

Tree-based Models and Ensembles

Lecture 17

Supervised Learning Techniques

Covered so far

K-Nearest Neighbors

Linear regression

Perceptron

Logistic Regression

Fisher's Linear Discriminant / Linear Discriminant Analysis

Quadratic Discriminant Analysis

Naïve Bayes

Decision Trees

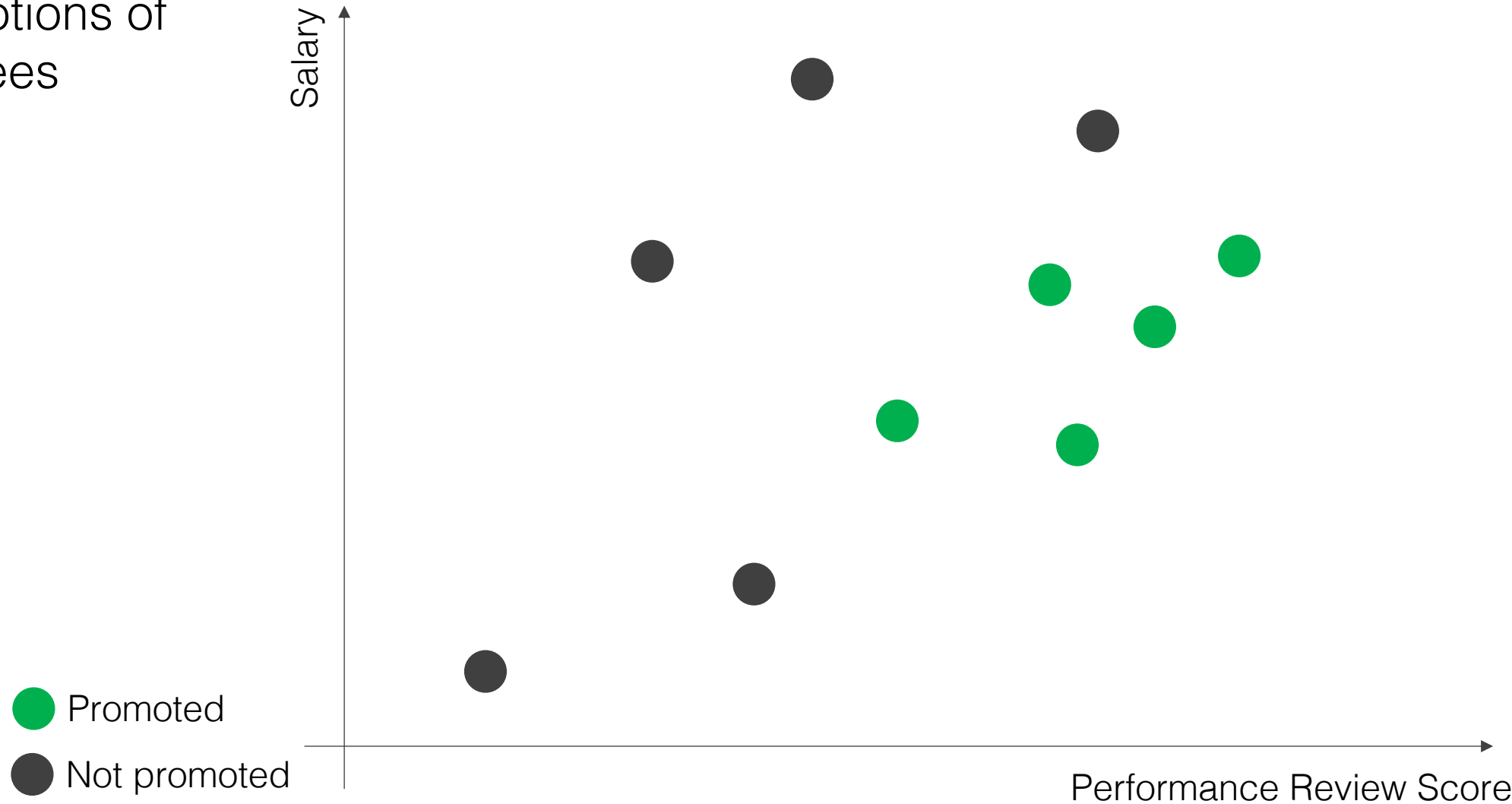
Ensemble methods (bagging and boosting)

Rely on a linear combination of weights and features: $\mathbf{w}^T \mathbf{x}$

Classification and Regression Trees (CART)

Classification trees = decision trees

Predicting promotions of
salaried employees



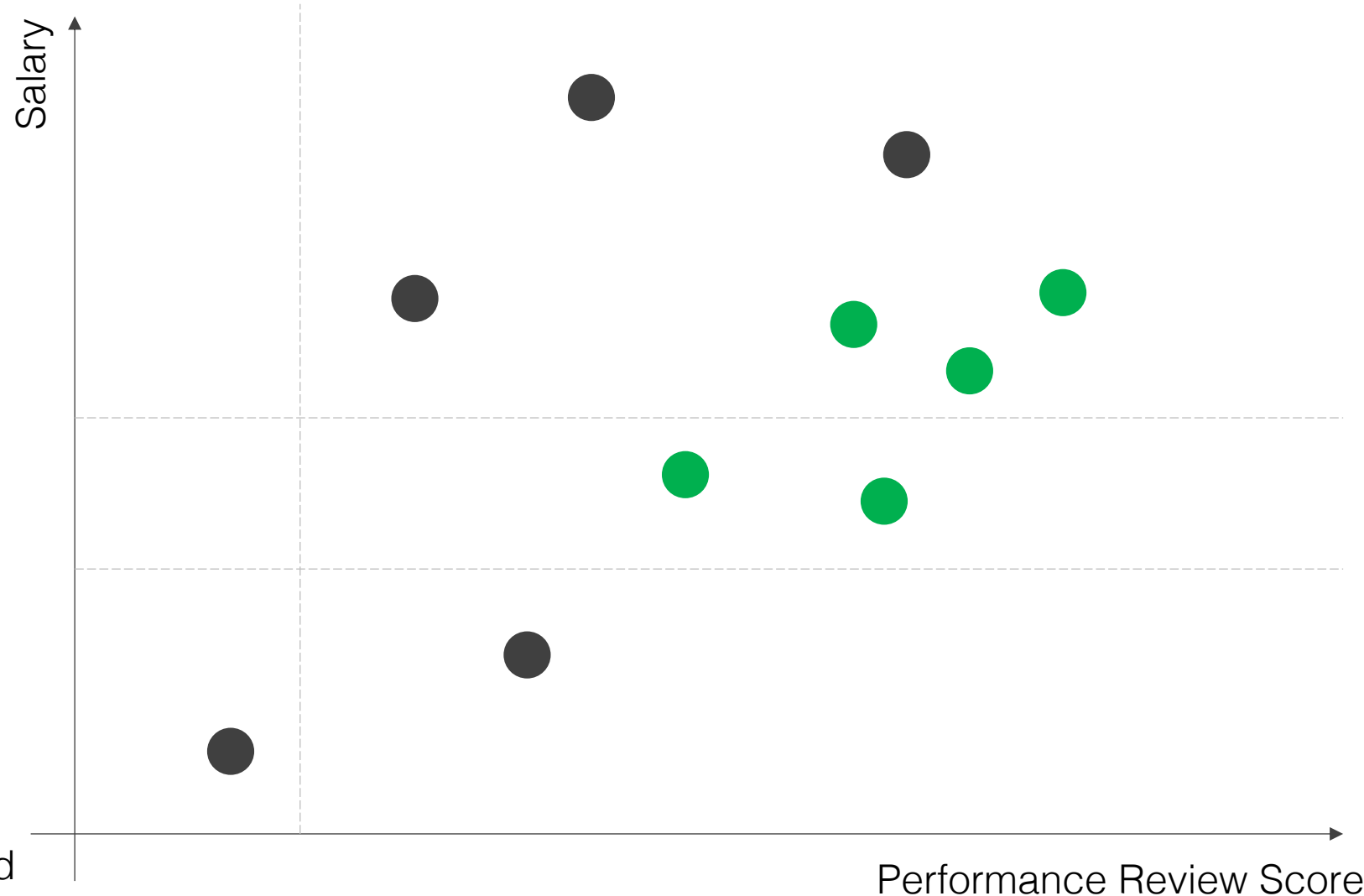
Classification and Regression Trees (CART)

Predicting promotions of salaried employees

1

Find the best “split” in any one feature (that best classifies the data) that divides the region in two

