

# Introduction to Machine Learning

Session 2a: Introduction to Classification and Regression Trees

Reto Wüest

Department of Political Science and International Relations  
University of Geneva

- ① The Basics of Decision Trees
- ② Regression Trees
  - Example: Baseball Salary Data
  - Terminology for Trees
  - Building a Regression Tree
  - Tree Pruning
- ③ Classification Trees
  - Building a Classification Tree

# The Basics of Decision Trees

# The Basics of Decision Trees

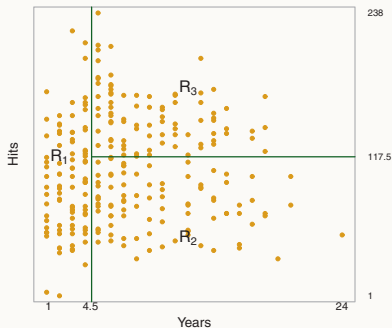
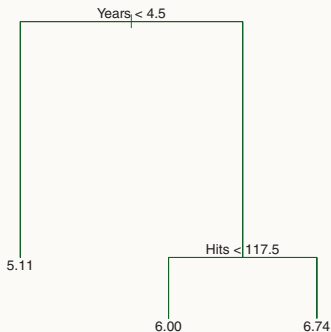
- Tree-based methods **stratify** or **segment** the predictor space into a number of simple regions.
- To make a **prediction** for a test observation, we use the mean or mode of the training observations in the region to which it belongs.
- These methods are called **decision-tree** methods because the splitting rules used to segment the predictor space can be summarized in a tree.
- Decision trees can be applied to both **regression** and **classification** problems.

# Regression Trees

## Example: Baseball Salary Data

The goal is to predict a baseball player's (log) salary based on the number of years played in the major leagues and the number of hits in the previous year.

Regression Tree Fit to Baseball Salary Data



(Source: James et al. 2013, 304f.)

- Regions  $R_1$ ,  $R_2$ , and  $R_3$  above are the **terminal nodes** or **leaves** of the tree.
- Points along the tree where the predictor space is split are the **internal nodes** (indicated above by  $\text{Years} < 4.5$  and  $\text{Hits} < 117.5$ ).
- Segments of the tree that connect the nodes are called **branches**.

Roughly speaking, there are two steps:

- 1 Divide the predictors space (i.e., the set of possible values for predictors  $X_1, X_2, \dots, X_p$ ) into  $J$  **distinct** and **non-overlapping** regions,  $R_1, R_2, \dots, R_J$ .
- 2 Make the **same prediction** for every test observation that falls into region  $R_j$ , which is the **mean of the response values** for the training observations in  $R_j$ .



Step 1 (more detailed):

- How do we construct the regions  $R_1, \dots, R_J$ ?
- We divide the predictor space into **high-dimensional rectangles** (boxes),  $R_1, \dots, R_J$ , so that they minimize the RSS

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

where  $\hat{y}_{R_j}$  is the **mean response** of the training observations in the  $j$ th box.

## Step 1 (more detailed):

- It is computationally **not feasible** to consider every possible partition of the predictor space into  $J$  boxes.
- Therefore, we take a **top-down, greedy** approach that is known as **recursive binary splitting**:
  - Top-down: we **begin at the top** of the tree (where all observations belong to a single region) and **successively split** the predictor space;
  - Greedy: we make the split **that is best at each particular step** of the tree-building process (i.e., we do not look ahead and pick a split that will lead to a better tree in some future step).

### Step 1 (more detailed):

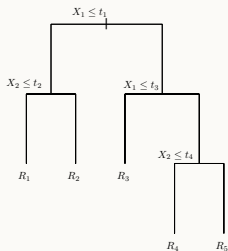
- How do we perform recursive binary splitting?
- We first select the predictor  $X_j$  and the cutpoint  $s$  such that splitting the predictor space into the regions  $\{X \mid X_j < s\}$  and  $\{X \mid X_j \geq s\}$  leads to the greatest possible reduction in RSS. (We now have two regions.)
- Next, we again select the predictor and the cutpoint that minimize the RSS, but this time we split one of the two previously identified regions. (We now have three regions.)

Step 1 (more detailed):

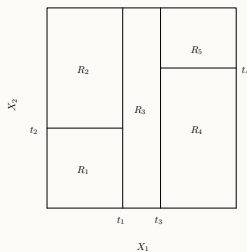
- Next, we split one of the three regions further, so as to minimize the RSS. (We now have four regions.)
- We continue this process until a **stopping criterion** is reached.
- Once the regions  $R_1, \dots, R_J$  have been created, we **predict the response** for a **test observation** using the mean of the **training observations** in the region to which the test observation belongs.

