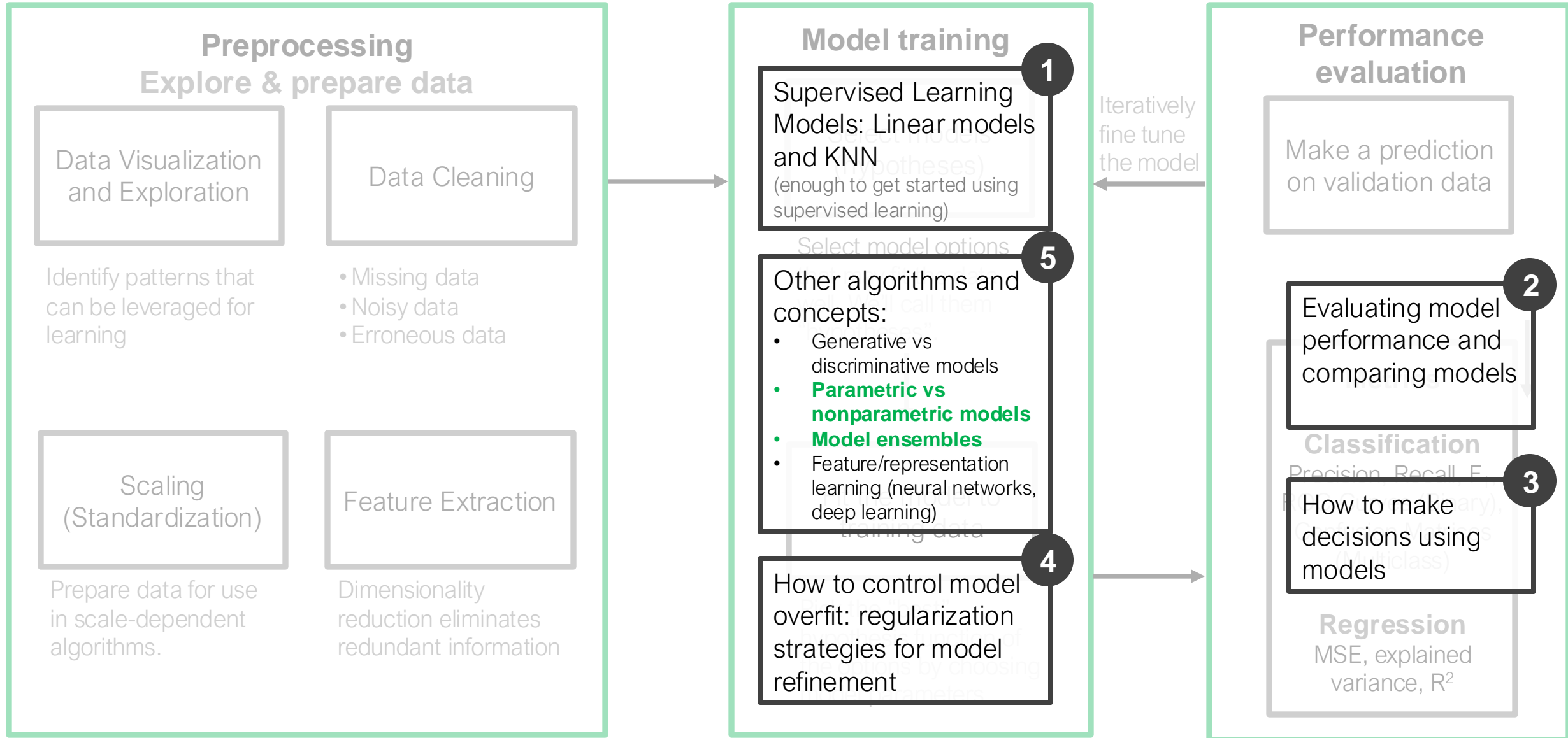


# Tree-based Models and Ensembles

# Supervised learning in practice



# Supervised Learning Techniques

Covered so far

Linear Regression

K-Nearest Neighbors

Perceptron

Logistic Regression

Linear Discriminant Analysis

Quadratic Discriminant Analysis

Naïve Bayes

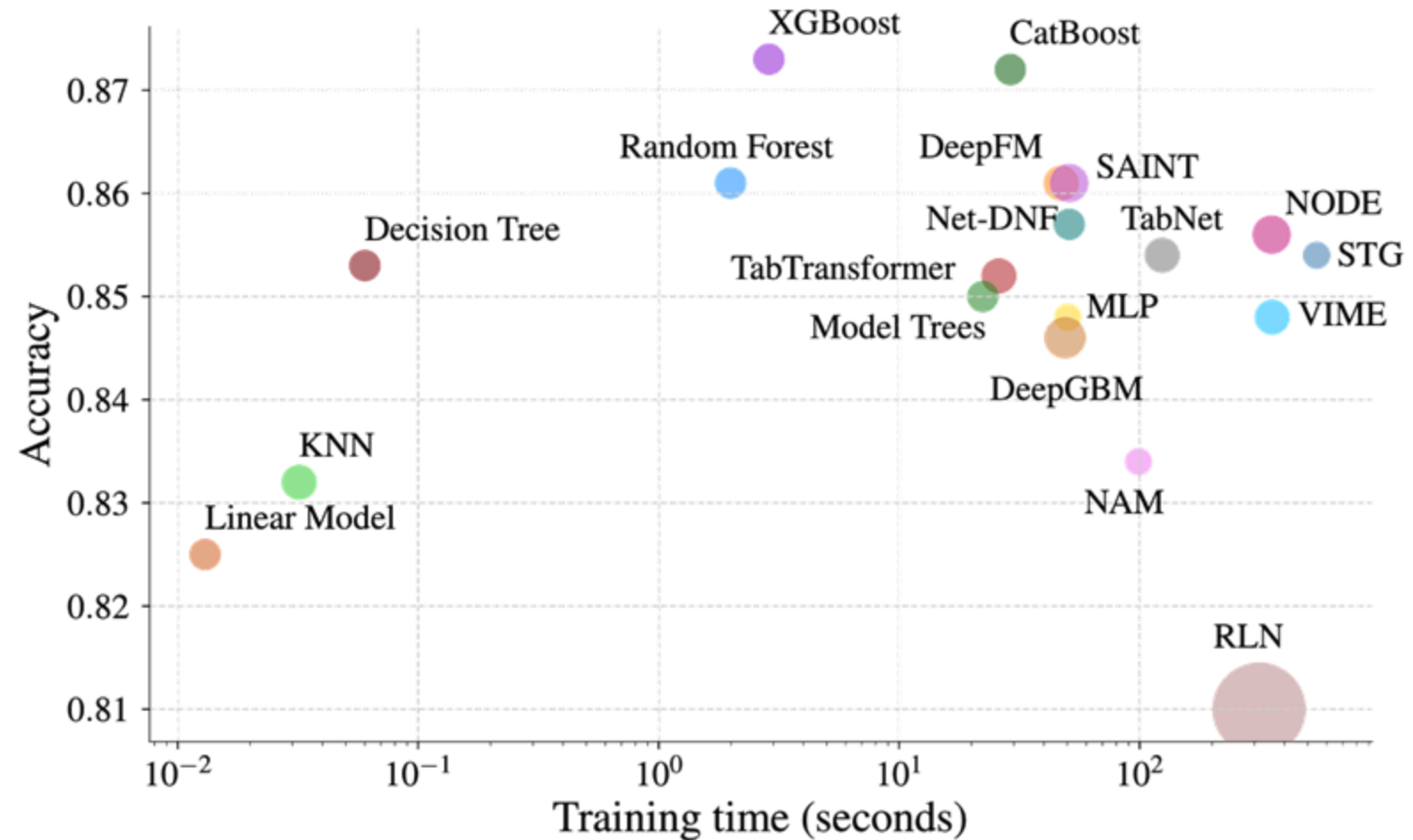
Decision Trees and Random Forests

Ensemble methods (bagging, boosting, stacking)

# Decision Tree Ensembles

"**gradient-boosted tree ensembles** still mostly outperform deep learning models on supervised learning tasks [on **heterogeneous tabular data**]"

a data set with a fixed number of features that are either continuous or categorical



Results on Adult Income dataset from UCI repository. Task: predict whether income exceeds \$50K/yr based on census data

Gradient-boosted decision tree ensembles	Benchmark datasets: HELOC <b>bold ones are the best ones by comparison</b> Adult HIGGS Coverttype Cal. Housing									
		Acc ↑	AUC ↑	Acc ↑	AUC ↑	Acc ↑	AUC ↑	Acc ↑	AUC ↑	MSE ↓
	Linear Model	73.0±0.0	80.1±0.1	82.5±0.2	85.4±0.2	64.1±0.0	68.4±0.0	72.4±0.0	92.8±0.0	0.528±0.008
	KNN [65]	72.2±0.0	79.0±0.1	83.2±0.2	87.5±0.2	62.3±0.1	67.1±0.0	70.2±0.1	90.1±0.2	0.421±0.009
	Decision Tree [197]	80.3±0.0	89.3±0.1	85.3±0.2	89.8±0.1	71.3±0.0	78.7±0.0	79.1±0.0	95.0±0.0	0.404±0.007
	Random Forest [198]	82.1±0.2	90.0±0.2	86.1±0.2	91.7±0.2	71.9±0.0	79.7±0.0	78.1±0.1	96.1±0.0	0.272±0.006
	XGBoost [53]	<u>83.5±0.2</u>	92.2±0.0	<u>87.3±0.2</u>	<u>92.8±0.1</u>	<u>77.6±0.0</u>	<u>85.9±0.0</u>	<b>97.3±0.0</b>	<b>99.9±0.0</b>	0.206±0.005
	LightGBM [78]	<u>83.5±0.1</u>	<u>92.3±0.0</u>	<b>87.4±0.2</b>	<b>92.9±0.1</b>	77.1±0.0	85.5±0.0	93.5±0.0	99.7±0.0	<b>0.195±0.005</b>
	CatBoost [79]	<b>83.6±0.3</b>	<b>92.4±0.1</b>	87.2±0.2	<u>92.8±0.1</u>	77.5±0.0	85.8±0.0	<u>96.4±0.0</u>	<u>99.8±0.0</u>	<u>0.196±0.004</u>
	Model Trees [199]	82.6±0.2	91.5±0.0	85.0±0.2	90.4±0.1	69.8±0.0	76.7±0.0	-	-	0.385±0.019
Deep Learning	MLP [200]	73.2±0.3	80.3±0.1	84.8±0.1	90.3±0.2	77.1±0.0	85.6±0.0	91.0±0.4	76.1±3.0	0.263±0.008
	DeepFM [15]	73.6±0.2	80.4±0.1	86.1±0.2	91.7±0.1	76.9±0.0	83.4±0.0	-	-	0.260±0.006
	DeepGBM [70]	78.0±0.4	84.1±0.1	84.6±0.3	90.8±0.1	74.5±0.0	83.0±0.0	-	-	0.856±0.065
	RLN [72]	73.2±0.4	80.1±0.4	81.0±1.6	75.9±8.2	71.8±0.2	79.4±0.2	77.2±1.5	92.0±0.9	0.348±0.013
	TabNet [5]	81.0±0.1	90.0±0.1	85.4±0.2	91.1±0.1	76.5±1.3	84.9±1.4	93.1±0.2	99.4±0.0	0.346±0.007
	VIME [88]	72.7±0.0	79.2±0.0	84.8±0.2	90.5±0.2	76.9±0.2	85.5±0.1	90.9±0.1	82.9±0.7	0.275±0.007
	TabTransformer [98]	73.3±0.1	80.1±0.2	85.2±0.2	90.6±0.2	73.8±0.0	81.9±0.0	76.5±0.3	72.9±2.3	0.451±0.014
	NODE [6]	79.8±0.2	87.5±0.2	85.6±0.3	91.1±0.2	76.9±0.1	85.4±0.1	89.9±0.1	98.7±0.0	0.276±0.005
	Net-DNF [57]	82.6±0.4	91.5±0.2	85.7±0.2	91.3±0.1	76.6±0.1	85.1±0.1	94.2±0.1	99.1±0.0	-
	STG [201]	73.1±0.1	80.0±0.1	85.4±0.1	90.9±0.1	73.9±0.1	81.9±0.1	81.8±0.3	96.2±0.0	0.285±0.006
	NAM [202]	73.3±0.1	80.7±0.3	83.4±0.1	86.6±0.1	53.9±0.6	55.0±1.2	-	-	0.725±0.022
	SAINT [9]	82.1±0.3	90.7±0.2	86.1±0.3	91.6±0.2	<b>79.8±0.0</b>	<b>88.3±0.0</b>	96.3±0.1	<u>99.8±0.0</u>	0.226±0.004

# Parametric vs Nonparametric techniques

## Non-parametric Models

Complexity of the model grows with the size of the training data

How is KNN more complex based on data size?