● Promoted
● Not promoted



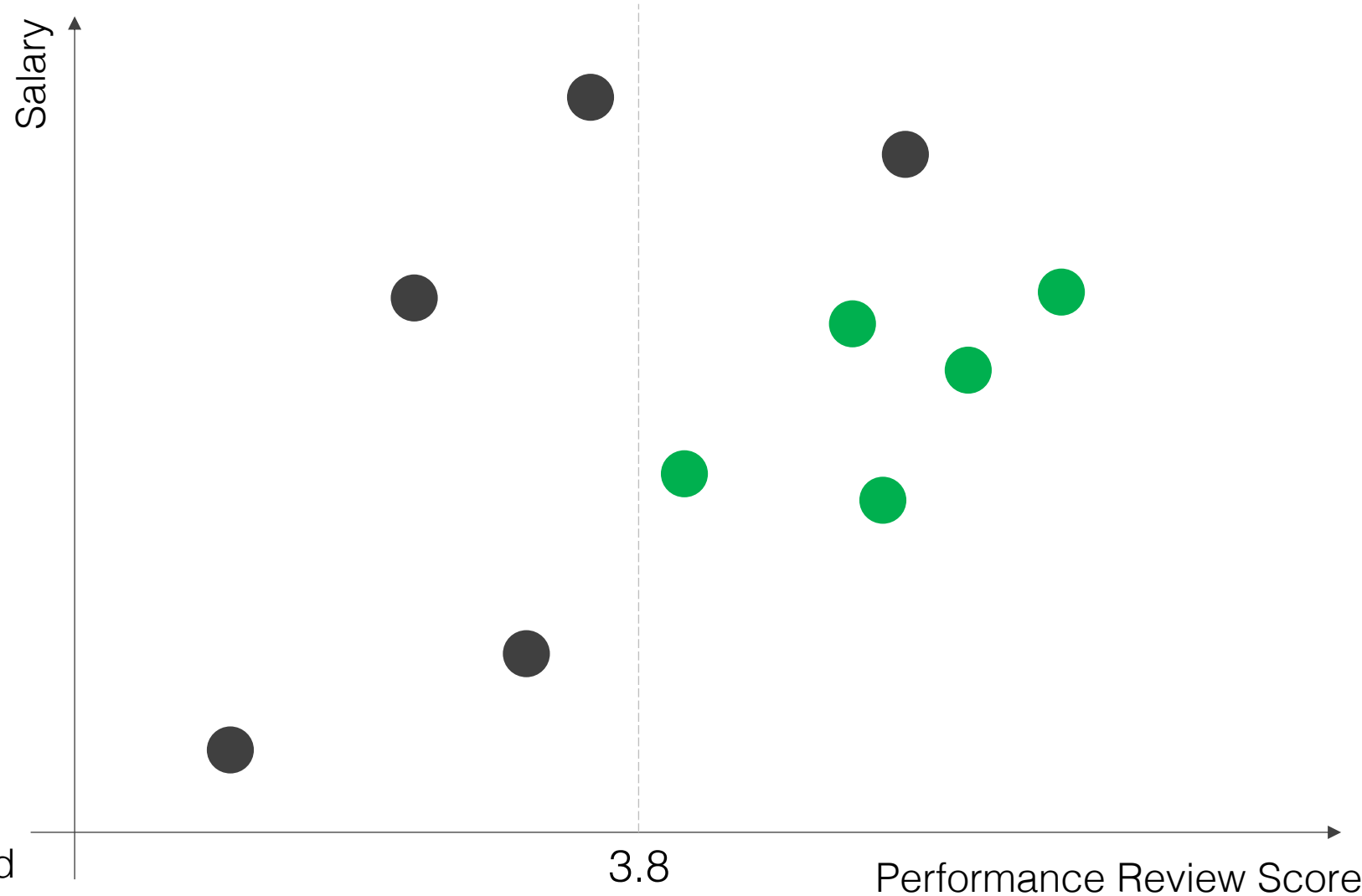
Classification and Regression Trees (CART)

Predicting promotions of salaried employees

1

Find the best “split” in any one feature (that best classifies the data) that divides the region in two

● Promoted
● Not promoted



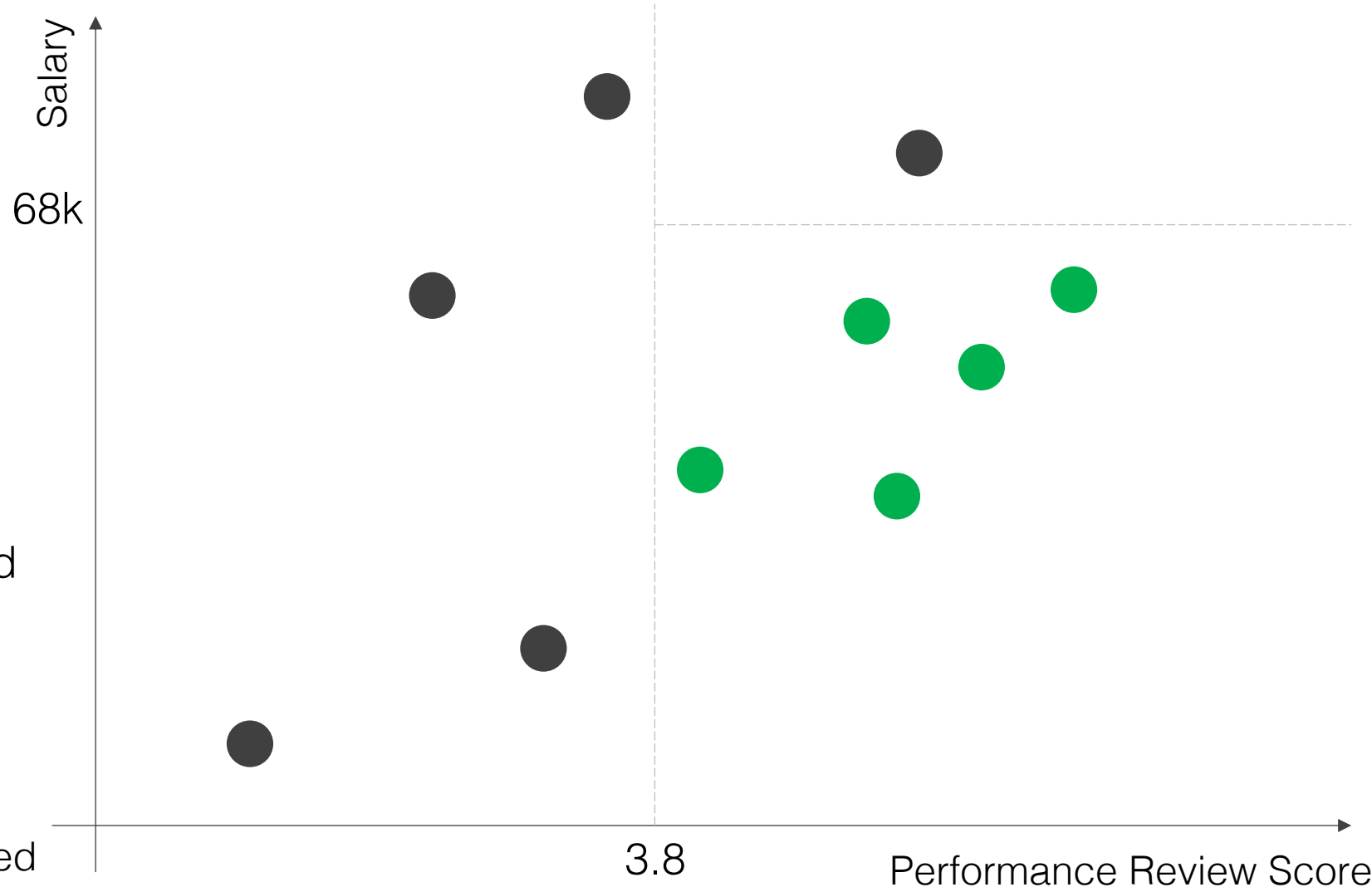
Classification and Regression Trees (CART)

Predicting promotions of salaried employees

- 1 Find the best “split” in any one feature (that best classifies the data) that divides the region in two
- 2 Continue splitting regions (1 feature at a time) until a stopping criterion is reached (e.g. there are at most N samples in any region)

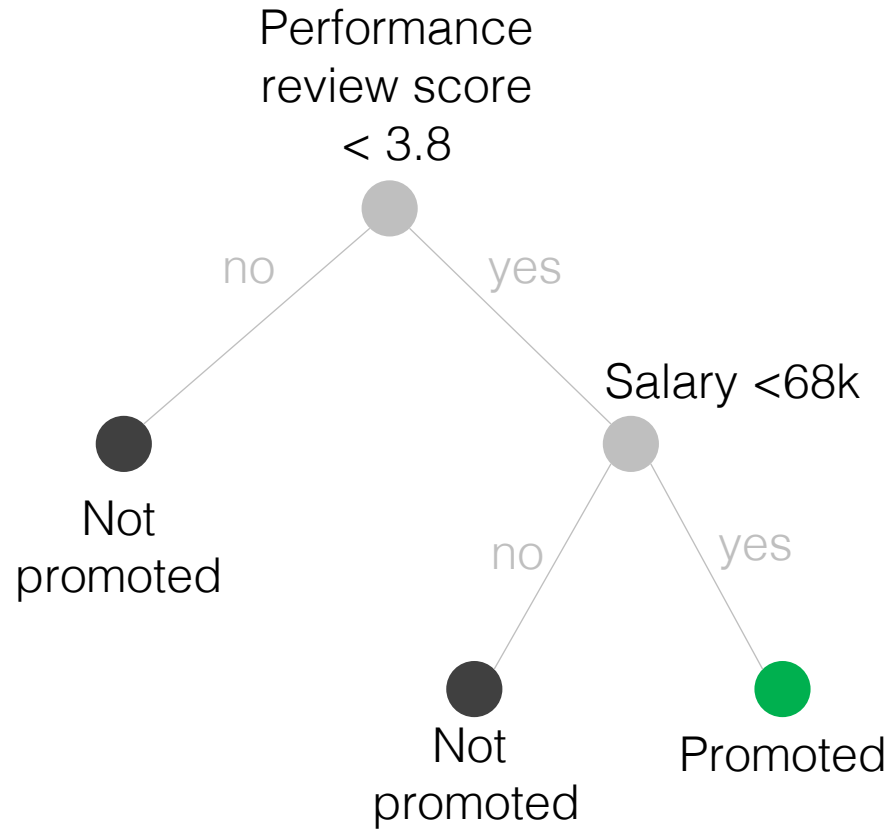
**Greedy, recursive
binary tree**

● Promoted
● Not promoted



Classification and Regression Trees (CART)

