

# Decision Trees

Jordi Villà i Freixa

Universitat de Vic - Universitat Central de Catalunya  
Study Abroad

*jordi.villa@uvic.cat*

course 2023-2024

# Index

- 1 Tree-based methods
- 2 Top down construction trees
- 3 Additional considerations
- 4 Controlling the tree shape
- 5 Bootstrapping aggregation
- 6 Random Forests
- 7 Boosting
- 8 Bibliography

# Preliminary note

The material in these slides is strongly based on [1]. When other materials are used, they are cited accordingly.

Mathematical notation follows as good as it can a [good practices proposal](#) from the Beijing Academy of Artificial Intelligence.

# What to expect?

In this session we will discuss:

- Decision trees
- Random Forests
- Boosting

# Tree-based methods

- Simple, intuitive and powerful for both regression and classification
- The method divides a feature space  $X$  into smaller regions and fit a simple prediction function for each region.

**Regression** eg, take the mean of the training responses associated with the training features that fall in the specific region

**Classification** eg, take the majority vote among corresponding response variables.

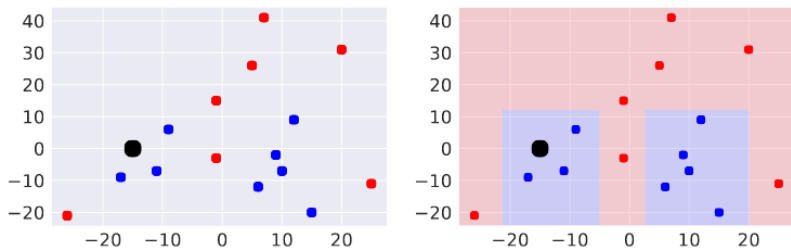
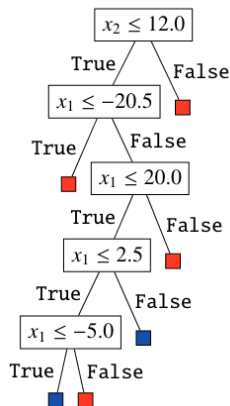


Figure 1: Left: training data and a new feature. Right: a partition of the feature space[1]

- In the example above, we cannot separate regions linearly, but we can create rectangles in the  $\mathbb{R}^2$  space.
- The classifier  $g$ , thus takes a color "blue" or "red" according to where the new black dot goes.
- Both the classification procedure and the partitioning can be represented by a binary *decision tree*.

# Partitioning of feature space $X$



The decision tree corresponding to Fig. 1[1]

- Each node  $v$  corresponds to a region  $\mathcal{R}_v$  of the feature space  $X$ .
- The root node is the feature space  $X$  itself.
- The final (undivided) leaves  $w_1, w_2, \dots$  (red/blue squares) form a partition of  $X$ , as they are disjoint and their union is  $X$ .
- Regional prediction functions  $g^w$  are associated with each leaf.



# Decision

Let us take the input  $\mathbf{x} = (x_1, x_2)^T$ :

- ① we start at the root node, which contains a condition  $x_2 \leq 12.0$  that the input data satisfies;
- ② We then proceed to the left child which contains the condition  $x_2 \leq -20.5$ , which our data does not satisfy, so we go to the next right leaf and so on.

More generally, a binary tree  $\mathbb{T}$  will partition the feature space into as many regions as leaf nodes, and the prediction function becomes:

$$g(\mathbf{x}) = \sum_{w \in \mathcal{W}} g^w(\mathbf{x}) \mathbb{1}\{\mathbf{x} \in \mathcal{R}_w\}; \quad \ell_\tau = \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, g(\mathbf{x}_i))$$

This depends on (1) how the regions  $\mathcal{R}_w$  are constructed and (2) how the regional prediction functions  $g^w$  of the leaf nodes are defined.  $g^w$  is category in a classification setting and real-valued in regression.

## Exercise 1

Training loss Can you provide an example of the fact that any training set  $\tau = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$  can be fitted via a tree with zero training loss? (Hint: imagine classifying the students in the class based on their age in days).

Discuss on the predictability of such model. Can you consider it overfitted?

# Construction of the tree

Let us take the training set  $\tau = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ . We consider the recursive algorithm:

- 1 First, we specify a *splitting rule* for each node  $v$  (a binary 1-0 or a True-False function).
- 2 Any subset of the training set can then be splitted into two sets.
- 3 We can define the recursive function `ConstructSubtree` to repeat the process.