- K-Nearest Neighbors
- Decision Trees

## Parametric Models

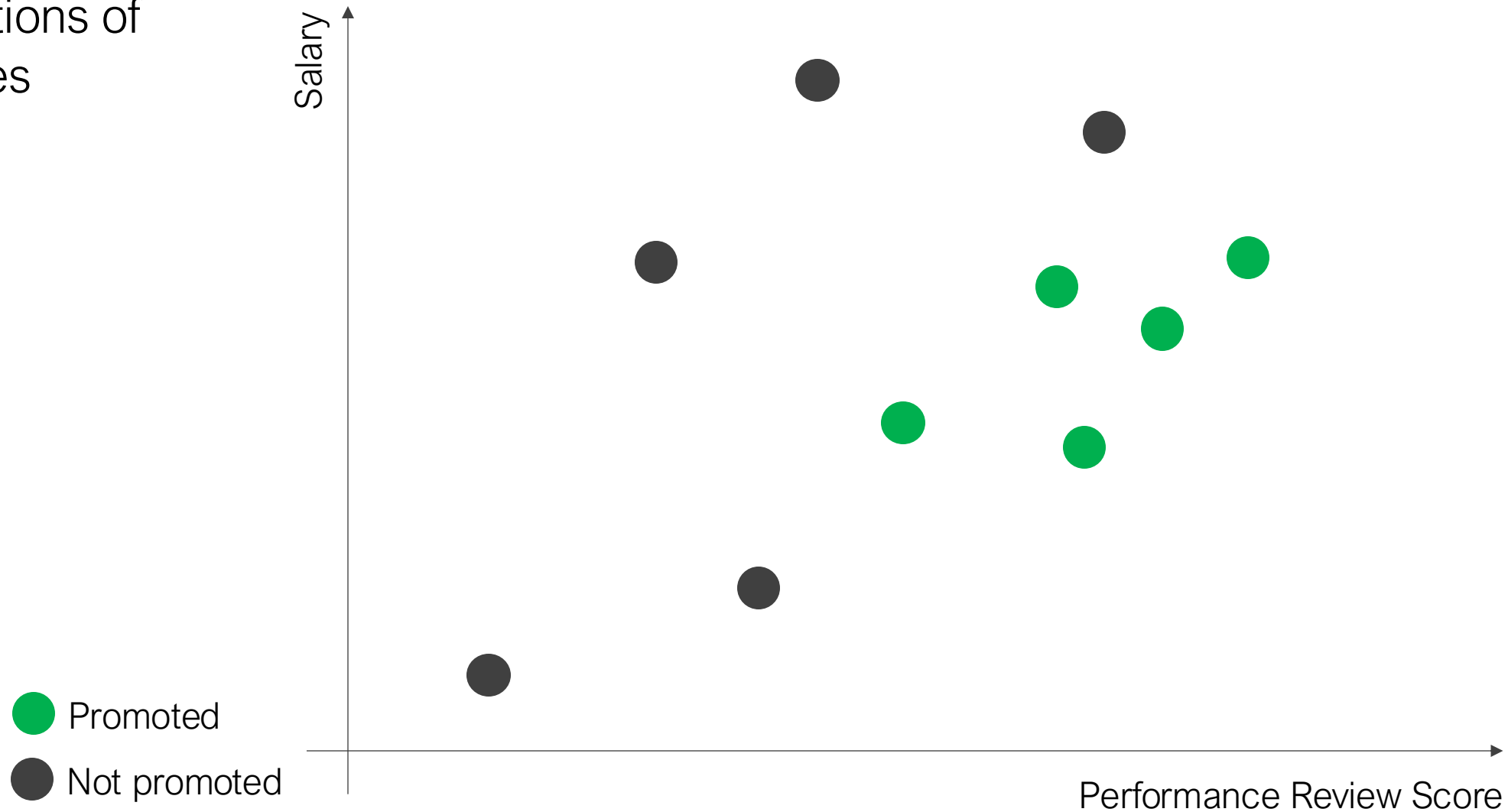
Fixed number of parameters (i.e. a fixed structure)

- Linear regression
- Logistic regression
- LDA, QDA
- Naïve Bayes with Gaussian likelihoods

# Classification and Regression Trees (CART)

Classification trees = decision trees

Predicting promotions of  
salaried employees



# Classification and Regression Trees (CART)

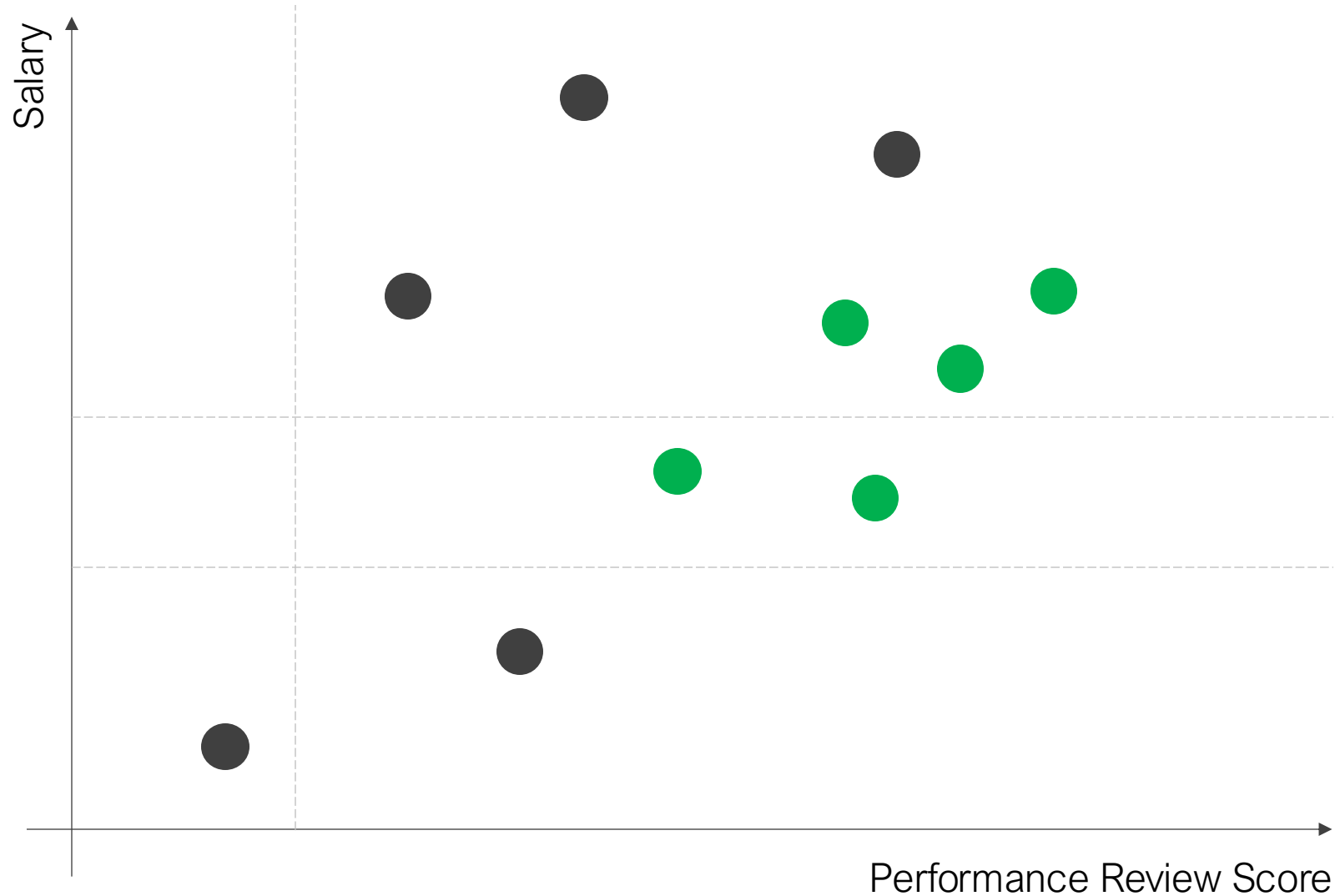
Predicting promotions of  
salaried employees

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)

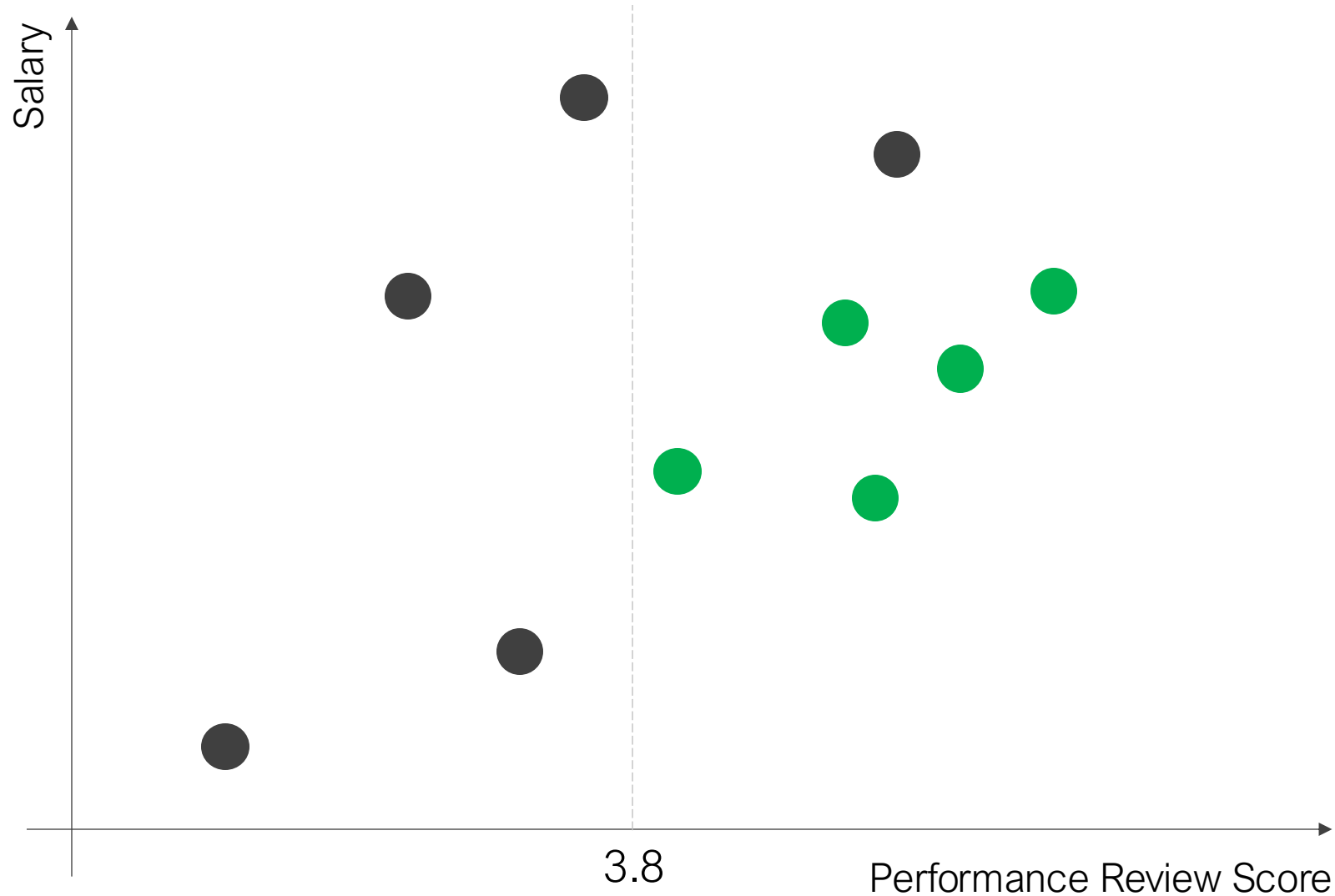
Predicting promotions of  
salaried employees