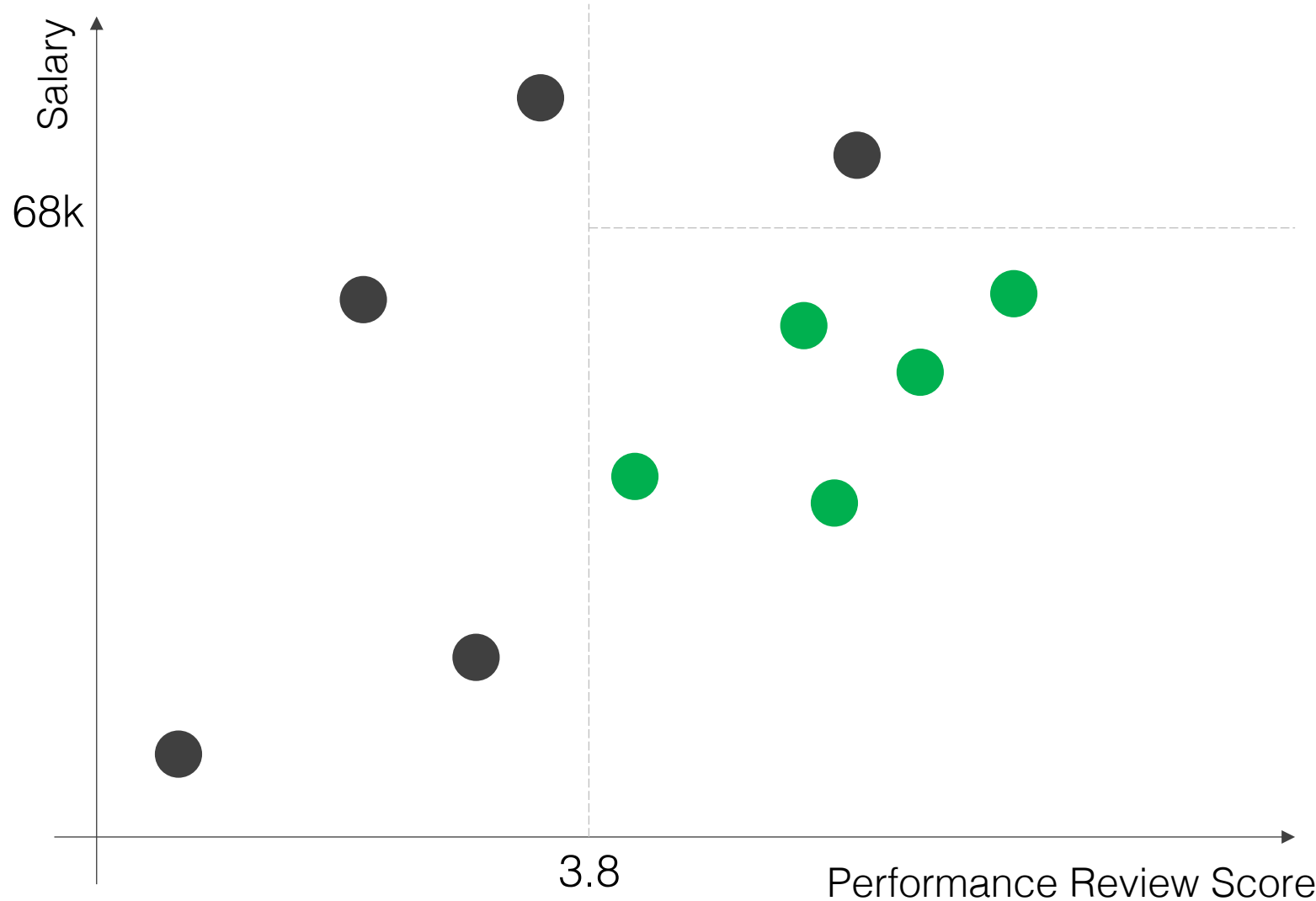
Tree representation:



● Splitting point

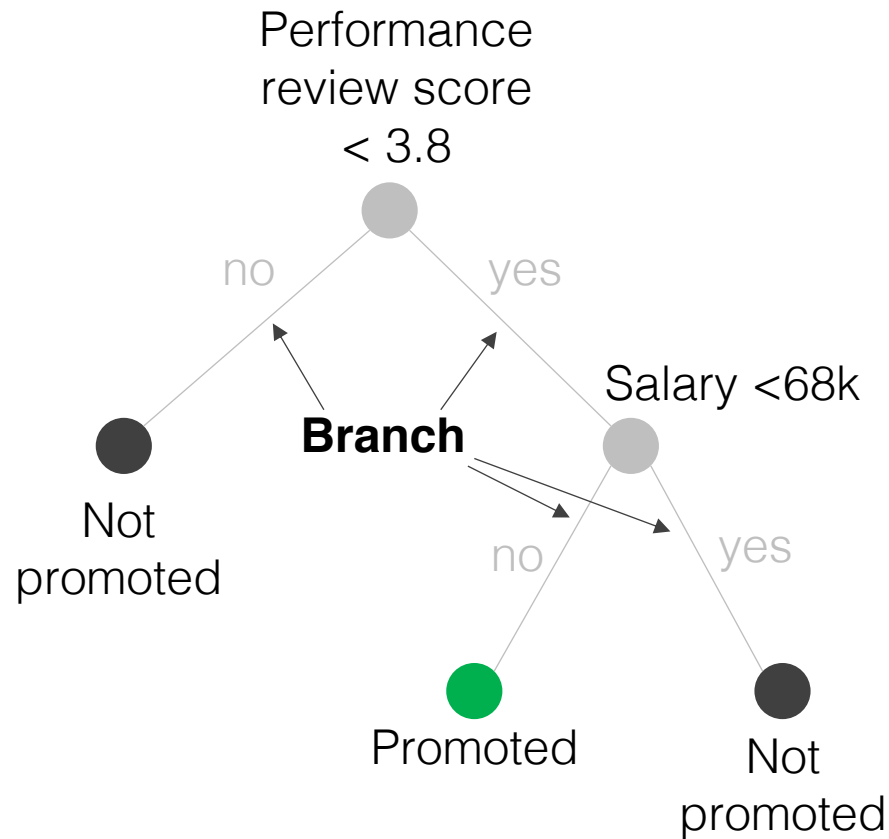
● Promoted

● Not promoted



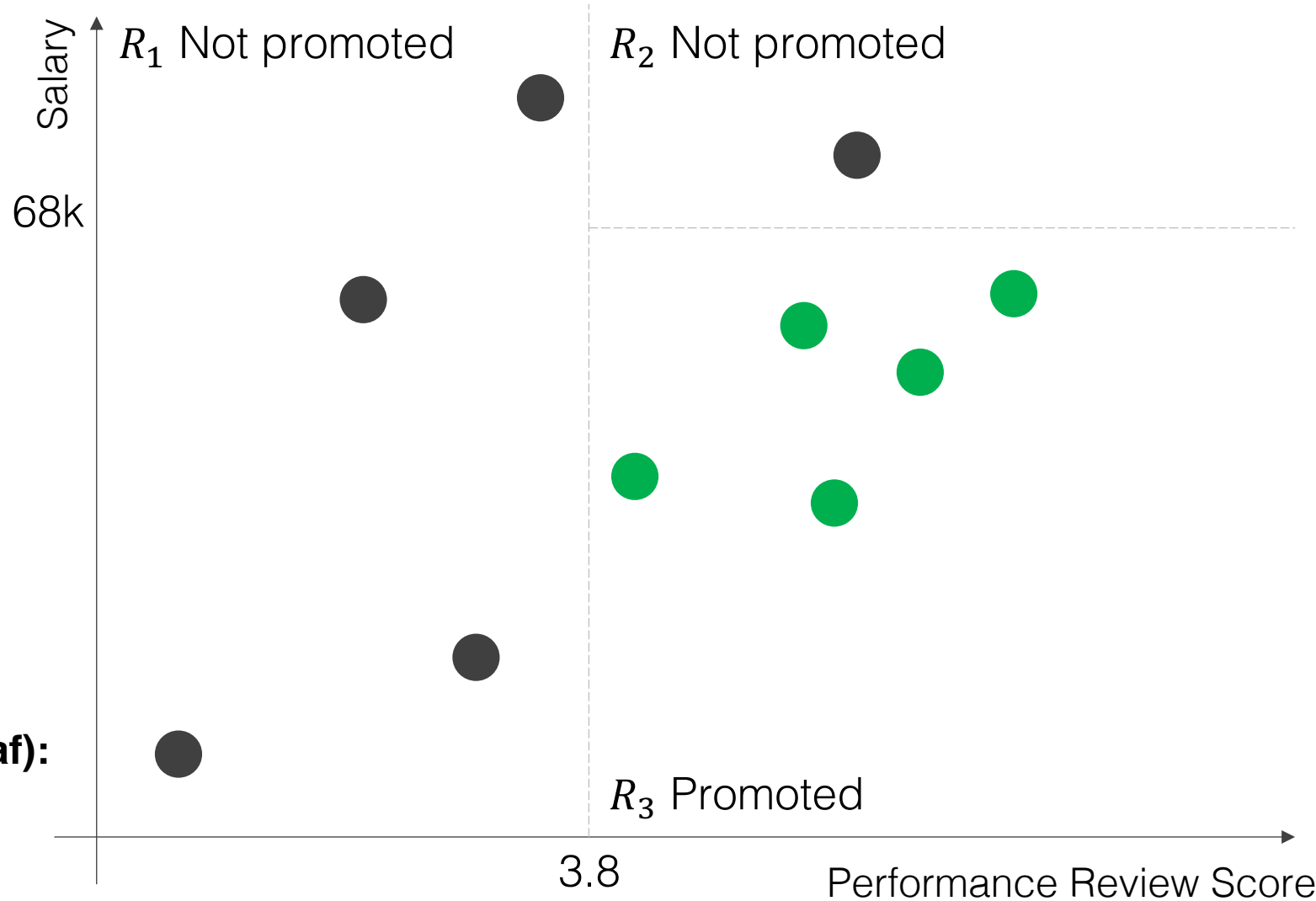
Classification and Regression Trees (CART)