# Building a Regression Tree: Example

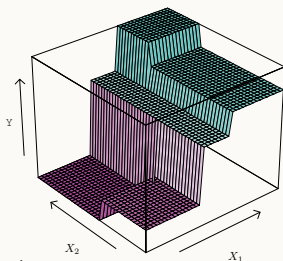
## Decision Tree



## Predictor Space



## Prediction Surface



- The above process may produce good predictions on the training set, but it likely to **overfit** the data, leading to poor **test set performance**.
- The reason is that the resulting tree might be too complex. A **less complex tree** (fewer splits) might lead to **lower variance** at the **cost of a little bias**.
- A less complex tree can be achieved by **tree pruning**: grow a **very large** tree  $T_0$  and then prune it back in order to obtain a **subtree**.

- How do we find the **best subtree**?
- Our goal is to select a subtree that leads to the **lowest test error rate**.
- For each subtree, we could estimate its test error using **cross-validation (CV)**.
- However, this approach is **not feasible** as there is a **very large** number of possible subtrees.
- **Cost complexity pruning** allows us to select only a **small set** of subtrees for consideration.

Cost complexity pruning:

- Let  $\alpha$  be a **tuning parameter**. For each value of  $\alpha$ , there is a subtree  $T \subset T_0$  that minimizes

$$\sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|, \quad (2)$$

where  $|T|$  is the number of terminal nodes of tree  $T$ .

- The tuning parameter  $\alpha$  controls the **trade-off** between the subtree's **complexity** and its **fit to the training data**.
- With increasing  $\alpha$ , quantity (2) will be minimized for a smaller subtree. (Note the similarity to the Lasso!)



Cost complexity pruning:

- We can then select the **optimal value** of  $\alpha$  using CV.
- Finally, we return to the **full data set** and obtain the subtree corresponding to the optimal value of  $\alpha$ .

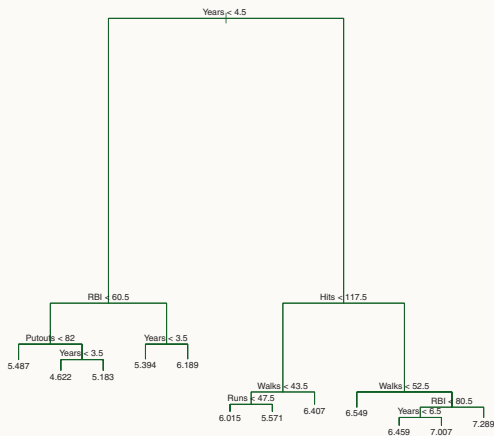
## Algorithm: Fitting and Pruning a Regression Tree

- 1 Use recursive binary splitting to grow a large tree on the training data.
- 2 Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of  $\alpha$ .
- 3 Use  $K$ -fold CV to choose  $\alpha$ . That is, divide the training observations into  $K$  folds. For each  $k = 1, \dots, K$ :
  - (a) Repeat Steps 1 and 2 on all but the  $k$ th fold of the training data.
  - (b) Evaluate the mean squared prediction error on the data in the left-out  $k$ th fold, as a function of  $\alpha$ .Average the results for each value of  $\alpha$ , and choose  $\alpha$  to minimize the average error.
- 4 Return the subtree from Step 2 that corresponds to the chosen value of  $\alpha$ .

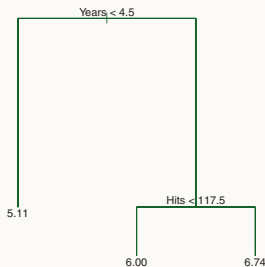
# Tree Pruning: Example

Fitting and Pruning a Regression Tree on the Baseball Salary Data

Unpruned Tree



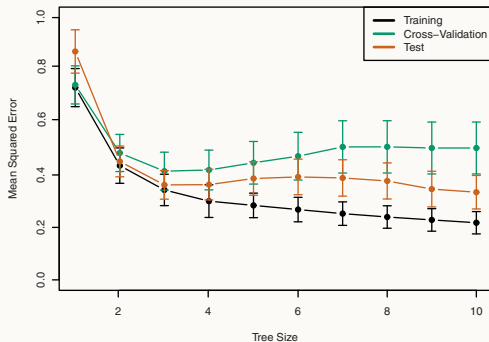
Pruned Tree



(Source: James et al. 2013, 304 & 310)

# Tree Pruning: Example

## Fitting and Pruning a Regression Tree on the Baseball Salary Data



(Source: James et al. 2013, 311)

The CV error is a reasonable approximation of the test error. The CV error takes on its minimum for a three-node tree (see previous slide).

# Classification Trees

- **Classification trees** are very similar to regression trees, except that they are used to predict a **qualitative** rather than a quantitative response.
- For a **regression tree**, the predicted response for an observation is given by the **mean response** of the training observations that belong to the same terminal node.
- For a **classification tree**, the predicted response for an observation is the **most commonly occurring class** among the training observations that belong to the same terminal node.

