

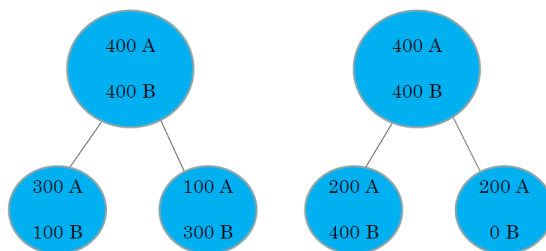
Universidad de la República
Facultad de Ingeniería

PRÁCTICO 3 : CART - MÉTODOS DE AGREGACIÓN

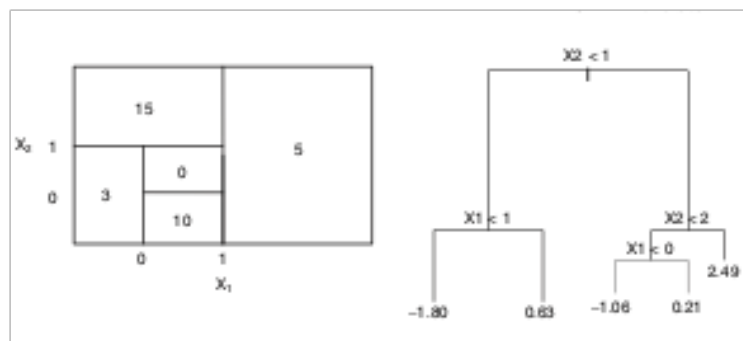
1. Prove that if X is categorical with m levels that there is $2^m - 1$ possible splits.
2. Prove that the three expressions of Gini index are the same.
3. Consider the following sample:

x_1	a	a	b	a	a	b	b	b
x_2	b	a	a	a	a	b	b	b
y	1	1	1	1	-1	-1	-1	-1

- a) Construct a classification tree, with the Gini index, from this sample
 - b) Compare with the tree obtained in R. What argument of the rpart function should you change to get the same tree?
4. Compute $\Delta i(t, s)$ for these two partitions using classification error, Gini index and entropy



5. Let consider the following figure:



- a) Sketch the tree corresponding to the partition of the predictor space illustrated in the left-hand panel of the figure.

- b) Create a diagram similar to the left-hand panel of the figure, using the tree illustrated in the right-hand panel of the same figure.
6. Suppose we produce 3 bootstrapped samples from a set containing black and white classes. We then apply a classification tree to each bootstrapped sample and, for a new value \mathbf{x}_0 produce 3 estimates of $\mathbb{P}(\text{black}|\mathbf{x}_0)$: 0.1, 0.1 and 0.9. What is the final prediction under majority vote approach? What is the final prediction if we average the probabilities?
7. From the dataset `spam` of the `kernlab` library, split it into train (2/3) and test sample (1/3) with `set.seed(2021)`.
- Compute and draw the default tree T provided by `rpart` and the decision stump. Look at `T$frame` and examine it.
 - 1) Compute and draw the optimal tree T_1 with associate `cp` parameter given by cross-validation error.
2) Compute and draw the optimal tree T_2 with associate `cp` parameter given by the 1-SE rule.
3) Compare T, T_{max}, T_1 and T_2 in learning and in test samples.
 - Apply Bagging and Random Forest (default) and compare the prediction errors with a single tree.
 - Study the evolution of the OOB error with respect to `ntree` using `do.trace`.
 - Calculate the variable importance of the spam variables for Random Forest (default).
 - Calculate the importance of spam variables for stumps Random Forest.
 - Illustrate the influence of the `mtry` parameter on the OOB error and on the variable importance.
8. Use the `Carseats` dataset of the `tree` library.
- Split the data set into a training set and a test set (2/3 -1/3 with `set.seed(2021)`).
 - Fit a regression tree to the training set. Plot the tree, and interpret the results. What test error rate do you obtain?
 - Use cross-validation in order to determine the optimal level of tree complexity. Does pruning the tree improve the test error rate?
 - Use the Bagging approach in order to analyze this data. What test error rate do you obtain?
 - Use Random Forests to analyze this data. What test error rate do you obtain? Use the `importance()` function to determine which variables are most important. Describe the effect of m , the number of variables considered at each split, on the error rate obtained.
 - Answer the same questions if the variable `Sales` is discretized as follows: 1 if the `Sales` variable is higher than 8, 0 otherwise.