Tree representation:



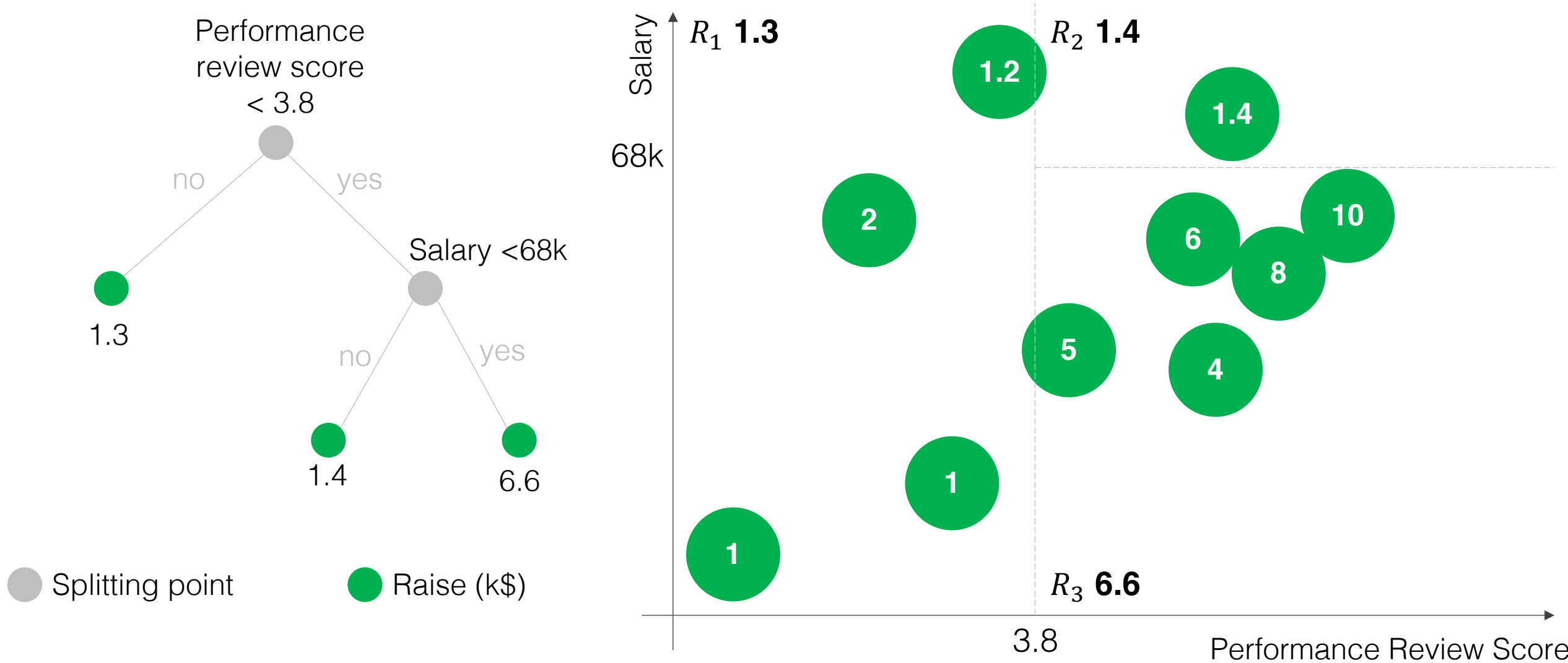
Internal node:
● Splitting point

Terminal node (leaf):
● Promoted
● Not promoted



The Regression Setting

In this case, each region is represented by an average of the values it contains



How do we determine which split to make?

Pick the split that reduces the error/cost most after the split

Regression

Mean square error

$$C_{MSE} = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

y_i = training data response i

\hat{y}_{R_j} = mean value in region j

Classification

Misclassification rate

$$C_{Misclass} = 1 - \max_k (\hat{p}_{jk})$$

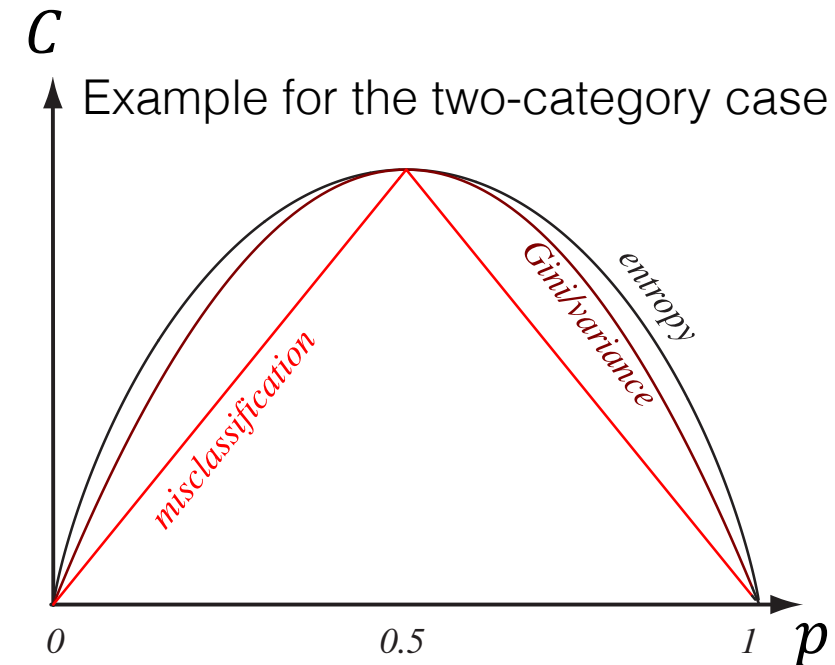
Gini impurity

$$C_{Gini} = \sum_{k=1}^K \hat{p}_{jk} (1 - \hat{p}_{jk})$$

Cross-entropy

$$C_{entropy} = - \sum_{k=1}^K \hat{p}_{jk} \log \hat{p}_{jk}$$

\hat{p}_{jk} = proportion of training observations in the j^{th} region from the k^{th} class



Duda, Hart, and Stork., Pattern Classification

Tree Pruning

Trees have the tendency to overfit the data

Consider the stopping rule: stop splitting once there is only 1 observation in each region (leads to complete overfit)

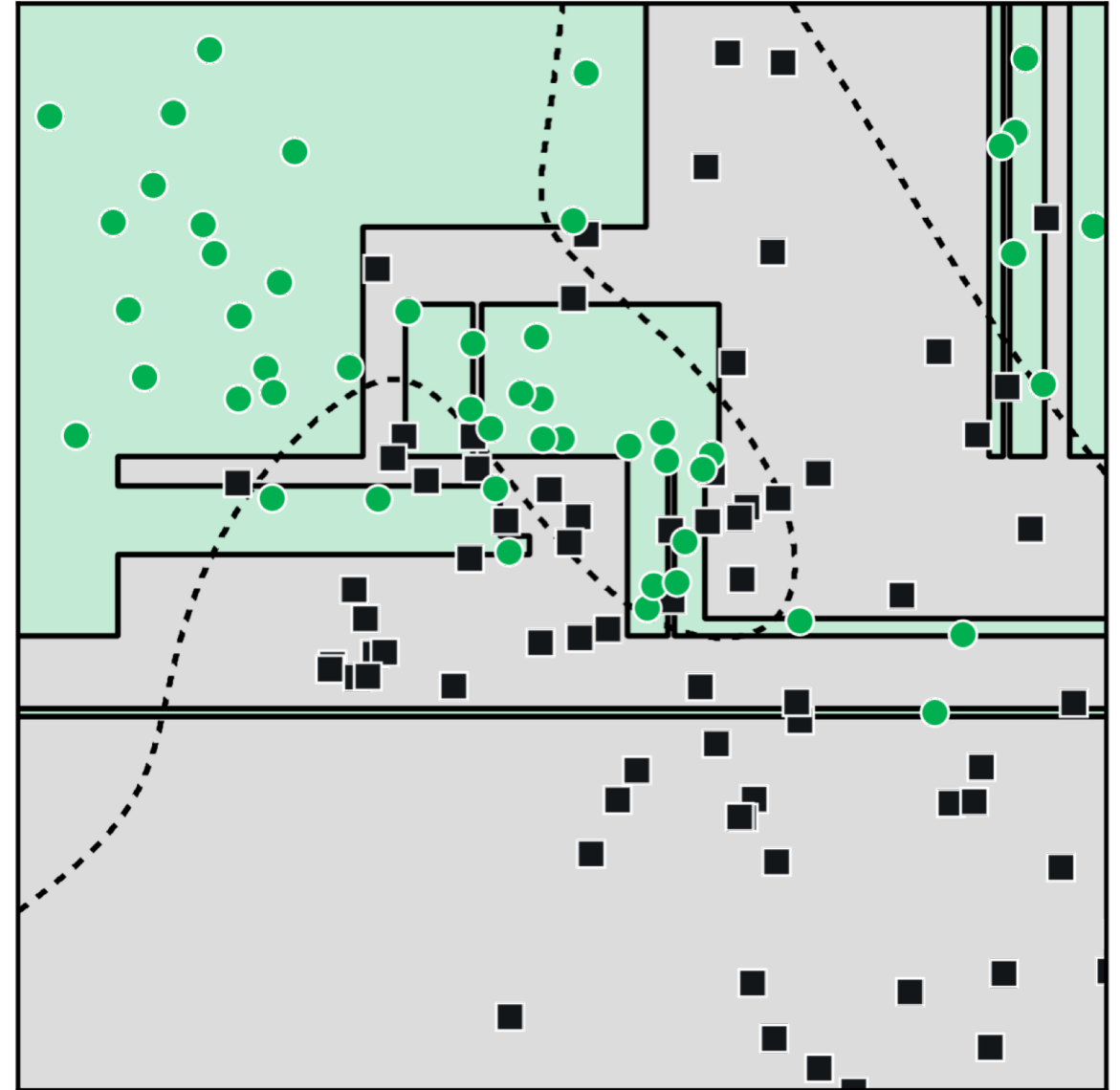
Pruning the tree back reduces this overfit (removing splits after the tree is formed)

