

TD 1: Decision Trees (corrected version)

Ensemble learning from theory to practice

1 Exercises and reminders

Exercise 1: Reminders (to do on paper):

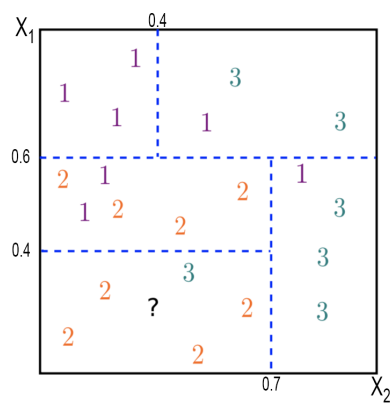


Figure 1: The spatial segmentation of a classification decision tree.

- Thanks to the spatial segmentation, draw the visual tree.

Correction:

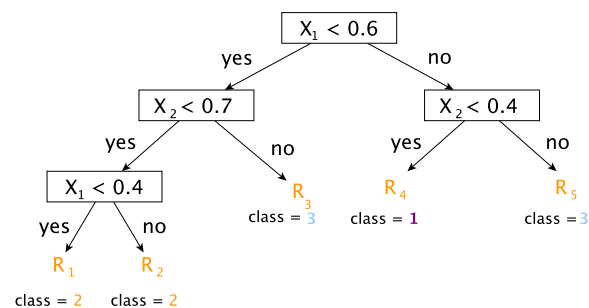


Figure 2: The visualized tree of the classification decision tree spatial segmentation.

2. What is the input space? What is the output space?

Correction:

- (a) The input space is $\mathcal{X} = \mathbb{R} \times \mathbb{R} = \mathbb{R}^2$ and another accepted response is $\mathcal{X} = [0, 1] \times [0, 1] = [0, 1]^2$.
 (b) The output space is $\mathcal{Y} = \{1, 2, 3\}$

3. Write the f formula associated.

Correction: Let's define x_i an observation (d -length vector) from the input space and $x_{i,j}$ a real observation associated to the variable X_j .

$$\begin{aligned} f(x_i) &= \mathbb{1}_{\{x_i \in R_4\}} + 2\mathbb{1}_{\{x_i \in R_1 \cup R_2\}} + 3\mathbb{1}_{\{x_i \in R_3 \cup R_5\}} \\ &= \sum_{k=1}^K \gamma_k \mathbb{1}_{\{x_i \in R_k\}} \end{aligned}$$

where,

$$\begin{aligned} \gamma_1 &= \gamma_2 = 2 \\ \gamma_3 &= \gamma_5 = 3 \\ \gamma_4 &= 1 \\ R_1 &= \{x_i : x_{i,1} < 0.4\} \cap \{x_i : x_{i,2} < 0.7\} \\ R_2 &= \{x_i : 0.4 < x_{i,1} < 0.6\} \cap \{x_i : x_{i,2} < 0.7\} \\ R_3 &= \{x_i : x_{i,1} < 0.6\} \cap \{x_i : x_{i,2} > 0.7\} \\ R_4 &= \{x_i : x_{i,1} > 0.6\} \cap \{x_i : x_{i,2} < 0.4\} \\ R_5 &= \{x_i : x_{i,1} > 0.6\} \cap \{x_i : x_{i,2} > 0.4\} \end{aligned}$$

4. How many split there are?

Correction: 4 splits.

5. How many nodes?

Correction: 9 nodes including 5 terminal nodes.

6. How many leaves?

Correction: 5 leaves.

7. What will be the predicted value associated to the new data symbolized by the question mark?

Correction: class 2 (or orange).

8. Is this tree accurate?

Correction: It looks accurate (with a little bit classification error which is preferable to overfitting) but we can have an interrogation about the relevance of the last split.

9. Compute the MSE (Mean Squared Error measure) to evaluate the quality of the decision tree predictor

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Correction: $MSE = \frac{11}{21} \approx 0.52$

10. What do you think about the result? What is disturbing in the previous question (clue: train data vs test data)?

Correction: First we can expect a better result (closest to zero). Then here we compute the error on the train set however to conclude about the accuracy of a decision tree it is better appropriate to consider the error on the test set. At last, here we are in a classification case therefore we have to use an error adapted (e.g. missclassification error). Indeed, MSE is only adapted for the regression task.

Exercise 2 (Questions about understanding basic concepts of the course):

1. A student has trained a decision tree and notices that its performance is better on the train set than the test set. Should he increase or decrease the depth of the tree?

Correction: He should decrease the depth of the tree because reducing the depth will give a model that is less prone to overfitting.

2. A student has dataset with n instances and p features.

- (a) What is the maximum number of leaves of a decision tree on this dataset?

Correction: Each leaf contains at least one observation, therefore the answer is n .

- (b) What is the maximum depth of a decision tree on this dataset?

Correction: Each leaf contains at least one observation, therefore the answer is $n - 1$. If on the path between the root and the deepest leaf, at each split a leaf is created which contains a single observation. Note that the root has a depth equal to zero.

Exercise 3 (Understand the splitting process idea by practice with Gini impurity):

Consider the following dataset containing for 10 plants the length and width of their sepals. We want to discriminate plants that belong to the species Iris virginica (+) from others (-).

Label	+	+	+	+	+	+	-	-	-	-
Length (cm)	6.7	6.7	6.3	6.5	6.2	5.9	6.1	6.4	6.6	6.8
Width(cm)	3.3	3	2.5	3	3.4	3	2.8	2.9	3	2.8

1. Calculate the Gini impurity for all possible separation points using the **length** of the sepals as the separating variable.

Note that, the Gini impurity of an R region is defined as:

$$Imp(R) := \sum_{c=1}^C p_c(R)(1 - p_c(R)) \quad (1)$$

where, $p_c(R) := \frac{1}{|R|} \sum_{i: \mathbf{x}_i \in R} \mathbb{1}_{\{y_i=c\}}$. Thus, if all the instances of a region belong to the same class, the impurity of this region is equal to 0; conversely, if a region contains as many instances of each of the C classes, the right part of the product is $1 - p_c(R) = 1 - \frac{1}{C}$, or $\frac{1}{2}$ in the case of a binary classification.

2. Calculate the Gini impurity for all possible separation points using the **width** of the sepals as the separating variable.

3. What is the first node of a decision tree trained on this dataset with the Gini impurity?

Exercise 4: Understand the splitting process idea (to do on paper): We will show that minimizing the quadratic risk amounts to minimizing the variance in each hyper-rectangle of the input space partition.

1. For all region R_k , write the minimization empirical quadratic risk problem where the predictor is a decision tree.
2. Write R_k in function of 2 subspaces $R_L(j, s)$ and $R_R(j, s)$ obtained after splitting a region R_k via a split (j, s) .

-
3. How are these two subspaces relative to each other?
 4. Write the new risk formula based on these two subspaces.
 5. Write the variance formula in a node k (corresponding to a region R_k), then in each child nodes of k .
 6. Show that minimizing the empirical quadratic risk formula amounts to minimizing the variance in each hyper-rectangle of the input space partition.