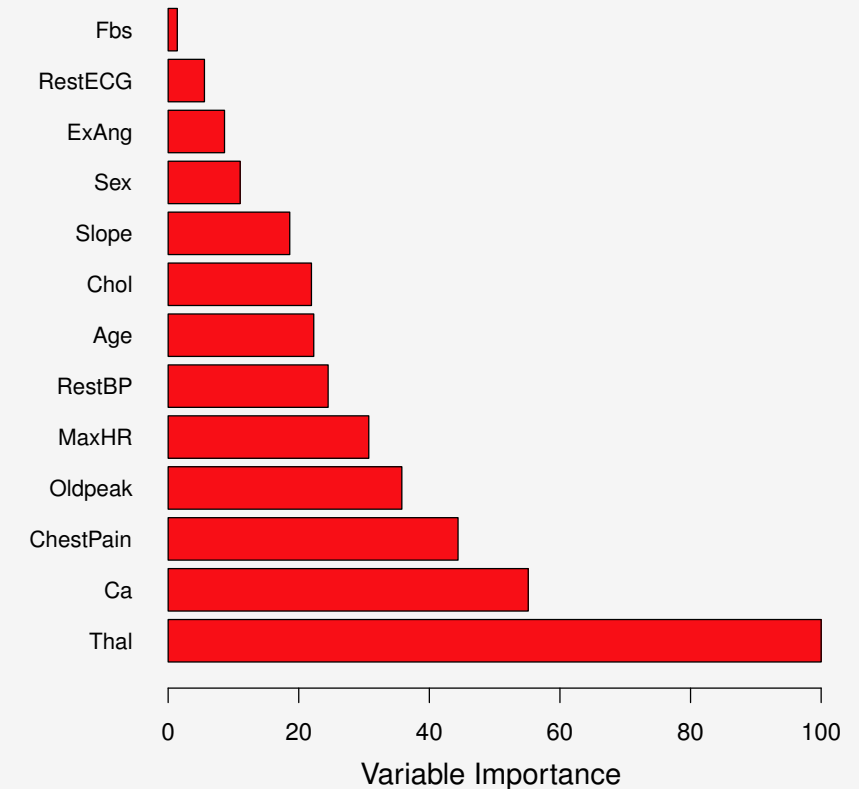But we can still get **insight** into what our model is doing using Bagging

# Bagging and Feature Importance

In a single Decision Tree Classifier we can get an idea of feature importance by

o Measuring Information Gain / reduction in Impurity is achieved by splitting on feature

o Compare these values at the end and rank features

With Bagging we can still do this.

We just compute these measures over all trees and average

# An Even Better Approach

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
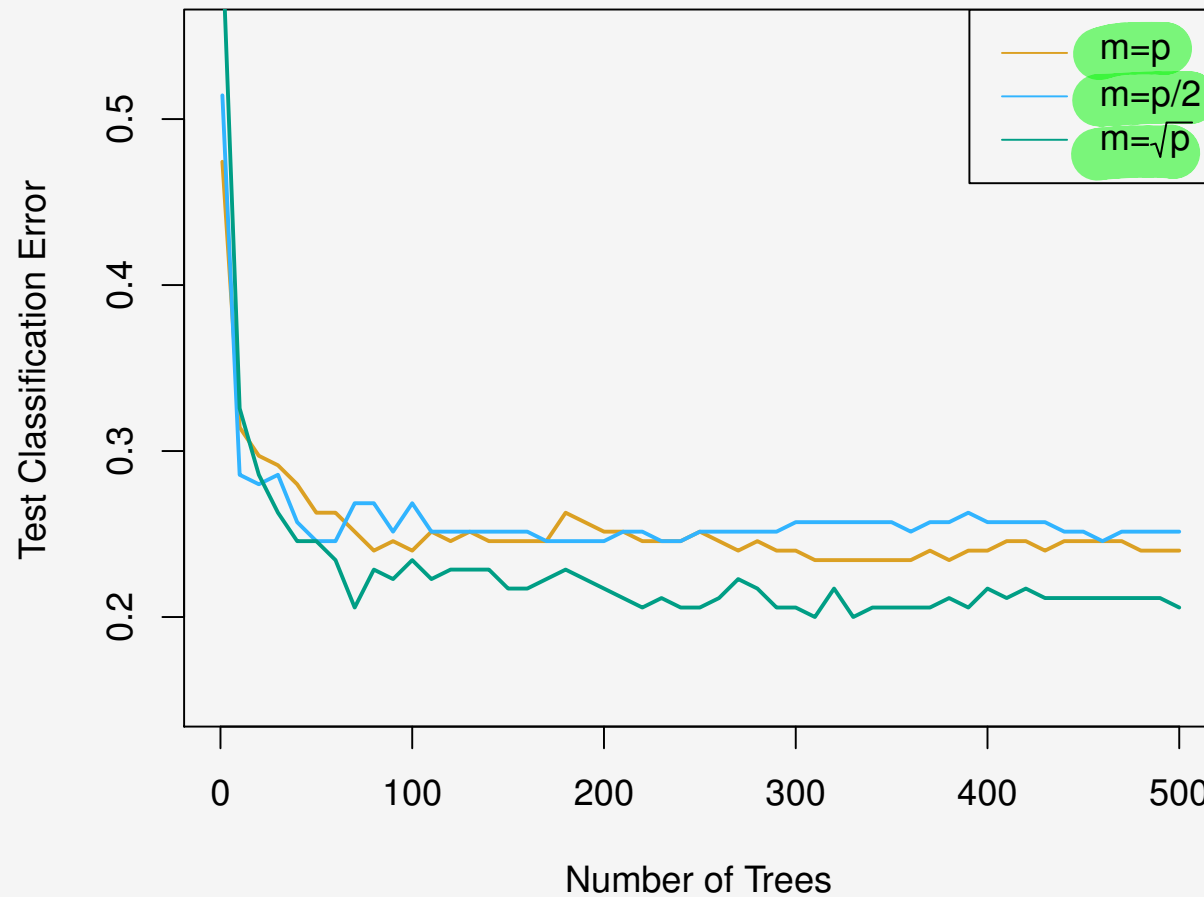
# Random Forests

o Still doing Bagging, in the sense that we fit bootstrapped resamples of training data

o BUT, every time we have to choose a feature to split on, don't consider all p features

o Instead, consider only $m \approx \sqrt{p}$ features to split on

**Advantages**:

o Since at each split, considering less than half of features, this gives very different trees

o Very different trees in your ensemble decreases correlation, and gives more robust results

o Also, slightly cheaper because you don't have to evaluate each feature to choose best split

# Practical Results

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



[James, et al]

# Where Do We Go From Here?

o Unfortunately, Decision Tree Classifiers are prone to overfitting

o On their own, tend to do worse than other classifiers we've seen

But **WITH THEIR POWERS COMBINED**!

Next up, Ensemble Methods.

o Take various weaker classifiers, and combine them to get strong classifiers

o Bagging and Random Forests

o Boosted Decision Trees and the AdaBoost Algorithm