Pruning can be optimized through a penalty on the number of terminal nodes:

$$C_{Prune} = \sum_{j=1}^T \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2 + \alpha T$$

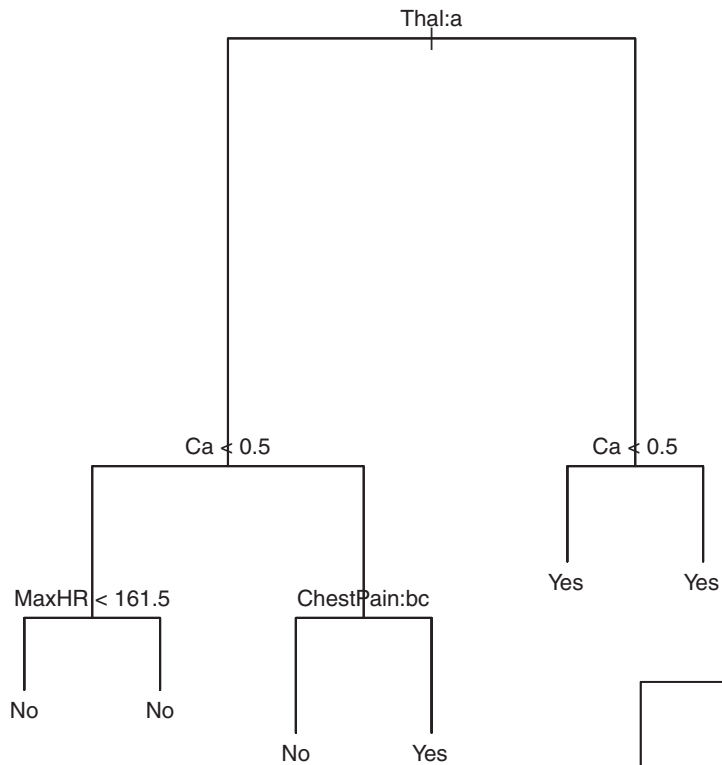
penalty on number of terminal nodes number of terminal nodes