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)

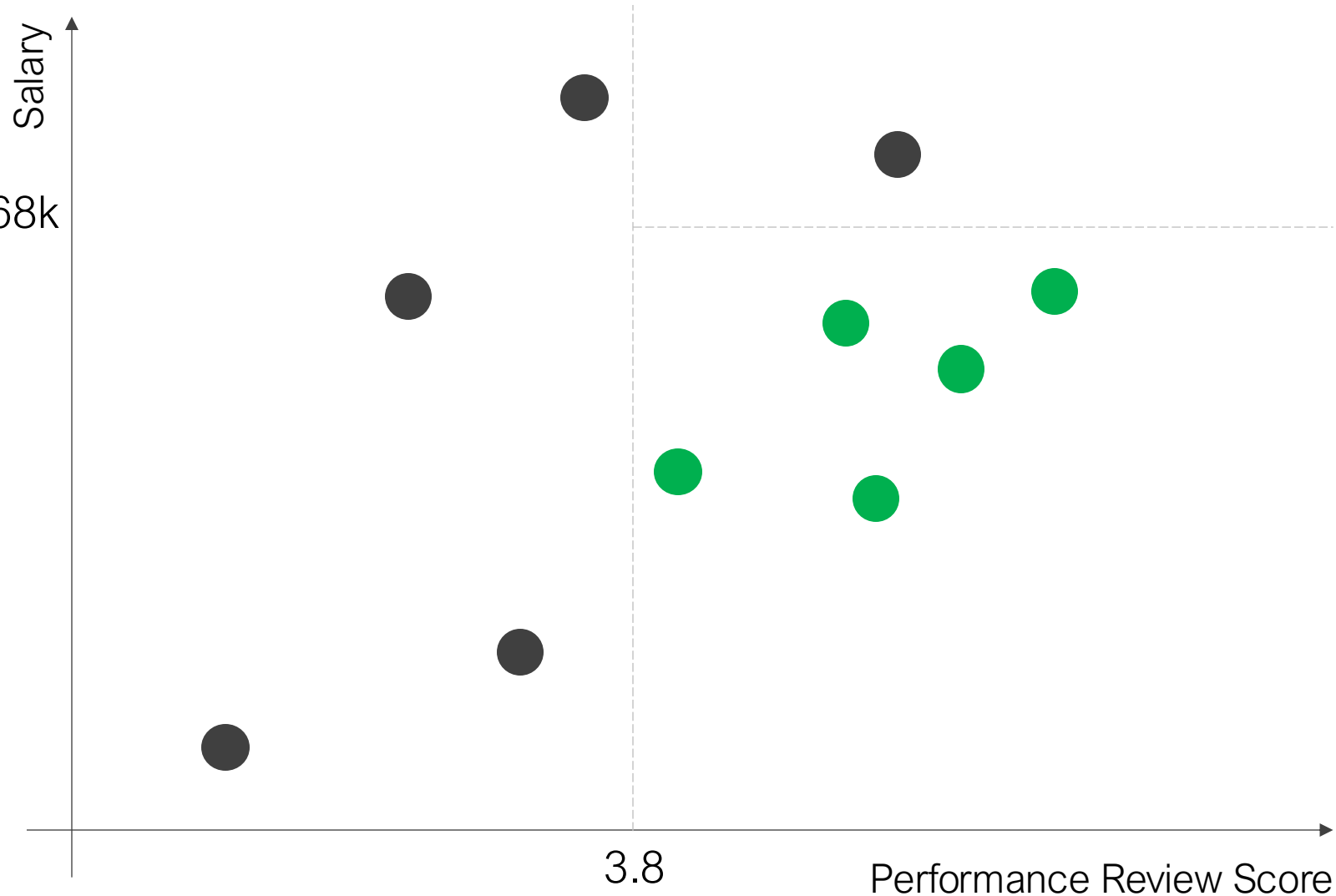
Predicting promotions of salaried employees

**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)

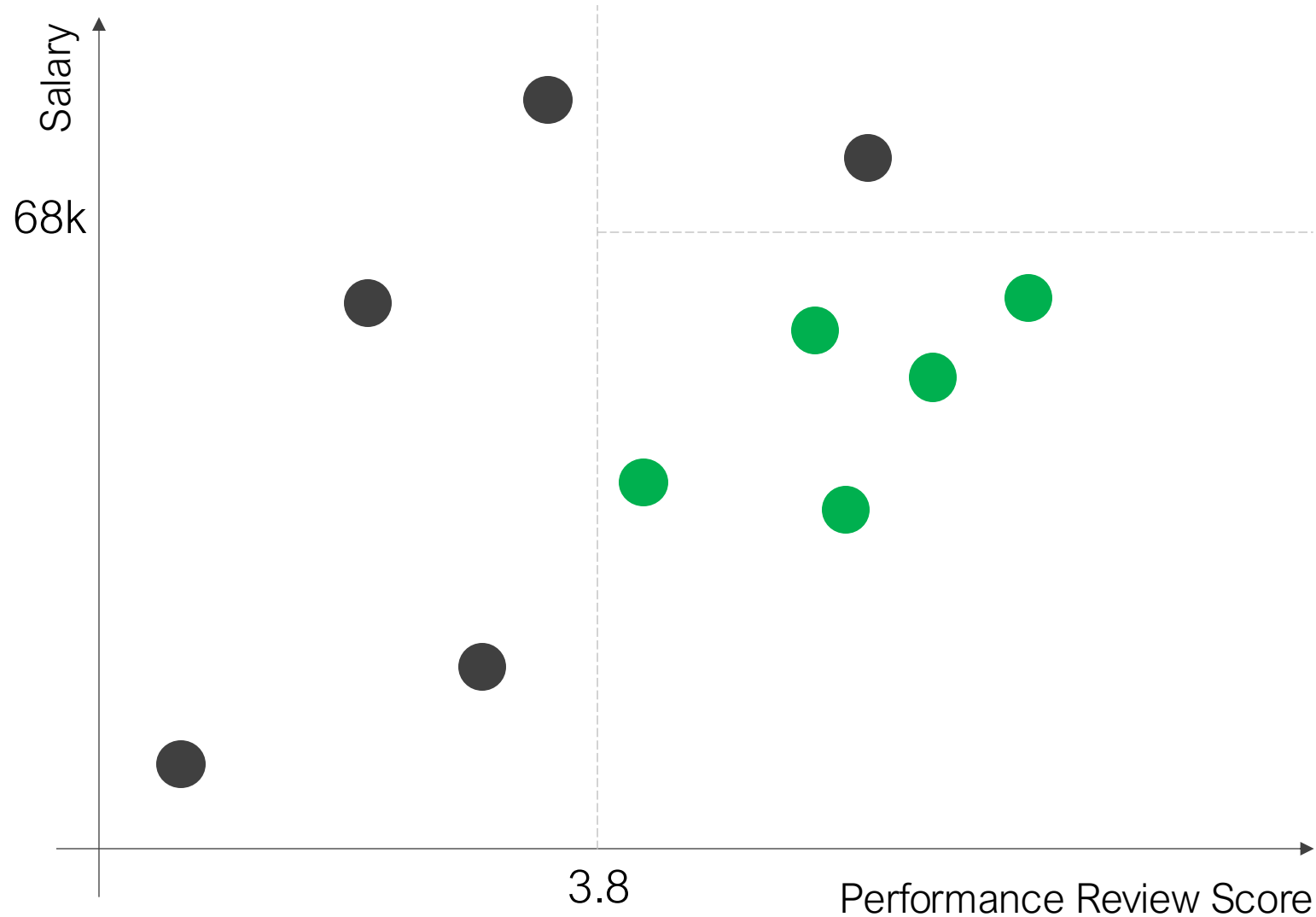
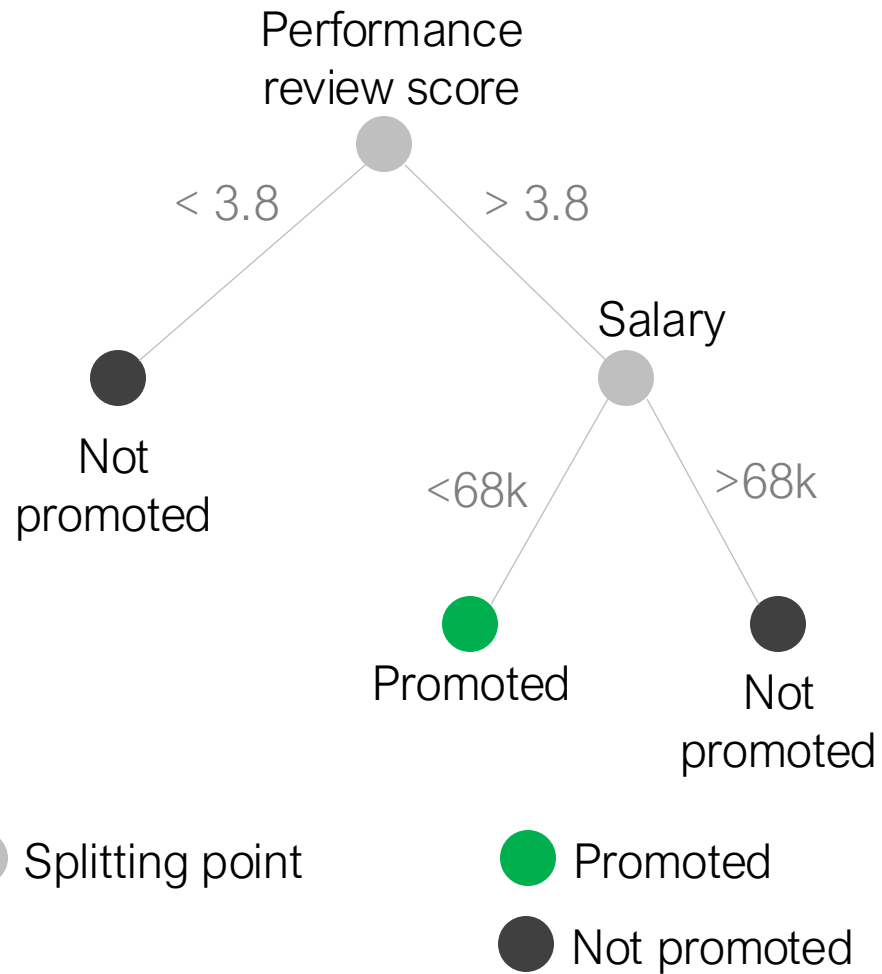
**Greedy, recursive  
binary tree**