## Building a Classification Tree

- Just as in the regression setting, we use **recursive binary splitting** to grow a classification tree.
- However, in the classification setting, RSS **cannot be used** as a criterion for making binary splits.
- We could use the **classification error rate**, which is the fraction of training observations in a terminal node that do not belong to the most common class

$$E = 1 - \max_k(\hat{p}_{mk}), \quad (3)$$

where  $\hat{p}_{mk}$  represents the proportion of training observations in the  $m$ th terminal node that are from the  $k$ th class.

## Building a Classification Tree

- However, it turns out that classification error is **not sufficiently sensitive** for tree-growing.
- Therefore, two other measures are preferable: the **Gini index** and **entropy**.
- The **Gini index** is a measure of total variance across the  $K$  classes:

$$G = \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}). \quad (4)$$

It takes on a small value if all of the  $\hat{p}_{mk}$ 's are close to 0 or 1. Therefore, a small value indicates that a node contains predominantly observations from a single class (node **purity**).



- An alternative to the Gini index is the **entropy**, given by

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}. \quad (5)$$

(Note that since  $0 \leq \hat{p}_{mk} \leq 1$ , it is  $0 \leq -\hat{p}_{mk} \log \hat{p}_{mk}$ .)

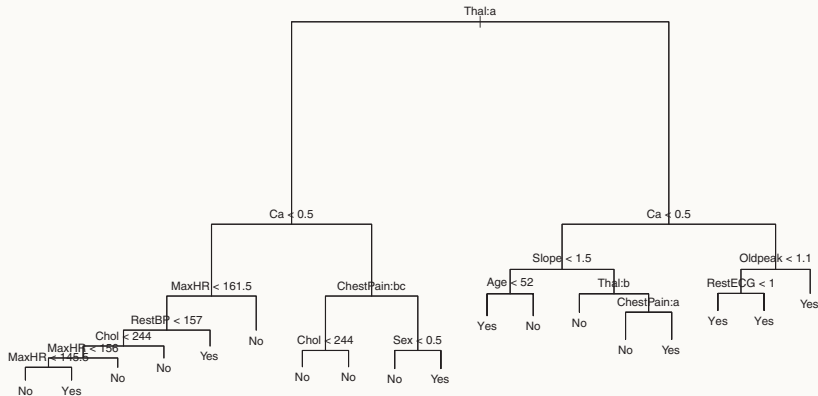
- The entropy will take on a value near 0 if the  $\hat{p}_{mk}$ 's are all near 0 or 1. Therefore, like the Gini index, the entropy will take on a small value if the  $m$ th node is **pure**.

## Building a Classification Tree

- **Building a classification tree:** either the **Gini index** or the **entropy** is used to evaluate the quality of a particular split, since these measures are more sensitive to node purity than the classification error rate.
- **Pruning the tree:** any of the three measures might be used, but the **classification error rate** is preferable if prediction accuracy of the final tree is the goal.

# Building a Classification Tree: Example

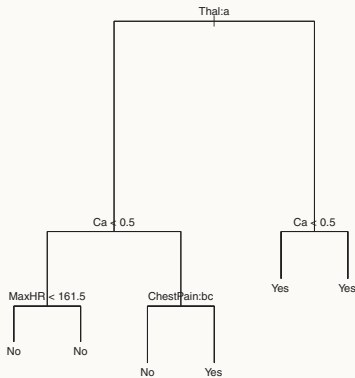
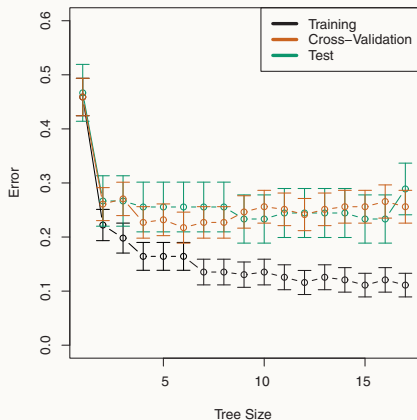
## Fitting and Pruning a Classification Tree on the Heart Disease Data



(Source: James et al. 2013, 313)

# Building a Classification Tree: Example

## Fitting and Pruning a Classification Tree on the Heart Disease Data



(Source: James et al. 2013, 313)

## Building a Classification Tree: Example

- Note that in the above example, some of the splits yielded two terminal nodes that have the **same predicted value**.
- Why are these splits performed at all?
- Such splits lead to **increased node purity** (they do not reduce the classification error, but they improve the Gini index and the entropy, which are more sensitive to node purity).
- Node purity is important because it tells us something about how **certain** we are when making a prediction.