Decision Tree



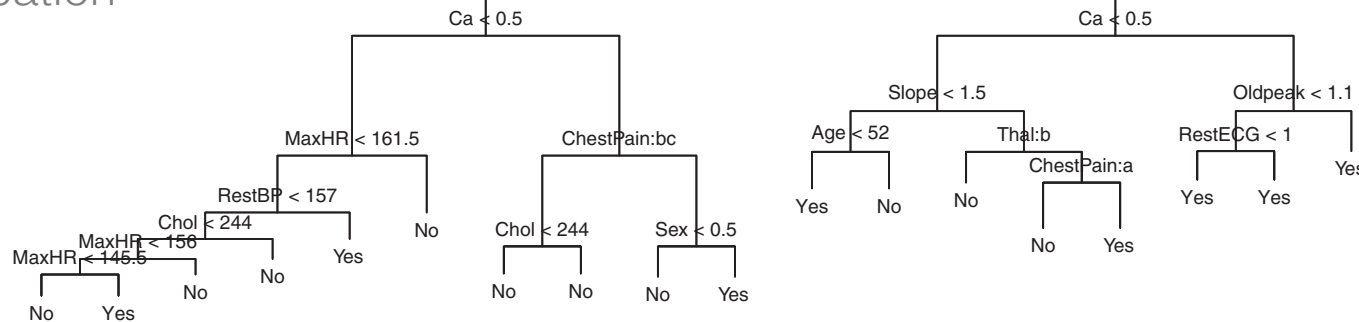
Pruning example

Pruned Tree

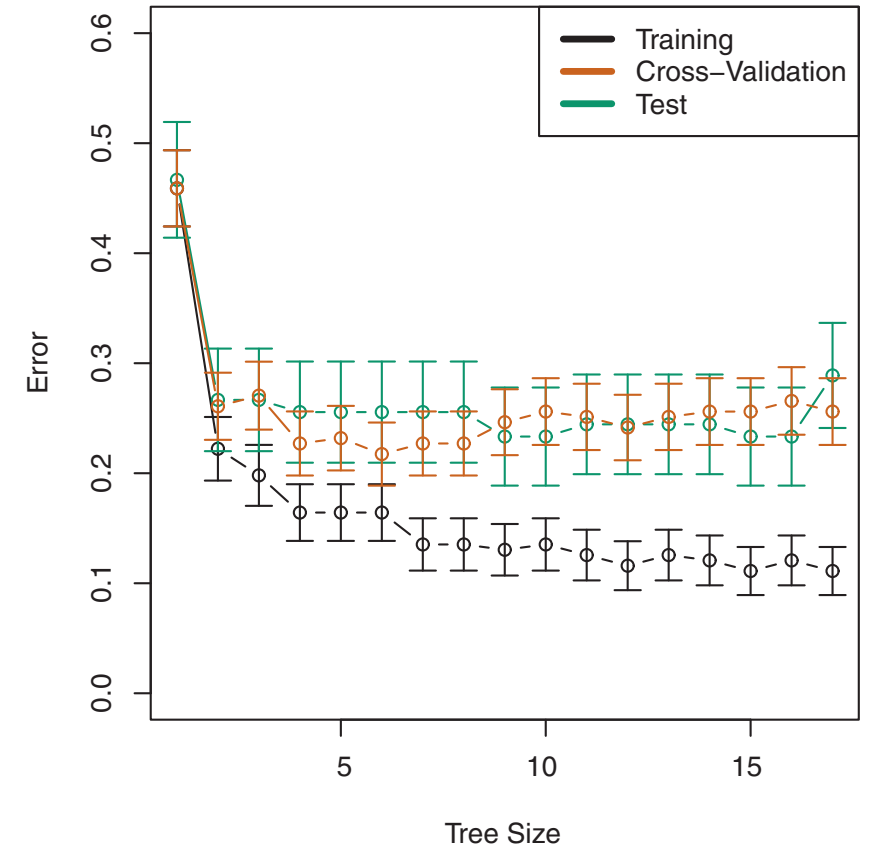


Original Tree

Example: heart disease classification

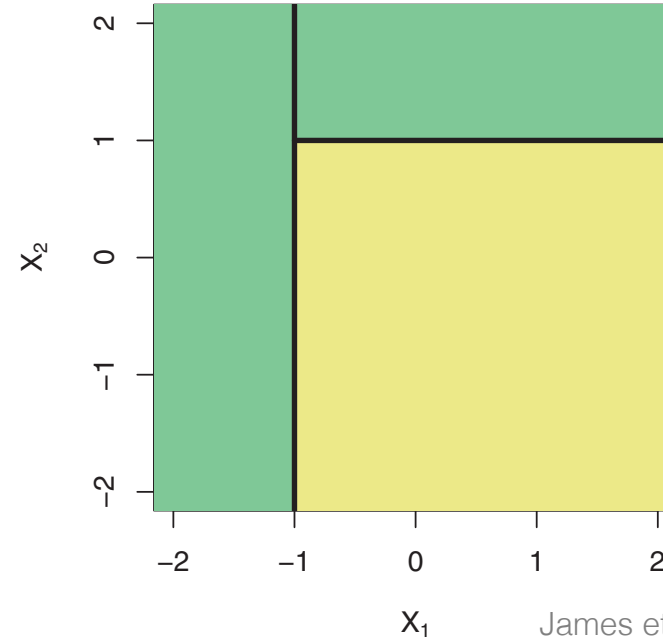
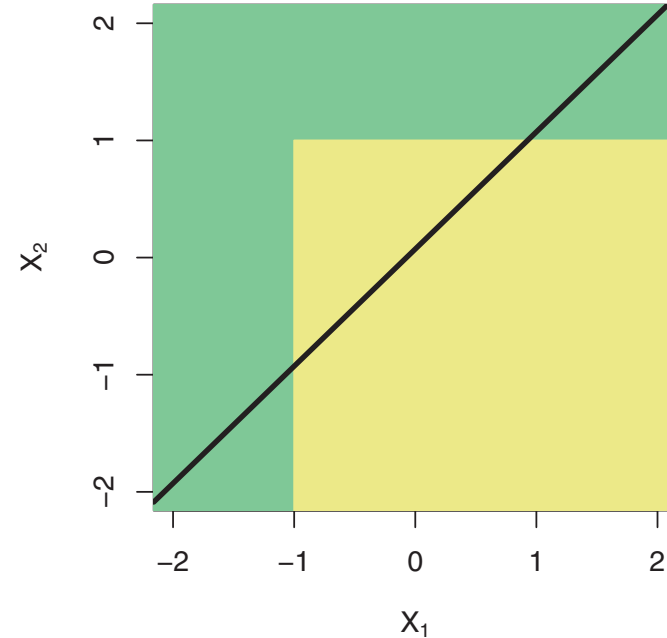
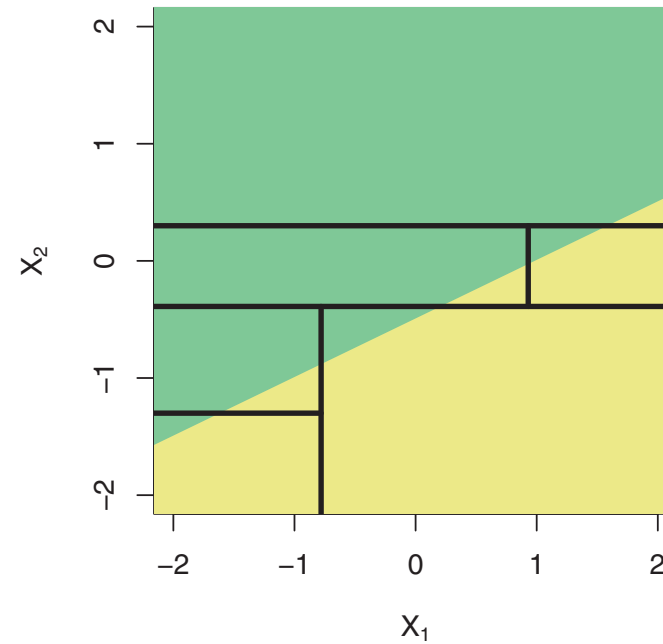
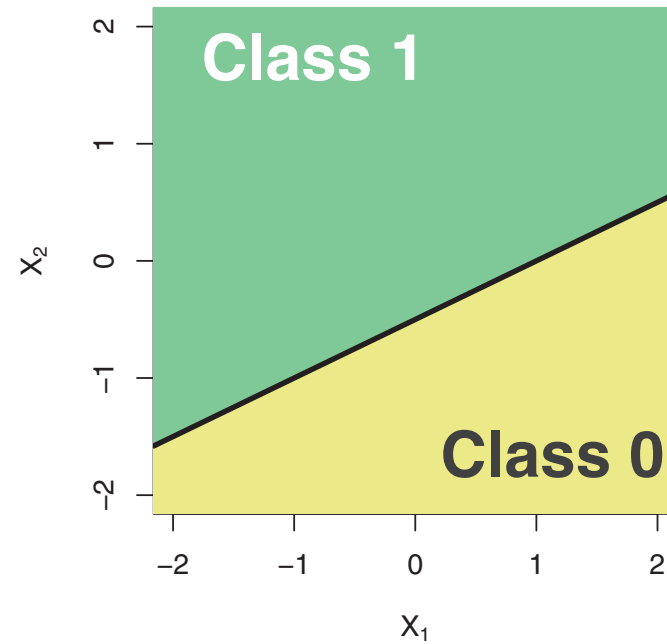


Performance



James et al., An Introduction to Statistical Learning

Linear model



Classification Tree

Struggle when the boundary is not parallel to an axis

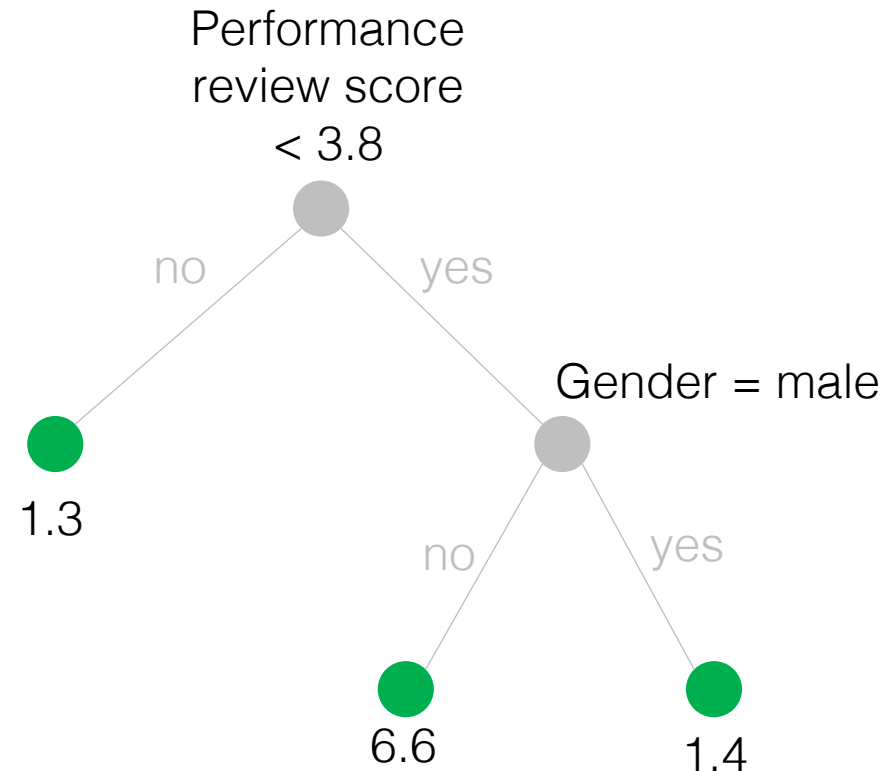
...nonlinear feature transforms could help...

James et al., An Introduction to Statistical Learning

Pros/Cons

Numerical data

Categorical data



Pros:

Trees easily handle multiple types of data

Trees are easy to interpret

Cons:

Trees do not typically have the same level of predictive accuracy of many methods

Tend to overfit
(have high variance)

Ensemble learning

How can we combine models to improve performance?

Bagging (bootstrap aggregation)

Random forests (tree-specific modification of bagging)

Gradient boosting

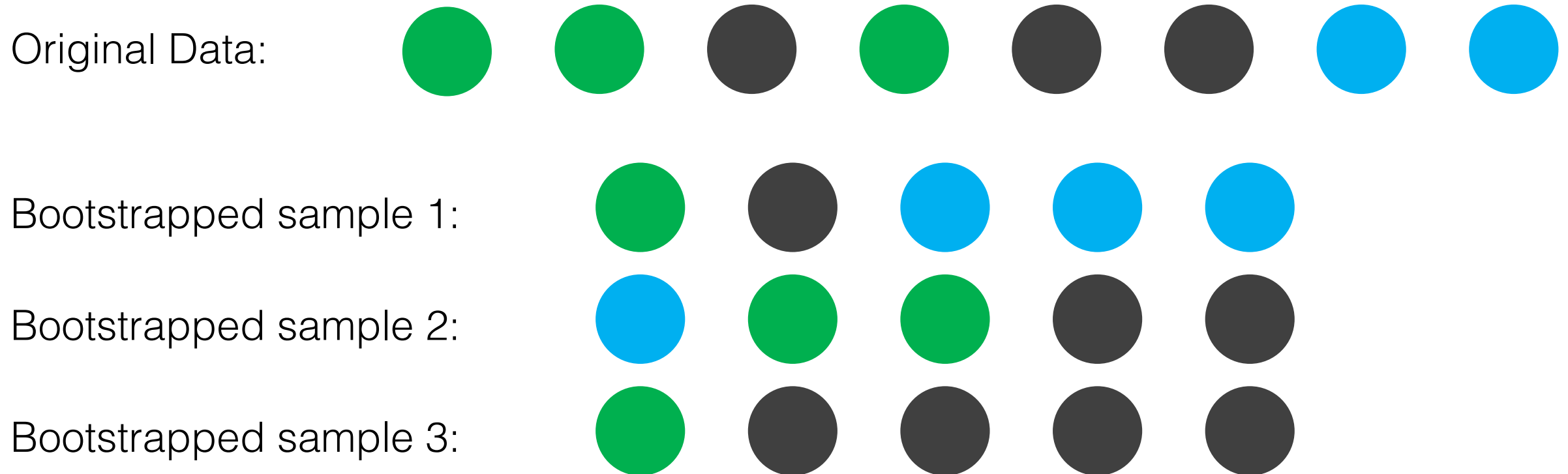
Applicable to
many models
(not just
trees)

Bagging

Bootstrap aggregation

Trees **overfit** (have high variance). Averaging over observations **reduces variance**

Recall bootstrap sampling (sampling with replacement):



Bagging

Bootstrap aggregation

Can be applied to many
machine learning
techniques!