● Promoted  
● Not promoted



# Classification and Regression Trees (CART)

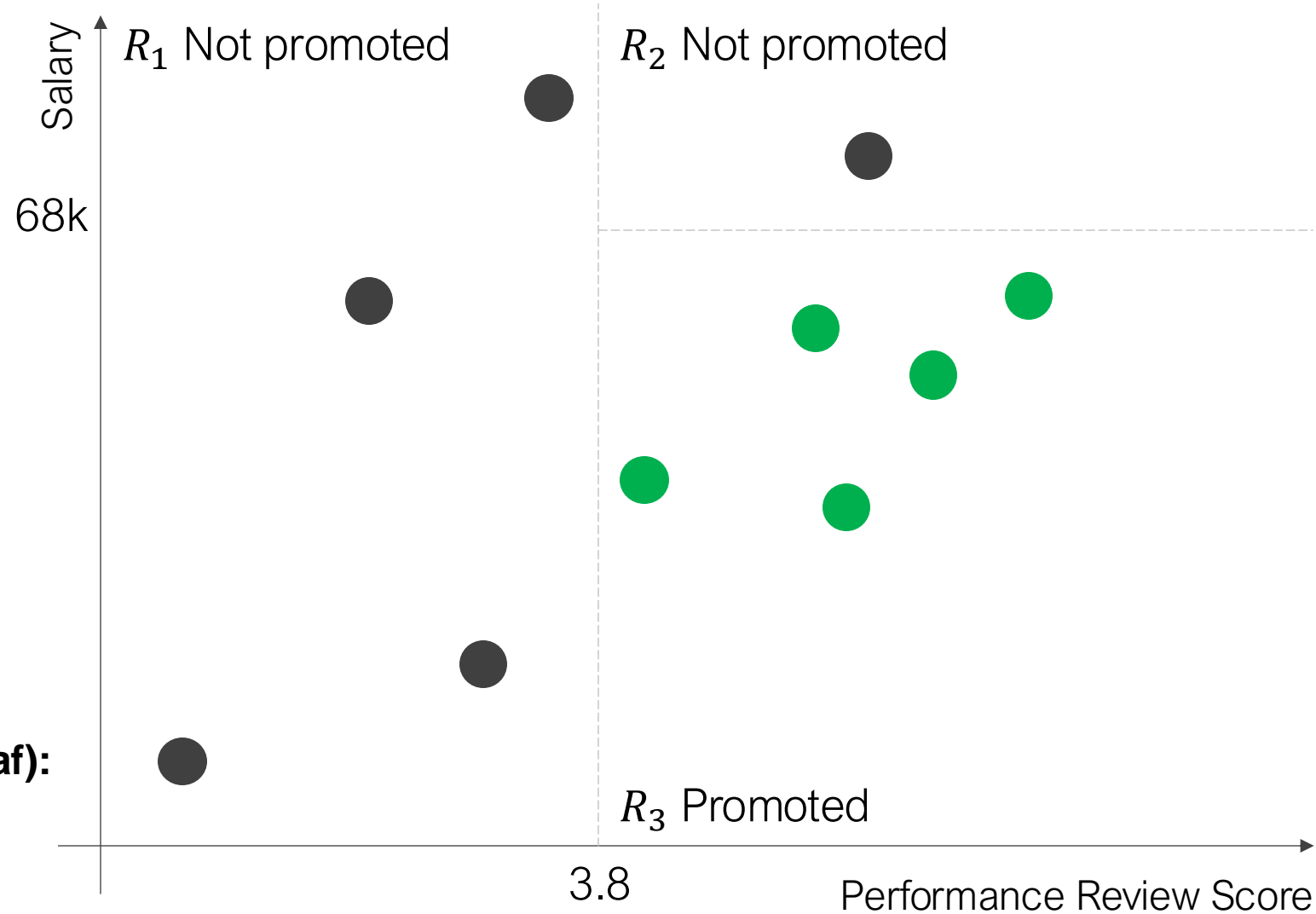
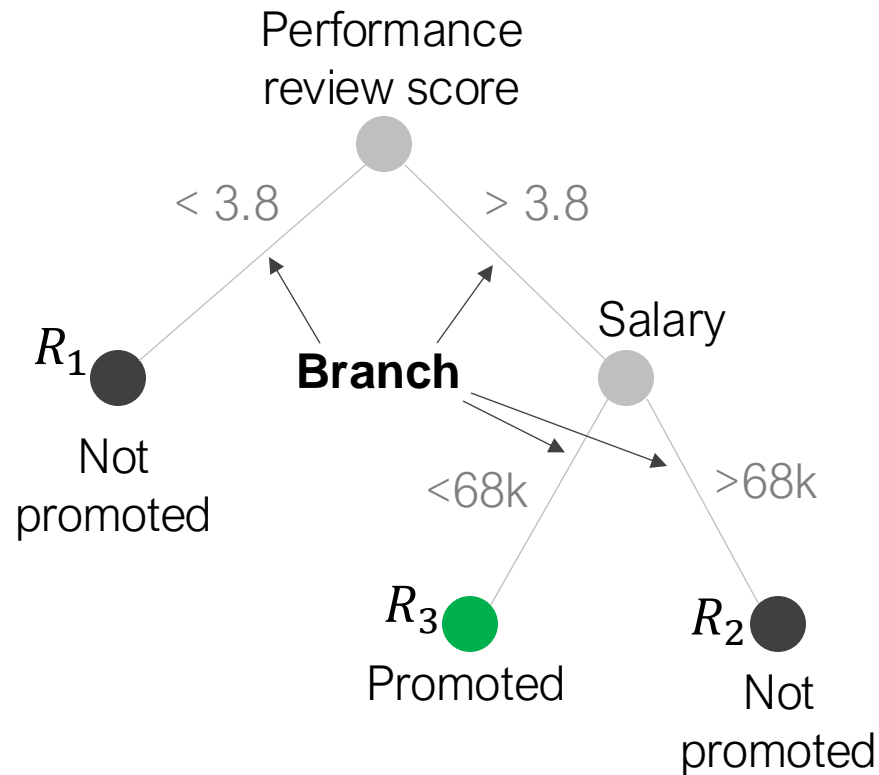
Tree representation:



# Classification and Regression Trees (CART)

Tree representation:

continue splitting until each region only has one class

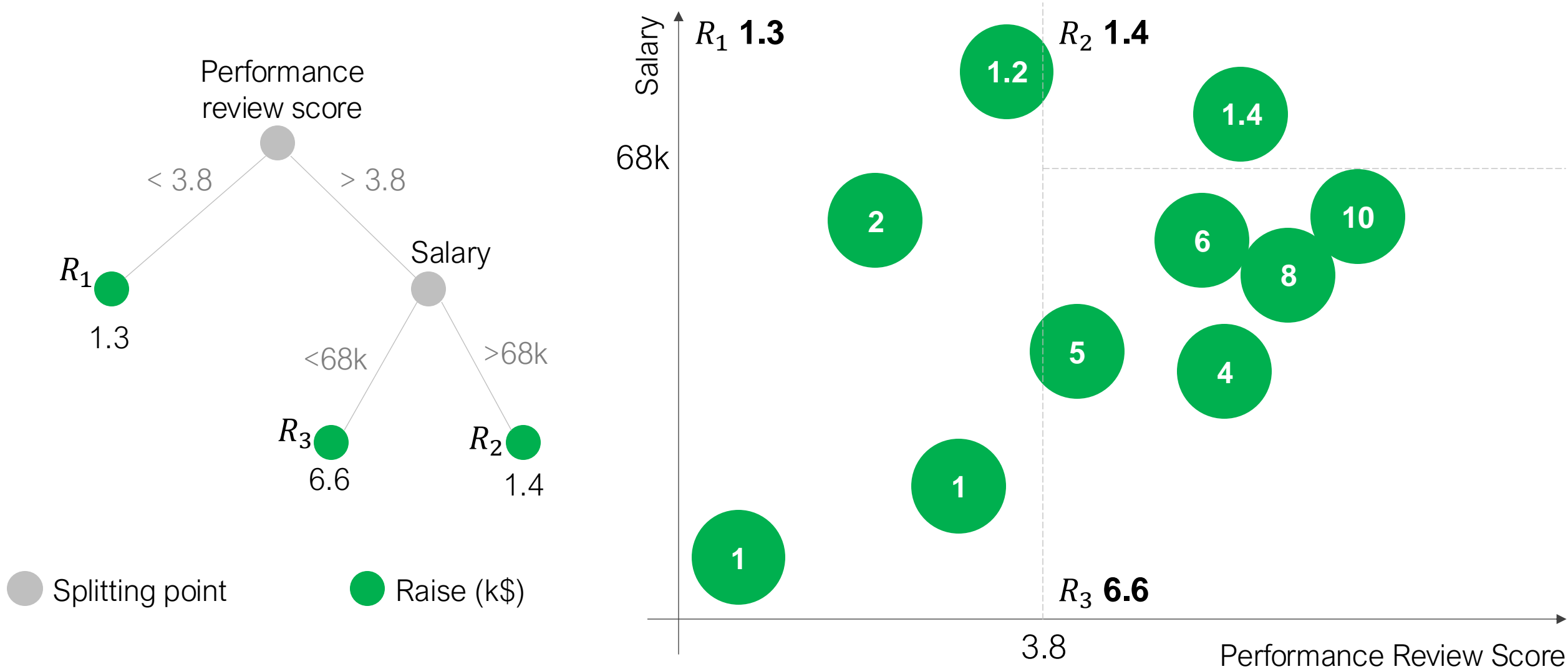


**Internal node:**  
● Splitting point

**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 criterion most after the split

## Splitting criterion

$$C = \sum_{r=1}^{R_{tot}} Q(r)$$

*(Note:  $R_{tot}$  is labeled "total" in the original image)*

## Regression

Mean square error

$$Q_{MSE}(r) = \sum_{i \in R_r} (y_i - \hat{y}_{R_r})^2$$

$y_i$  = training data response  $i$   
 $\hat{y}_{R_r}$  = mean value in region  $r$ , (where  $R_r$  is the set of samples in region  $r$ )

## Classification

Misclassification rate

$$Q_{Misclass}(r) = 1 - \max_k (\hat{p}_{rk})$$

Gini impurity

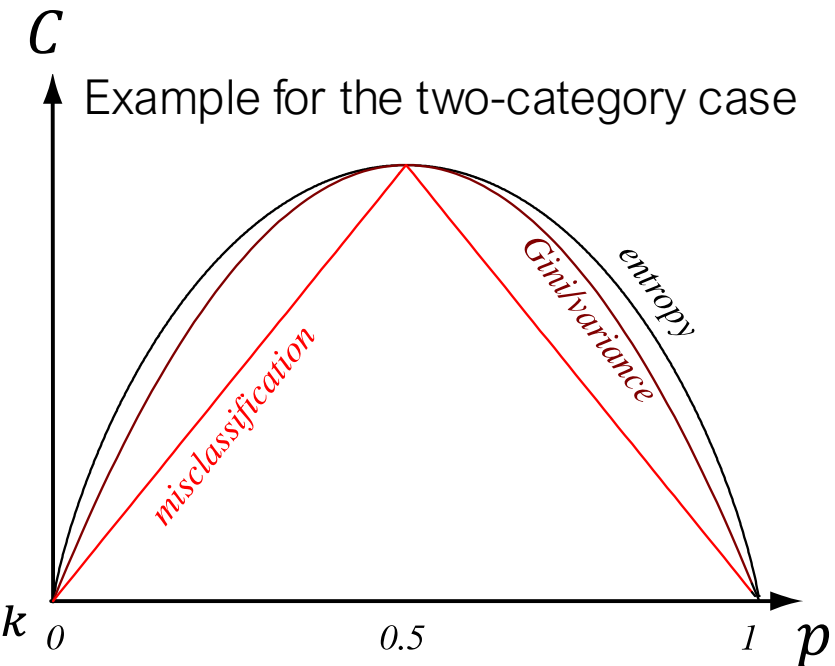
Measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled

$$Q_{Gini}(r) = \sum_{k=1}^K \hat{p}_{rk}(1 - \hat{p}_{rk})$$

Cross-entropy

$$Q_{entropy}(r) = - \sum_{k=1}^K \hat{p}_{rk} \log \hat{p}_{rk}$$

$\hat{p}_{rk}$  = proportion of training observations in the  $r^{\text{th}}$  region from the  $k^{\text{th}}$  class



Duda, Hart, and Stork., Pattern Classification

# How to measure quality of split for classification?

$\hat{p}_{rk}$  = proportion of training observations in the  $r^{\text{th}}$  region from the  $k^{\text{th}}$  class