**Algorithm 8.2.1:** Construct\_Subtree

---

**Input:** A node  $v$  and a subset of the training data:  $\sigma \subseteq \tau$ .  
**Output:** A (sub) decision tree  $\mathbb{T}_v$ .

```

1 if termination criterion is met then                                //  $v$  is a leaf node
2   |   Train a regional prediction function  $g^v$  using the training data  $\sigma$ .
3 else                                                                // split the node
4   |   Find the best splitting rule  $s_v$  for node  $v$ .
5   |   Create successors  $v_T$  and  $v_F$  of  $v$ .
6   |    $\sigma_T \leftarrow \{(x, y) \in \sigma : s_v(x) = \text{True}\}$ 
7   |    $\sigma_F \leftarrow \{(x, y) \in \sigma : s_v(x) = \text{False}\}$ 
8   |    $\mathbb{T}_{v_T} \leftarrow \text{Construct\_Subtree}(v_T, \sigma_T)$            // left branch
9   |    $\mathbb{T}_{v_F} \leftarrow \text{Construct\_Subtree}(v_F, \sigma_F)$        // right branch
10 return  $\mathbb{T}_v$ 

```

---

Needed tools:

- Regional prediction functions  $g^w$
- Specification of the splitting rule
- Termination criterion

# Regional prediction functions $g^w$

Any function may be used for  $g^w$  of node in a leaf  $v = w$ , but, typically:

**Classification**  $g^w$  is usually constant and equal to the most common class label of  $\tau$  in the associated region  $\mathcal{R}_w$

**Regression**  $g^w$  taken as the mean response in the region.

# Splitting rules

In the algorithm in slide 12, we split into two subsets. For any splitting rule, the contribution to the loss is always greater if we stop the splitting at some point (see Exercise ??). So, it makes sense to *minimize* the training loss.

**ession** the squared-error loss of binary splitting into two child nodes  $v_T$  and  $v_F$  is:

$$\frac{1}{n} \sum_{(\mathbf{x}, y) \in \sigma_T} 1\{y \neq y_T^*\} + \frac{1}{n} \sum_{(\mathbf{x}, y) \in \sigma_F} 1\{y \neq y_F^*\} \quad (1)$$

**cation** the loss can be obtained from

$$\frac{1}{n} \sum_{(\mathbf{x}, y) \in \tau; x_j \leq \xi} (y - \bar{y}_T)^2 + \frac{1}{n} \sum_{(\mathbf{x}, y) \in \tau; x_j > \xi} (y - \bar{y}_F)^2 \quad (2)$$

To minimize the expression ?? we need to evaluate over all  $j$  and  $\xi$ , for each of the  $m \times p$  values  $x_{j,k}$  and then take the minimizing pair  $(j, x_{j,k})$ . For 2 we can follow a similar approach.

# Regional prediction functions $g^w$

## Exercise 2

Splitting Can you rationalize how to generate a termination criteria based on the previous fact?

# Measuring impurity

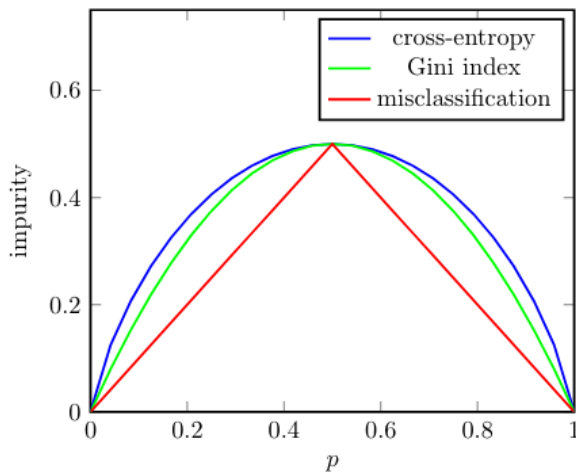
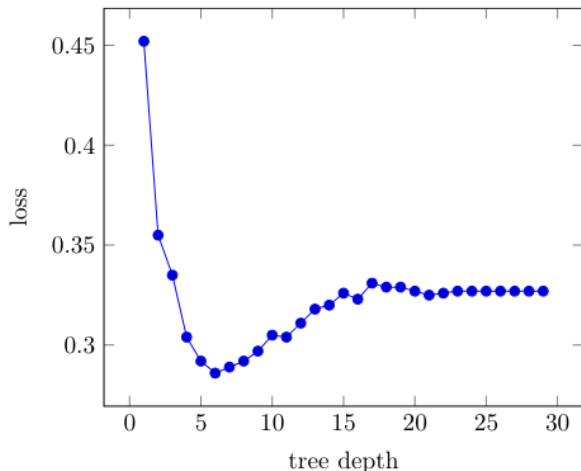