- 1 Create a random bootstrap sample from the training data
- 2 Train a model on that bootstrap sample and call it $\hat{f}_i(\mathbf{x})$
- 3 Repeat 1 and 2 until we have B models trained on different bootstrap samples
- 4 Take the average of the output for our new model estimate:

$$\hat{f}_{bag}(\mathbf{x}) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i(\mathbf{x})$$

(for classification models we can take a majority vote instead)

Bagging

Tree Number:

1

2

3

4

Observations
Included:
(out of 1-9)

[1,2,3,3,8]

[1,2,4,7,7]

[1,5,6,8,9]

[2,2,2,4,9]

Features list:

[A, B, C, D]

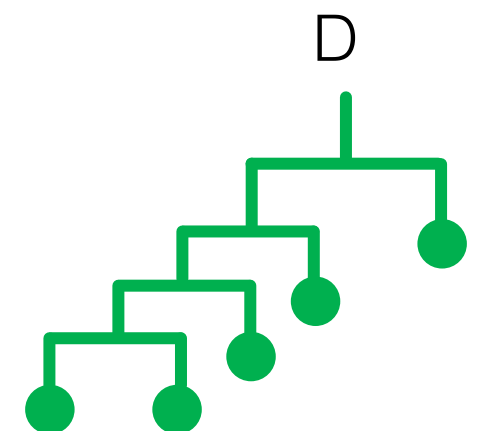
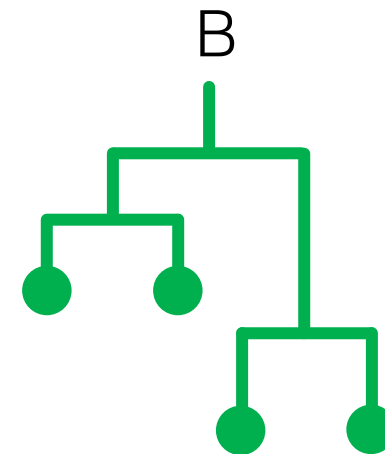
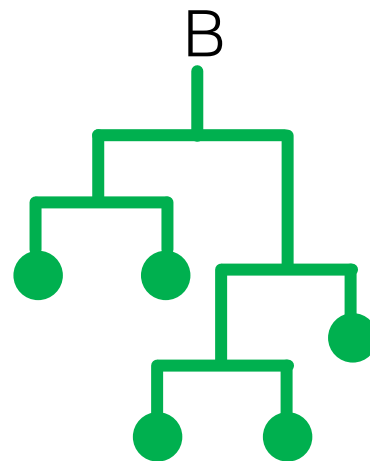
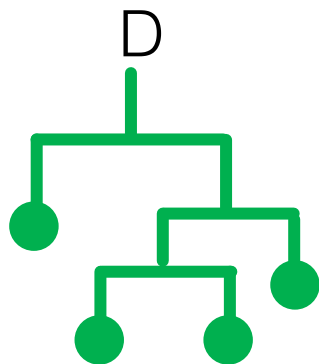
[A, B, C, D]

[A, B, C, D]

[A, B, C, D]

First split:

Trees:



Variable Importance

Decision trees are very interpretable, but this is lost with bagging

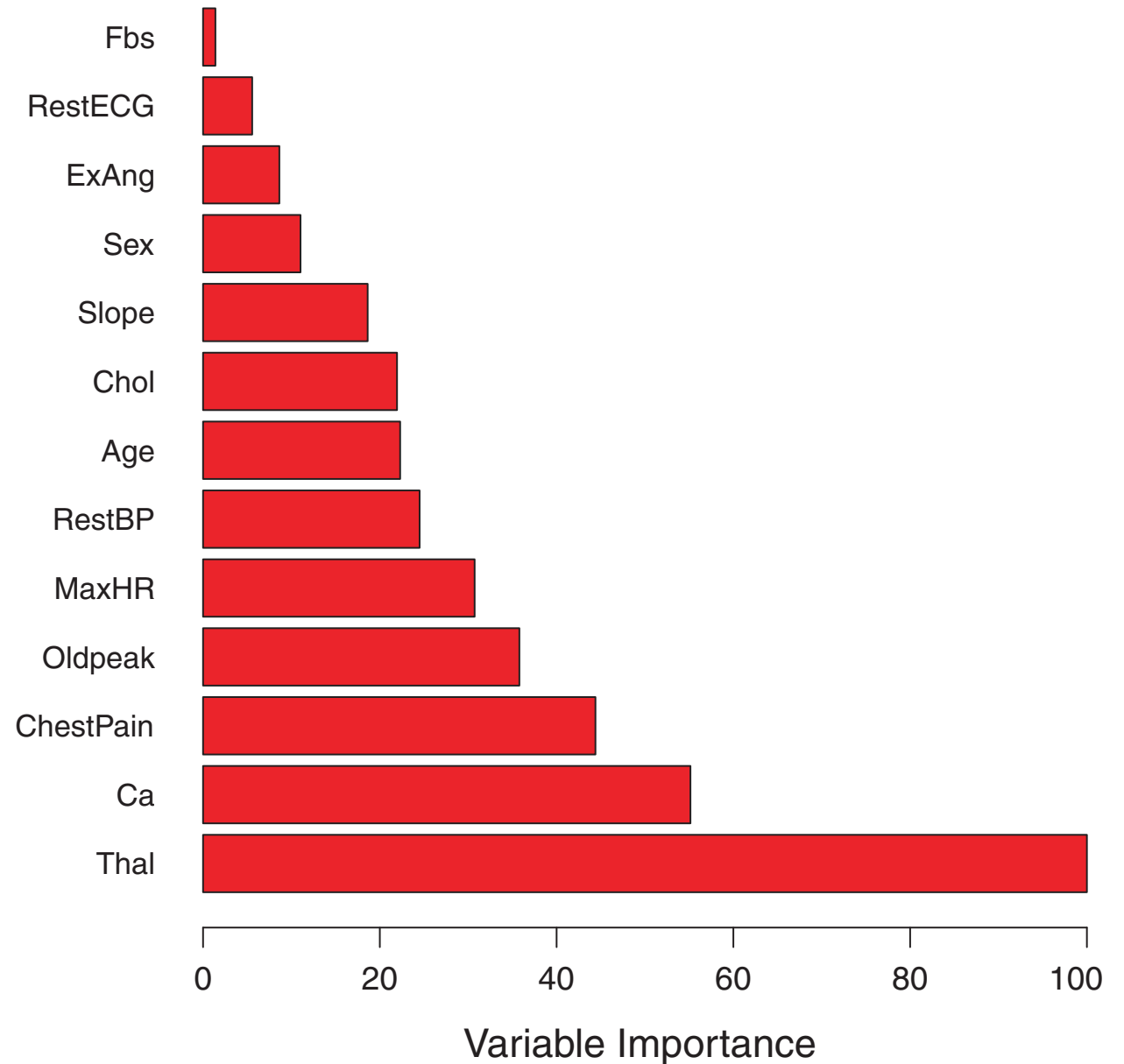
We can construct another measure called “variable importance” to **compare feature contributions**

1

Calculate the total amount the error (or impurity) decreased by splitting on each feature.

2

Average over all the trees resulting from bagging



Random Forests

A **small tweak on bagging**

Random forests
decorrelate
the bagged trees

Decision trees are constructed greedily

This can lead to highly correlated trees

“Strong” features will typically be split before moderately strong predictors.

Each time a split is considered, a **random subset of m features** is selected as candidates from the full set of p features

Typically chose: $m = \sqrt{p}$

Bagging

Random forests

Observations
Included:
(out of 1-9)

$[1, 2, 3, 3, 8]$

$[1, 2, 3, 3, 8]$

Feature list:

$[A, B, C, D]$

$[A, B, C, D]$

Feature options
for each split:

$[A, B, C, D]$

$[A, B, C, D]$

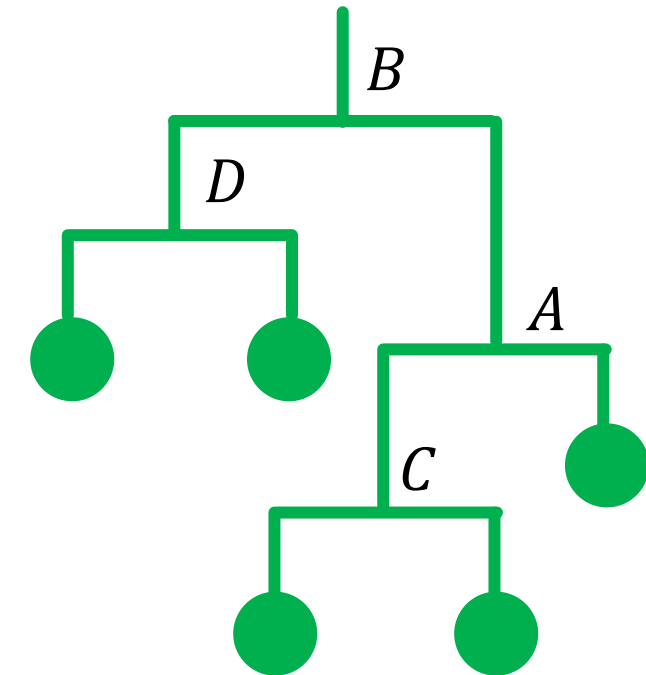
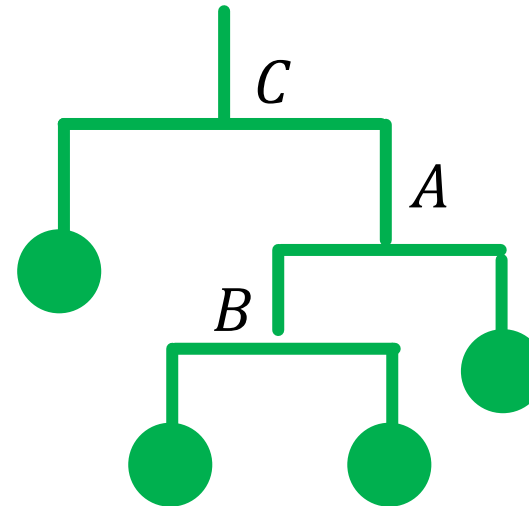
$[A, B, C, D]$