Class 1 ●  
Class 2 ●

**For each region:**

Misclassification rate

$$Q_{\text{Misclass}}(r) = 1 - \max_k (\hat{p}_{rk})$$

	1	2
	0.333	0.167

Gini impurity

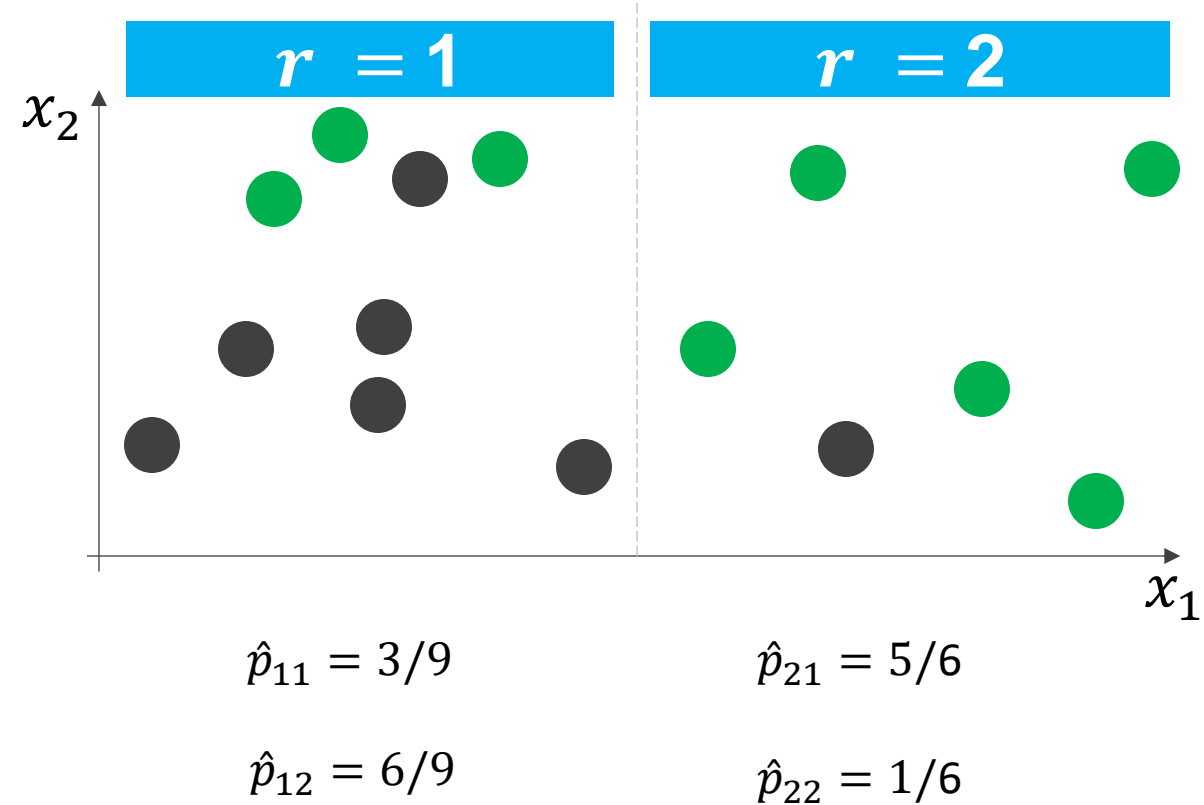
$$Q_{\text{Gini}}(r) = \sum_{k=1}^K \hat{p}_{rk}(1 - \hat{p}_{rk})$$

	0.444	0.278
--	-------	-------

Cross-entropy

$$Q_{\text{entropy}}(r) = - \sum_{k=1}^K \hat{p}_{rk} \log \hat{p}_{rk}$$

	0.637	0.450
--	-------	-------



# Tree Pruning

Trees have the tendency to overfit the data

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

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

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

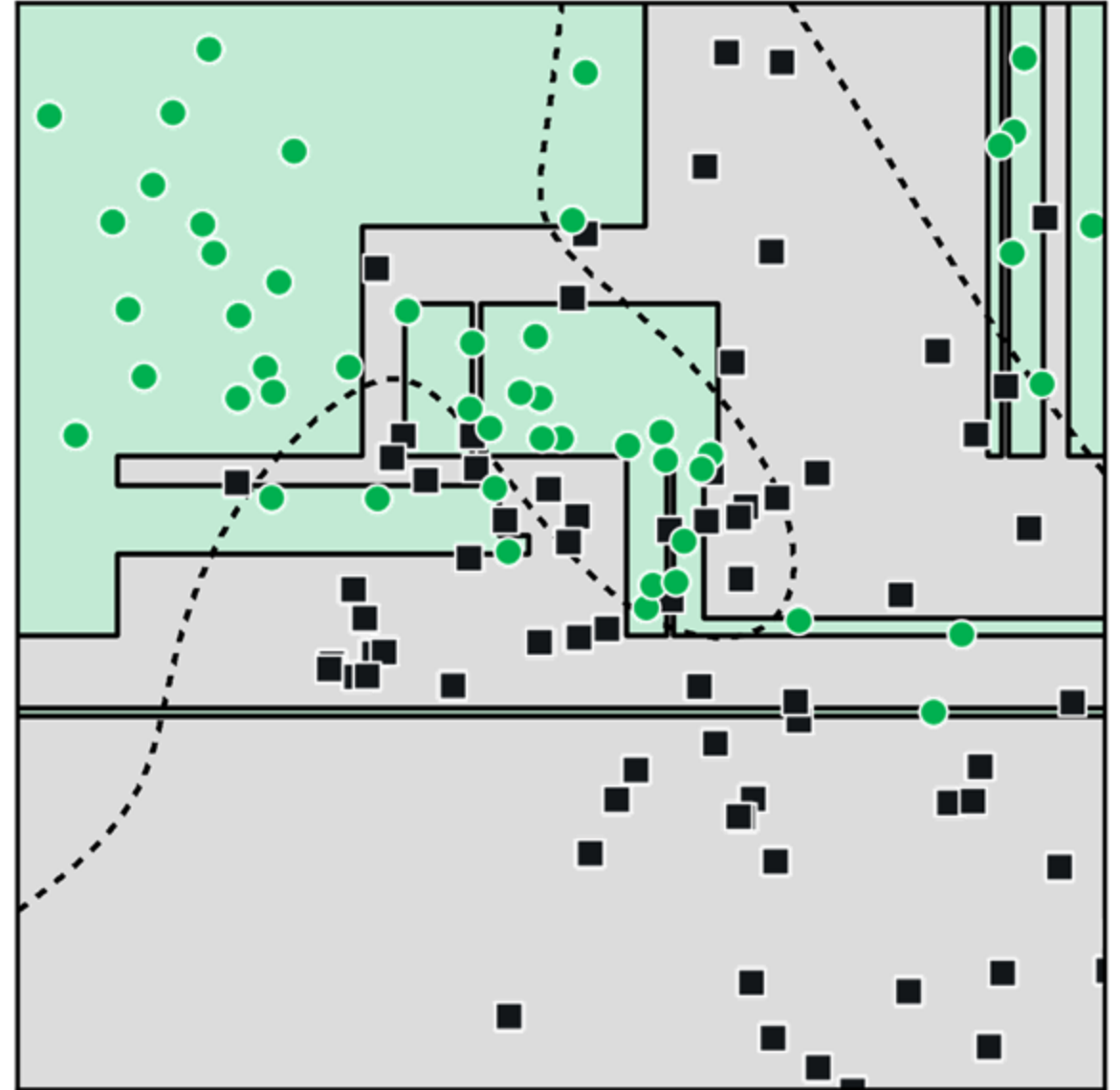
$$C_{Prune} = \sum_{j=1}^T \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2 + \alpha T$$