Figure 2: Entropy(normalized, divided by 2), Gini and missclassification impurities in binary classification (with probabilities  $p$  and  $1 - p$ ). They are maximal when



# Termination criterion

We can use a simple approach (number of data points in a tree node, for example) or when there is no advantage with more splits:



- Binary (preferred) vs non binary trees. Too few data points in nodes close to the tree root.
- Alternative splitting rules (see Slide 19). Perhaps use several variables?
- Categorical variables ( $2^k$  subsets for optimal partition with  $k$  different categorical variables).
- Missing values (add a category of *missing values*)

# Alternative splitting rules

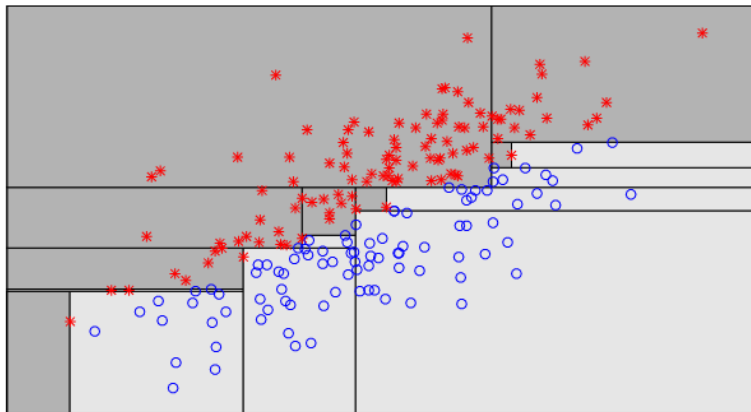
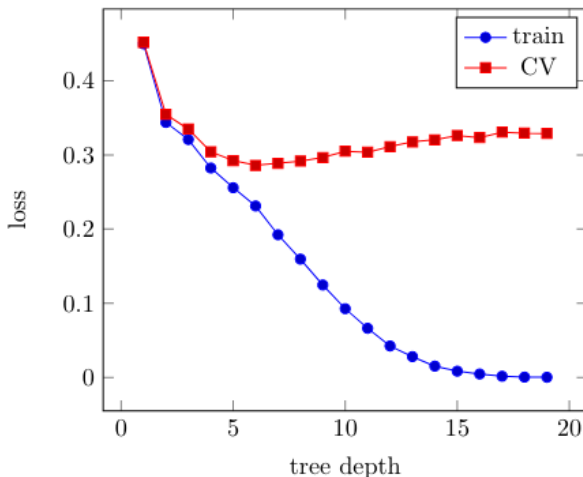
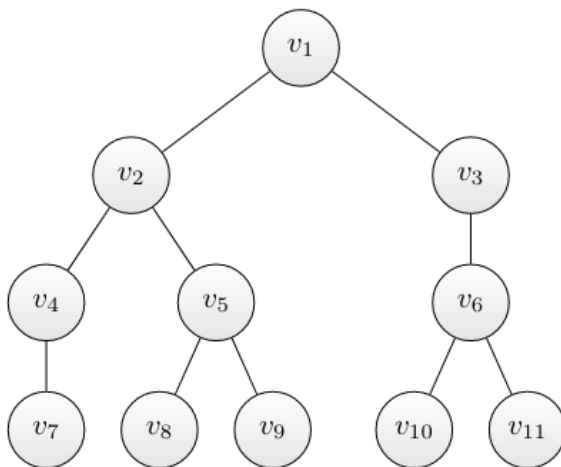


Figure 4: The two groups are separated, in decision tree, with a collection of rectangles, leading to an unnecessary classification procedure some times.[1]



**Figure 5:** Shallow trees tend to underfit and deep trees to overfit. Cross-validation and training loss as a function of the tree depth for a binary classification problem.[1]. Relate this plot with Exercies 1