$[A, B]$

$[C, D]$

$[A, B]$

$[B, C]$

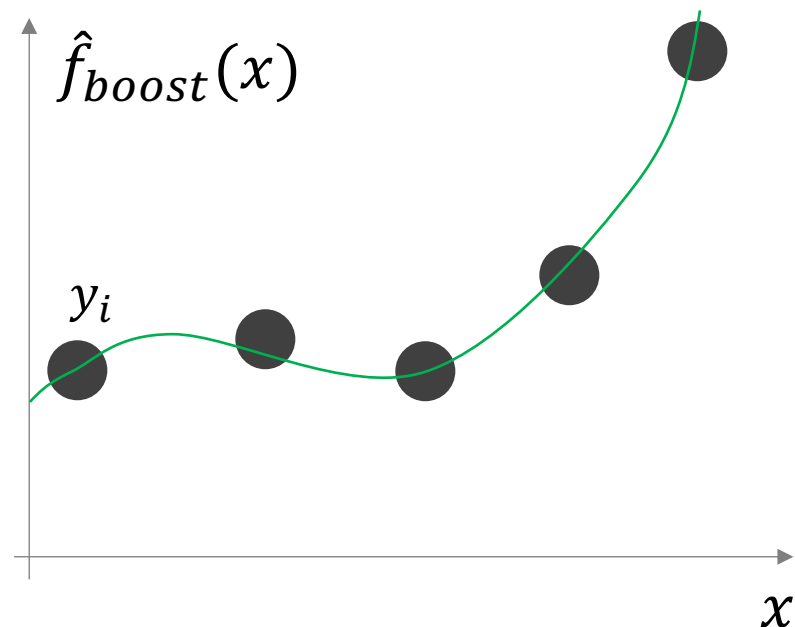
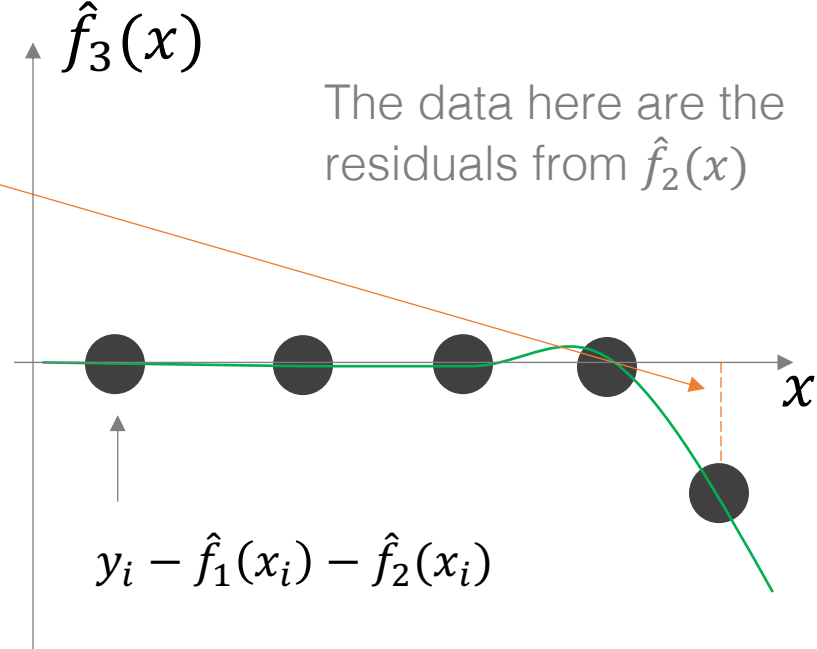
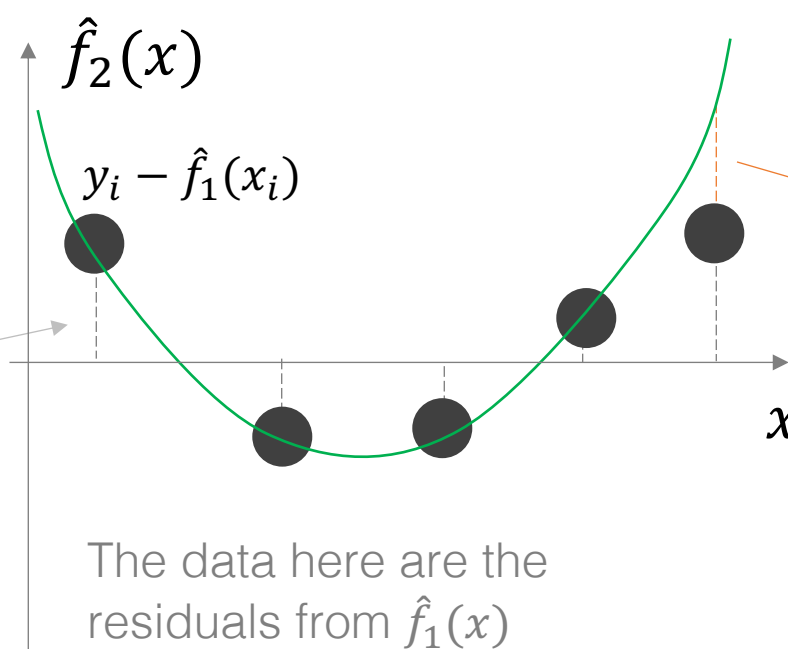
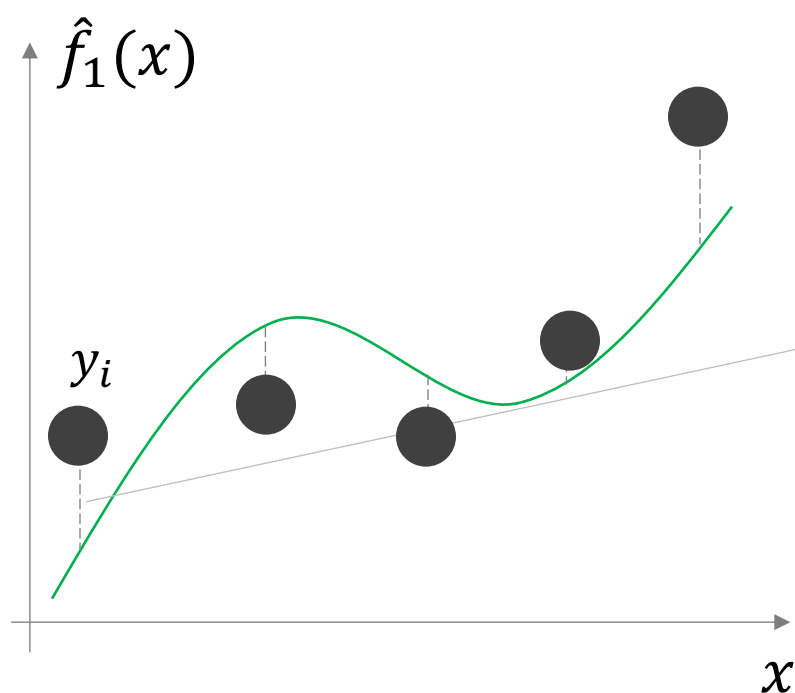


Boosting

Can be applied to many
machine learning
techniques!

Bagging created trees that were designed to be as independent as possible

Boosting involves building trees **sequentially**, each building on the errors of the last



We build consecutive models, each fit to the residuals of the last model

We sum models output to get the boosted prediction

$$\hat{f}_{boost}(x) = \hat{f}_1(x) + \hat{f}_2(x) + \hat{f}_3(x)$$

Boosting

Boosting for regression trees

- 1 Select the number of models to train, B , and learning rate λ
- 2 Set $\hat{f}(\mathbf{x}) = 0$ and $r_i = y_i$ for all the training data
- 3 Fit a tree, $\hat{f}_i(\mathbf{x})$ to the residuals, r_i (with d splits)
- 4 Update $\hat{f}(\mathbf{x}) = \hat{f}(\mathbf{x}) + \lambda \hat{f}_i(\mathbf{x})$
- 5 Update the residuals $r_i = r_i - \lambda \hat{f}_i(\mathbf{x}_i)$
- 6 Output the boosted model:
$$\hat{f}(\mathbf{x}) = \sum_{i=1}^B \lambda \hat{f}_i(\mathbf{x})$$

λ slows down the learning process to avoid overfitting

Often this is just a “stump” with $d = 1$ split

Repeat B times

Model Stacking

Train multiple supervised learning techniques (could be different models)

THEN Train a supervised learning technique that includes the **outputs** of the other models

Supervised Learning Techniques

Covered so far

K-Nearest Neighbors

Linear regression

Perceptron

Logistic Regression

Fisher's Linear Discriminant / Linear Discriminant Analysis

Quadratic Discriminant Analysis

Naïve Bayes

Decision Trees and Random Forests

Ensemble methods (bagging, boosting, stacking)

Neural Networks

Rely on a linear combination of weights and features: $\mathbf{w}^T \mathbf{x}$

Can also be used for regression

Can be used with many machine learning techniques (and regression)