not used for complex problem

penalty on number of terminal nodes

number of terminal nodes

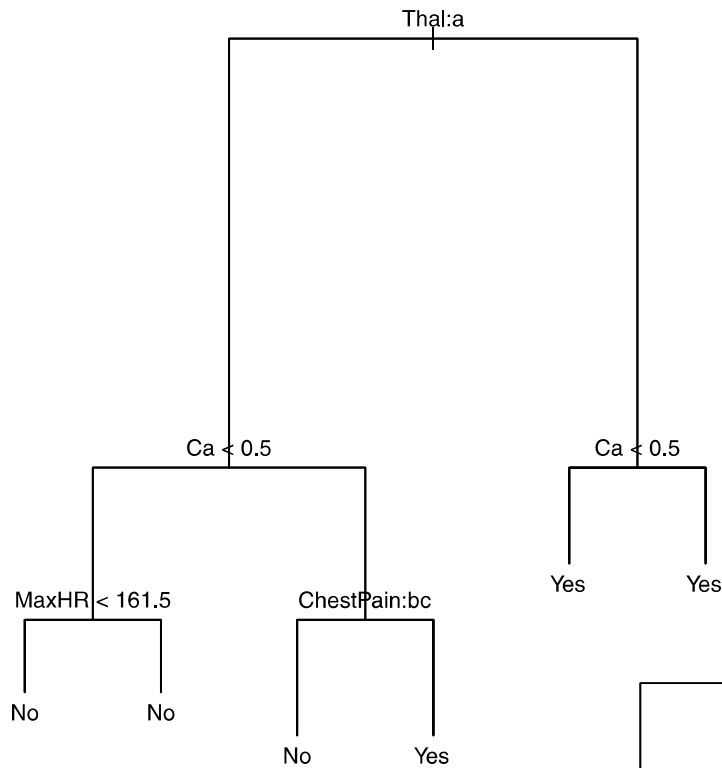
Decision Tree





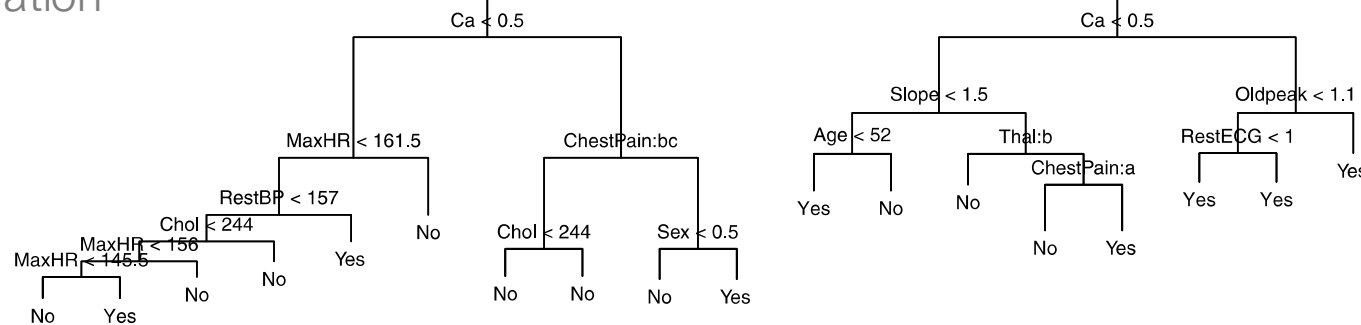
# Pruning example

Pruned Tree

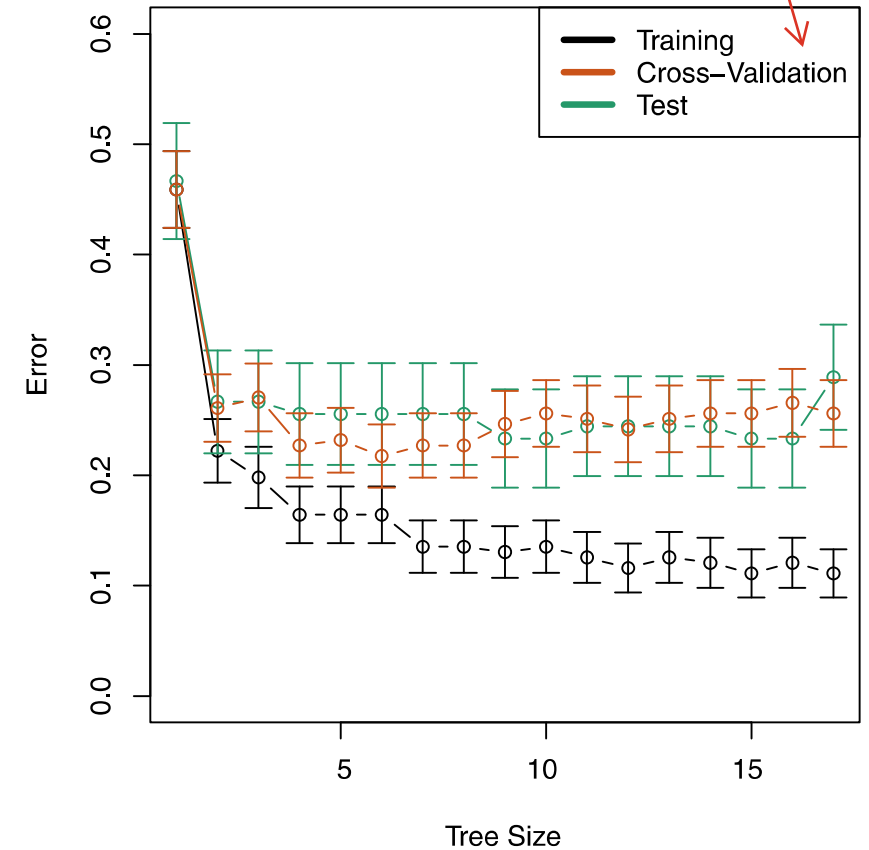


Original Tree

Example: heart disease classification

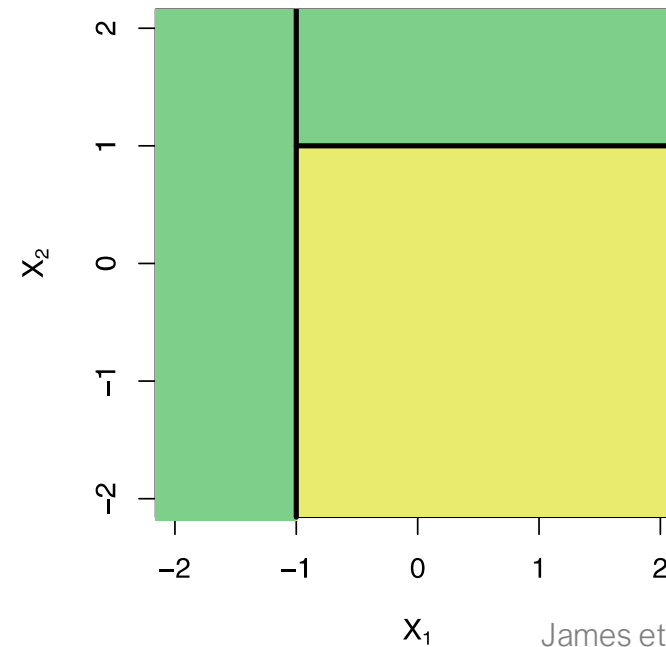
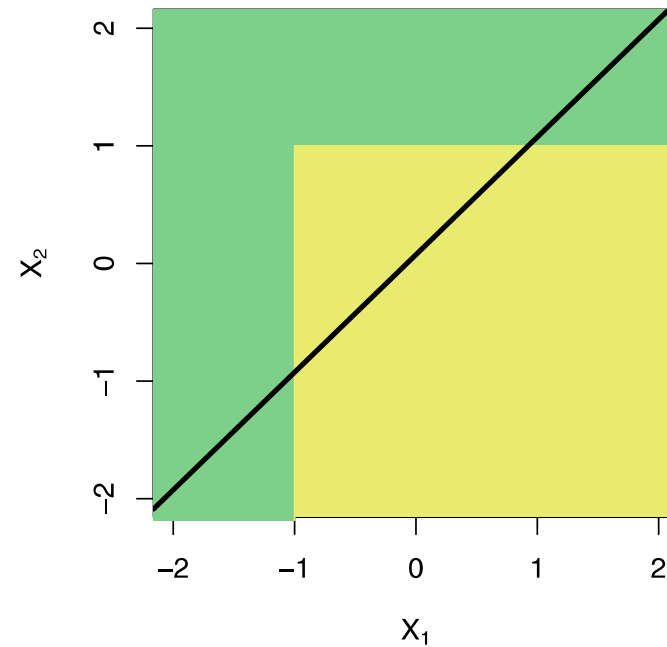
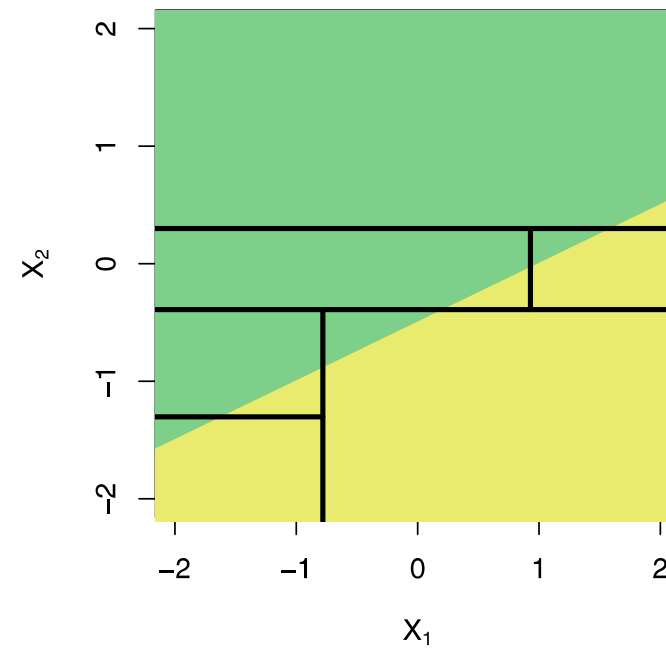
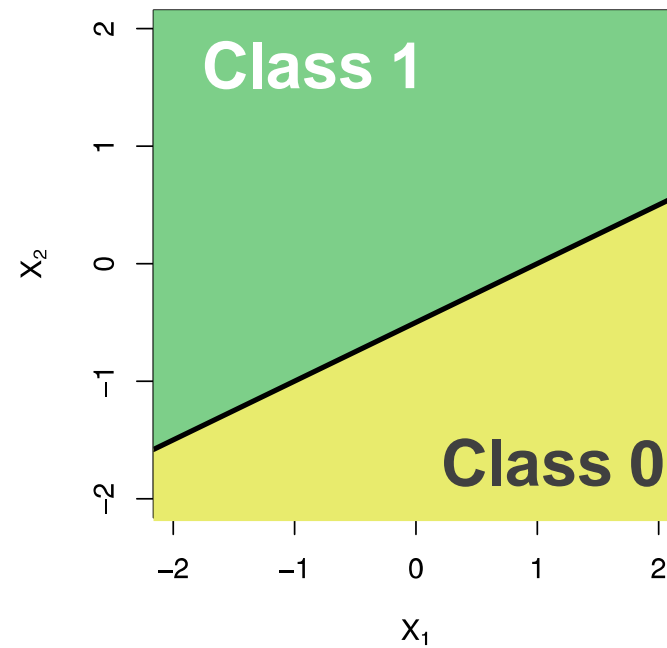


Performance



James et al., An Introduction to Statistical Learning

# Linear model



# Classification Tree

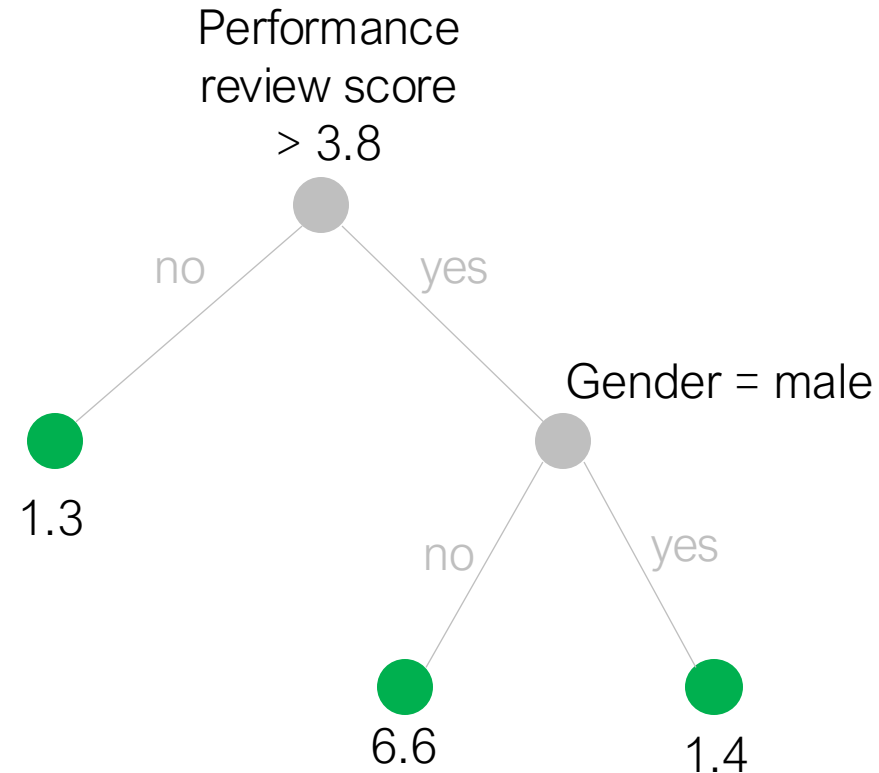
Struggle when the boundary is not parallel to an axis

...nonlinear feature transforms could help...

# 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 other methods

Tend to overfit  
(have high variance)

# Ensemble learning

Combining models to improve performance beyond any individual model alone

Bagging (bootstrap aggregation)

Random forests (tree-specific modification of bagging)

Gradient boosting

# Reducing Variance or Bias through ensembles

## Bagging

Models in ensemble:

high variance, low bias  
(i.e. overfit models)

Effect of aggregating:

Reduce variance through averaging output

## Boosting

high bias, low variance  
(i.e. underfit models, “weak learners”)

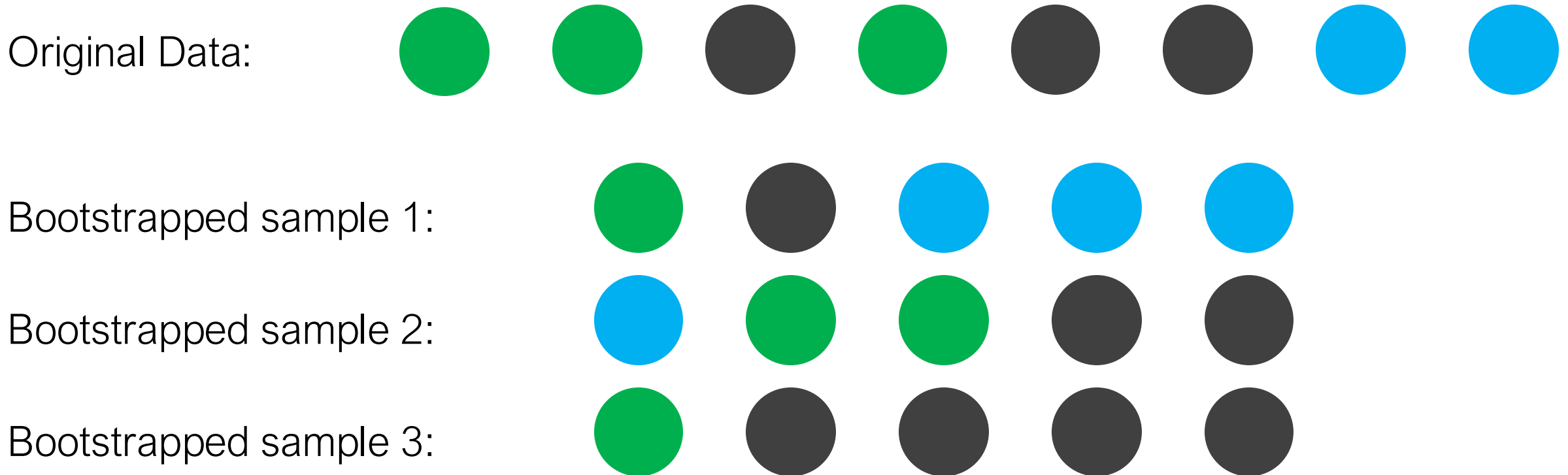
Reduce bias through sequentially fitting models to previous model errors