# Pruning a deep tree



## Pruning a deep tree

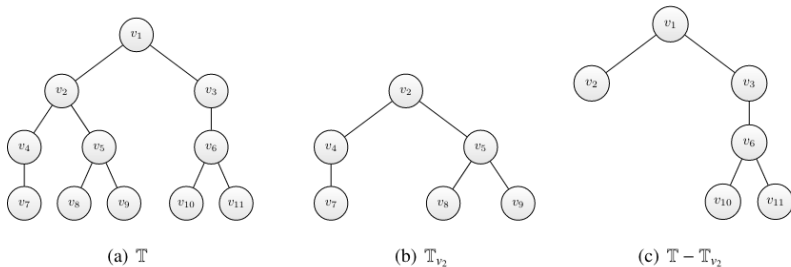


Figure 7: Pruning branches of the tree.[1]

**Algorithm 8.5.1:** Bootstrap Aggregation Sampling**Input:** Training set  $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and resample size  $B$ .**Output:** Bootstrapped data sets.

```

1 for  $b = 1$  to  $B$  do
2    $\mathcal{T}_b^* \leftarrow \emptyset$ 
3   for  $i = 1$  to  $n$  do
4     Draw  $U \sim \mathcal{U}(0, 1)$ 
5      $I \leftarrow \lceil nU \rceil$  // select random index
6      $\mathcal{T}_b^* \leftarrow \mathcal{T}_b^* \cup \{(\mathbf{x}_I, y_I)\}$ .
7 return  $\mathcal{T}_b^*, b = 1, \dots, B$ .
```

**Algorithm 8.5.2:** Out-of-Bag Loss Estimation

**Input:** The original data set  $\tau = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , the bootstrapped data sets  $\{\mathcal{T}_1^*, \dots, \mathcal{T}_B^*\}$ , and the trained predictors  $\{g_{\mathcal{T}_1^*}, \dots, g_{\mathcal{T}_B^*}\}$ .

**Output:** Out-of-bag loss for the averaged model.

```

1 for  $i = 1$  to  $n$  do
2    $C_i \leftarrow \emptyset$  // Indices of predictors not depending on  $(\mathbf{x}_i, y_i)$ 
3   for  $b = 1$  to  $B$  do
4     if  $(\mathbf{x}_i, y_i) \notin \mathcal{T}_b^*$  then  $C_i \leftarrow C_i \cup \{b\}$ 
5      $Y'_i \leftarrow |C_i|^{-1} \sum_{b \in C_i} g_{\mathcal{T}_b^*}(\mathbf{x}_i)$ 
6      $L_i \leftarrow \text{Loss}(y_i, Y'_i)$ 
7  $L_{\text{OOB}} \leftarrow \frac{1}{n} \sum_{i=1}^n L_i$ 
8 return  $L_{\text{OOB}}$ .
```



We build a number of trees, including a subset of features in each tree construction. This way, strong predictors will have a smaller chance to be considered at the root levels.

---

### Algorithm 8.6.1: Random Forest Construction

---

**Input:** Training set  $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of trees in the forest  $B$ , and the number  $m \leq p$  of features to be included, where  $p$  is the total number of features in  $\mathbf{x}$ .

**Output:** Ensemble of trees.

- 1 Generate bootstrapped training sets  $\{\mathcal{T}_1^*, \dots, \mathcal{T}_B^*\}$  via Algorithm 8.5.1.
  - 2 **for**  $b = 1$  **to**  $B$  **do**
  - 3     Train a decision tree  $g_{\mathcal{T}_b^*}$  via Algorithm 8.2.1, where each split is performed using  $m$  randomly selected features out of  $p$ .
  - 4 **return**  $\{g_{\mathcal{T}_b^*}\}_{b=1}^B$ .
-

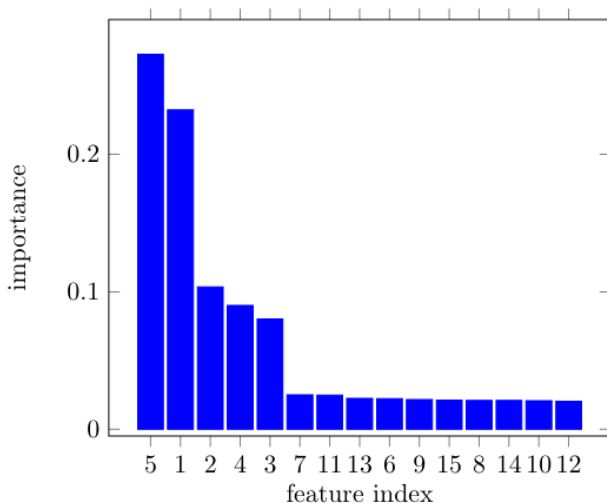


Figure 8: Importance measure for a 15-feature data set with only informative features  $x_1, x_2, x_3, x_4, x_5$  [1].

Similar to bootstrapping but now the predicting functions are learned sequentially, and not randomly.

---

**Algorithm 8.7.1:** Regression Boosting with Squared-Error Loss

---

**Input:** Training set  $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of boosting rounds  $B$ , and a shrinkage step-size parameter  $\gamma$ .

**Output:** Boosted prediction function.

- 1 Set  $g_0(\mathbf{x}) \leftarrow n^{-1} \sum_{i=1}^n y_i$ .
  - 2 **for**  $b = 1$  **to**  $B$  **do**
  - 3     Set  $e_i^{(b)} \leftarrow y_i - g_{b-1}(\mathbf{x}_i)$  for  $i = 1, \dots, n$ , and let  $\tau_b \leftarrow \{(\mathbf{x}_i, e_i^{(b)})\}_{i=1}^n$ .
  - 4     Fit a prediction function  $h_b$  on the training data  $\tau_b$ .
  - 5     Set  $g_b(\mathbf{x}) \leftarrow g_{b-1}(\mathbf{x}) + \gamma h_b(\mathbf{x})$ .
  - 6 **return**  $g_B$ .
- 

Figure 9:  $\gamma$  controls the speed of the fitting process. For small values, boosting takes smaller steps to the training loss minimization.

# Effect of $\gamma$

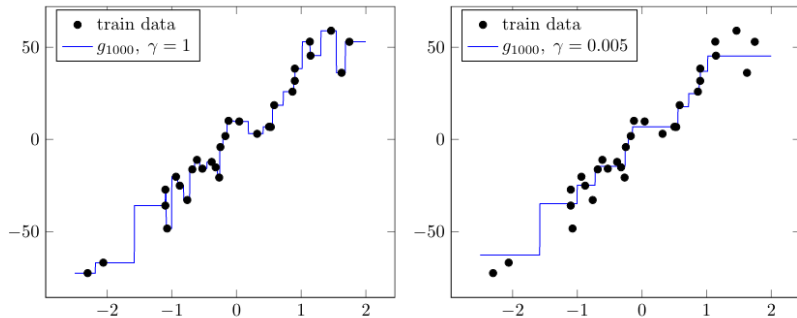


Figure 10: Fitted boosting regression model with two different values of  $\gamma$ . [1]

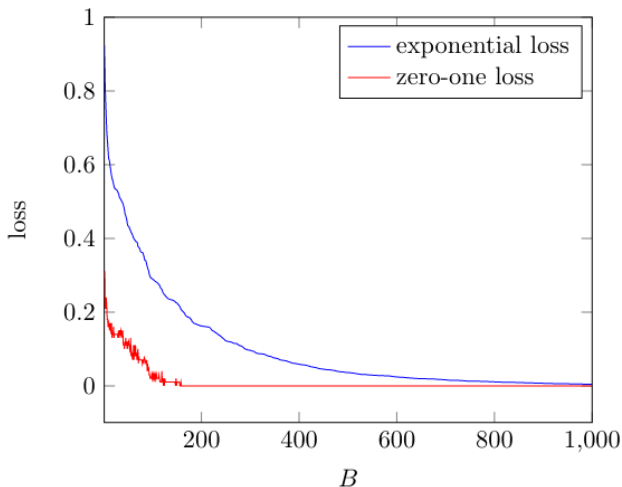


Figure 11: Exponential and zero-one training loss as a function of the number of boosting rounds  $B$  for a binary classification problem. [1]



Dirk P. Kroese, Zdravko Botev, Thomas Taimre, and Radislav Vaisman.

*Data Science and Machine Learning: Mathematical and Statistical Methods.*

Machine Learning & Pattern Recognition. Chapman & Hall/CRC, 2020.