# Bagging

Bootstrap aggregation

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

Recall bootstrap sampling (sampling with replacement):



# Bagging

Decision trees are a popular  
base model

Bootstrap aggregation

- 1 Create a random bootstrap sample from the training data
- 2 Train a model on that bootstrap sample and call it  $\hat{f}_b(\mathbf{x})$
- 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:

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

(for classification models we can get the average class confidence or take a majority vote)

# 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]

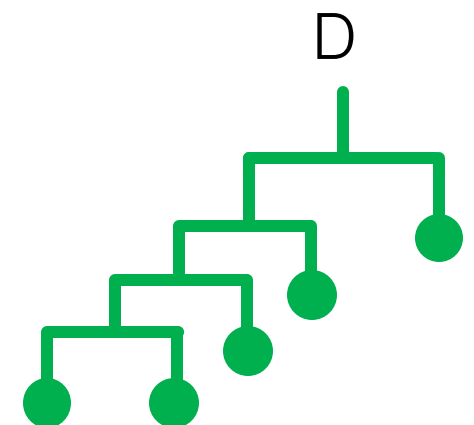
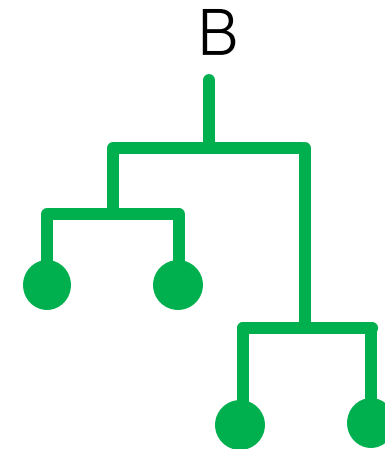
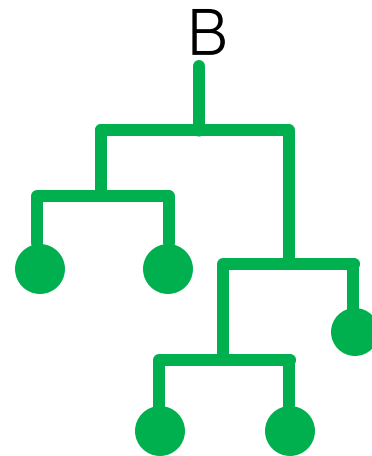
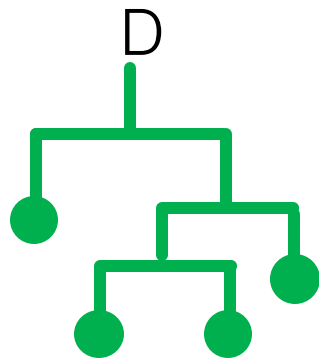
[A, B, C, D]

[A, B, C, D]

[A, B, C, D]

First split:

Trees:





# Variable Importance

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

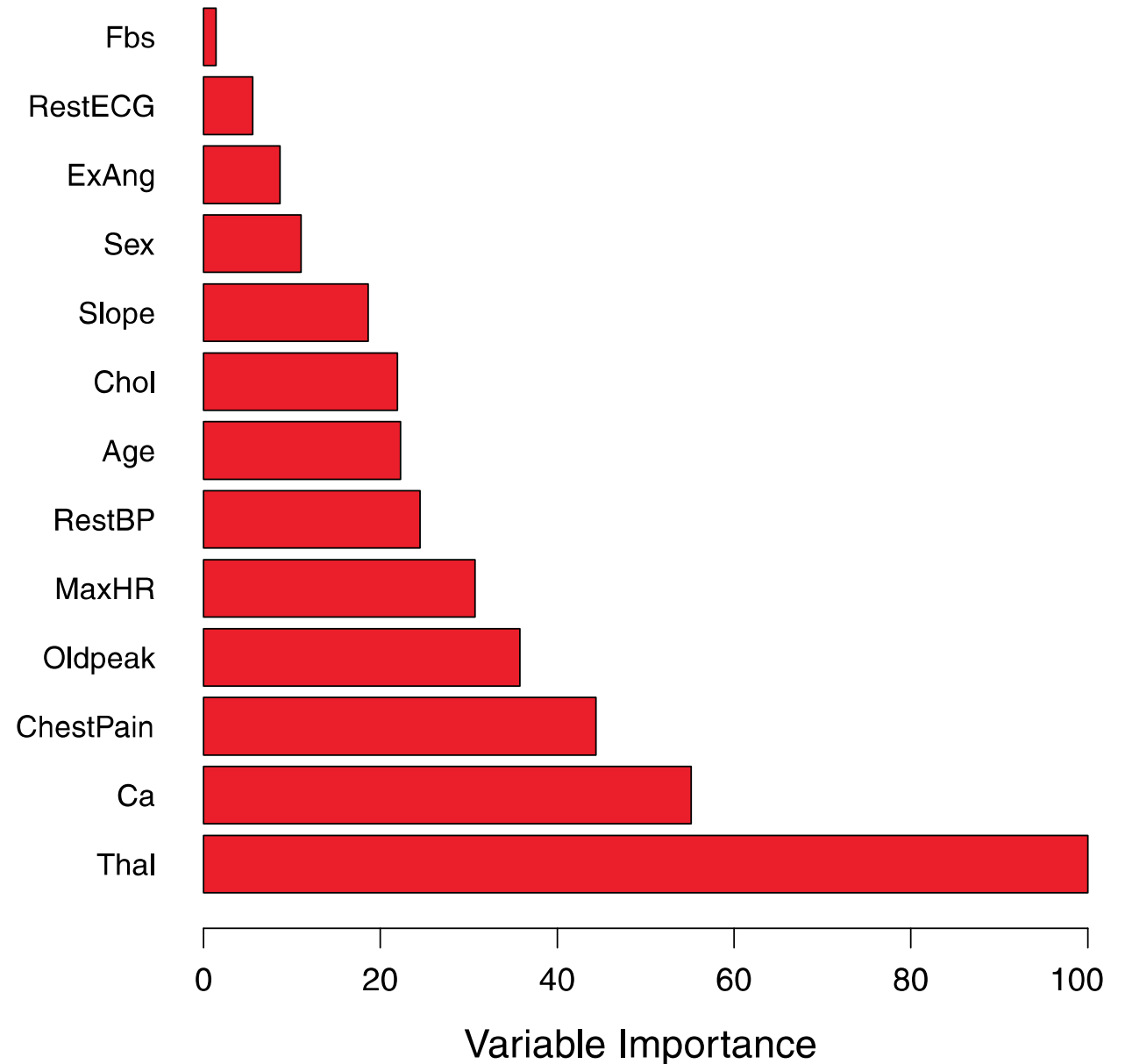
We can construct another measure called “variable importance” to **compare feature contributions**

1

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

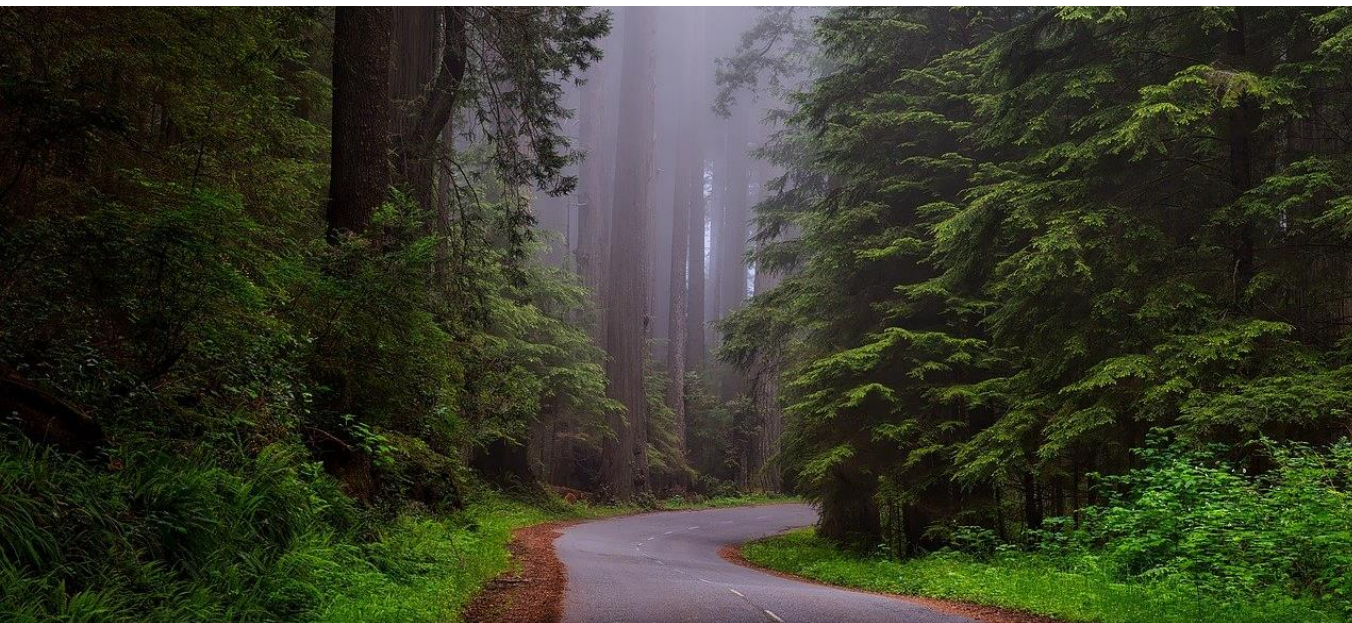
2

Average over all the trees resulting from bagging





# Random Forests





# 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:  $m = \sqrt{p}$

(If  $m = p$ , then we would be back to the bagging approach)

# Bagging

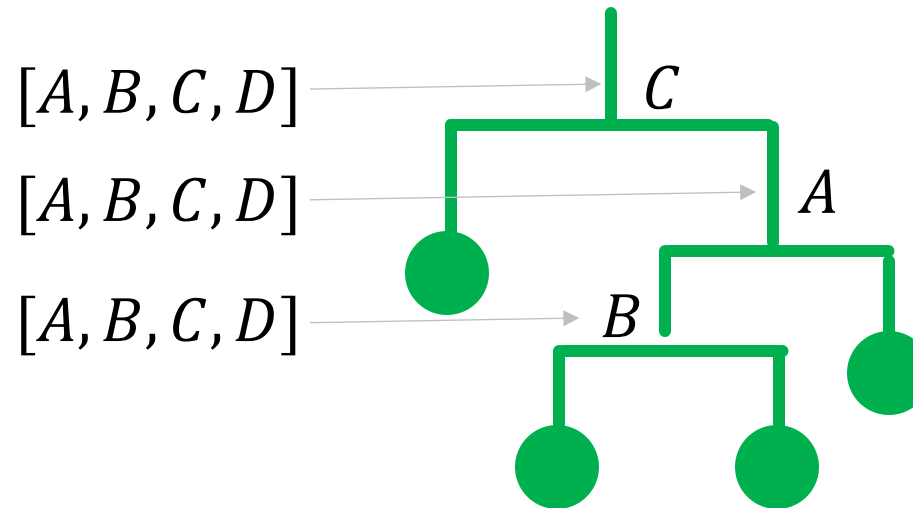
Observations  
Included:  
(out of 1-9)

Features list:

Feature options for  
each split:

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

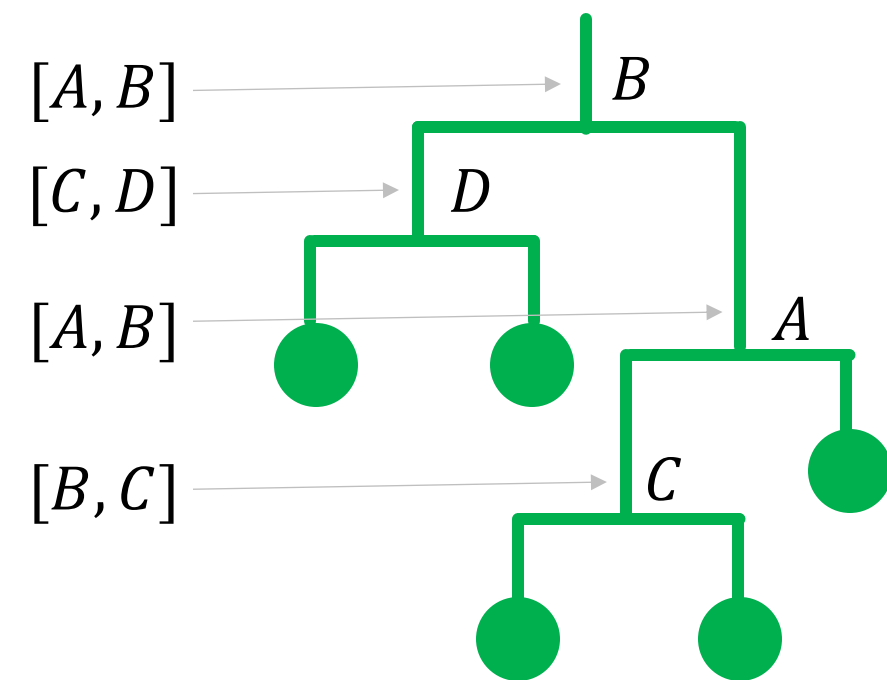
$[A, B, C, D]$



# Random forests

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

$[A, B, C, D]$

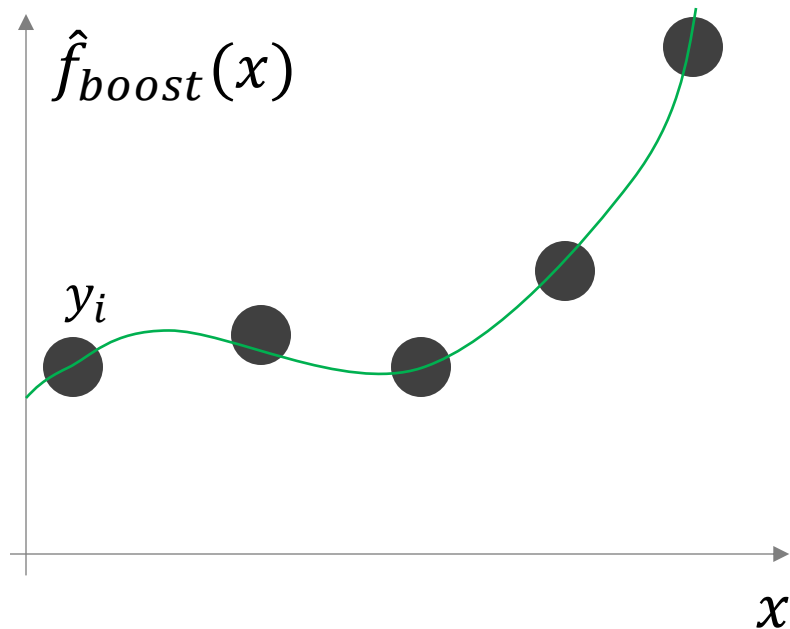
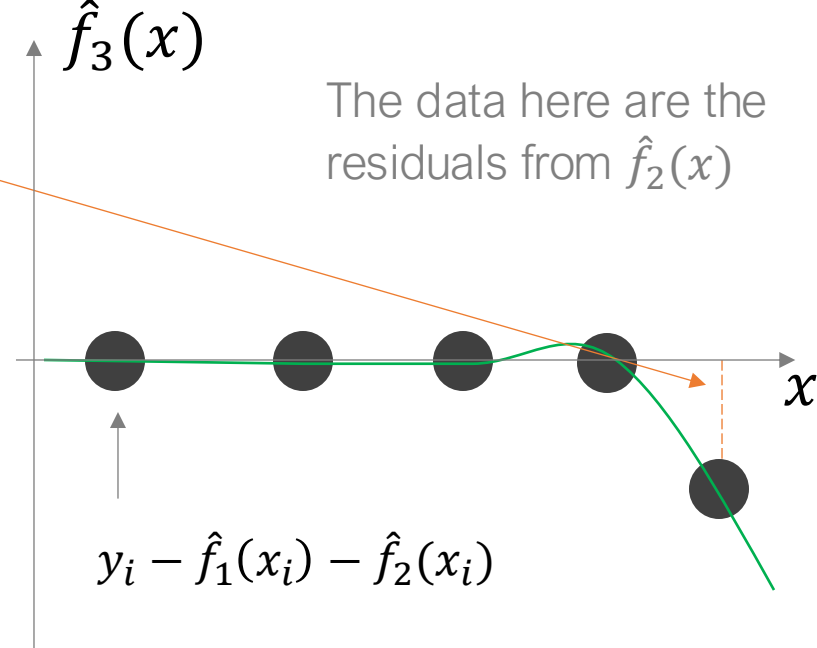
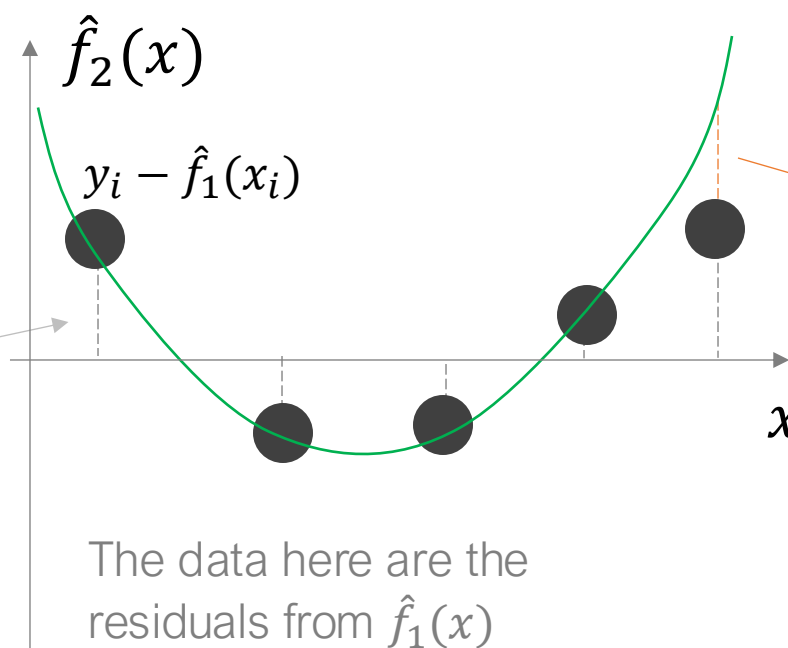
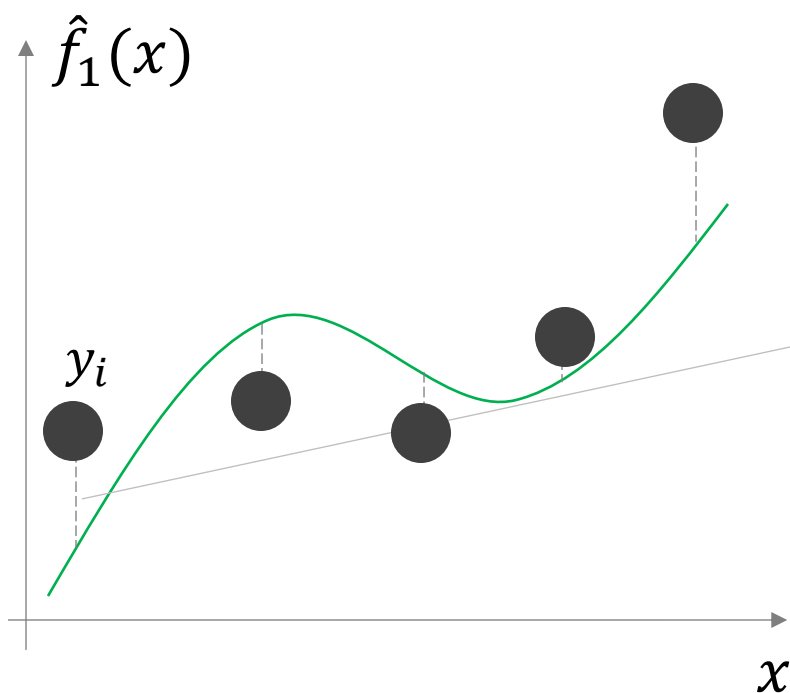


# Boosting

Decision trees are a popular  
base model

**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  $\lambda$
- 2 Set  $\hat{f}(\mathbf{x}) = 0$  and  $r_i = y_i$  for all the training data
- 3 Fit a tree,  $\hat{f}_b(\mathbf{x})$  to the residuals,  $r_i$  (with  $d$  splits)
- 4 Update  $\hat{f}(\mathbf{x}) \leftarrow \hat{f}(\mathbf{x}) + \lambda \hat{f}_b(\mathbf{x})$
- 5 Update the residuals  $r_i \leftarrow r_i - \lambda \hat{f}_b(\mathbf{x}_i)$
- 6 Output the boosted model: 
$$\hat{f}(\mathbf{x}) = \sum_{b=1}^B \lambda \hat{f}_b(\mathbf{x})$$

$\lambda$  slows down the learning process to avoid overfitting

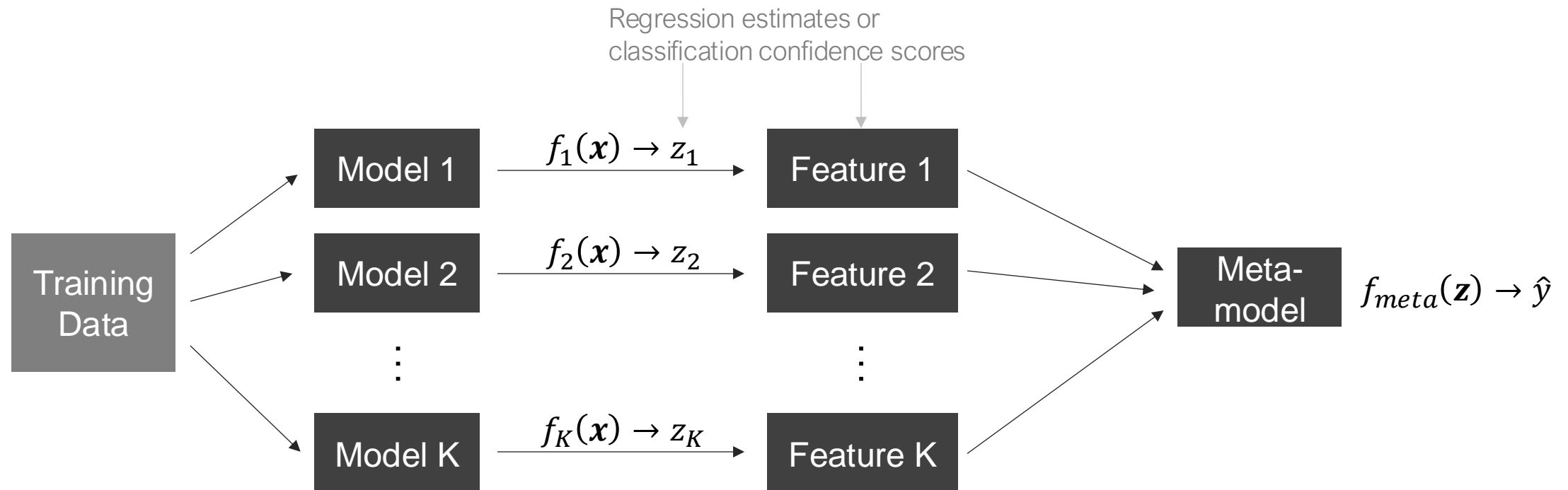
Often this is just a small number of splits (a stump)

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 as **features**





# Supervised Learning Techniques

Covered so far

- Linear Regression
- K-Nearest Neighbors
- Perceptron
- Logistic Regression
- Linear Discriminant Analysis
- Quadratic Discriminant Analysis
- Naïve Bayes
- Decision Trees and Random Forests
- Ensemble methods (bagging, boosting, stacking)

Can be used with numerous machine learning techniques, often CART

Appropriate for:  
● Classification  
● Regression