

Decision Trees

Pietro Gori

Maître de conférences
Equipe IMAGES - Télécom ParisTech
pietro.gori@telecom-paristech.fr

March 02, 2018



Une école de l'IMT

1 Introduction

- Reminder about classification
- Constant piecewise model/classifier

2 Decision trees

- Cost functions
- Impurity function
- Stopping criteria and variations
- Model selection

Supervised Learning - Probabilistic framework

X : input data x_i^j , random variable in $\mathcal{X} = \mathbb{R}^p$ with $i = 1, \dots, n$ and $j = 1, \dots, p$ where n and p are the number of observations and variables respectively

Y : response (to predict), random variable in $\mathcal{Y} = \{C_1, \dots, C_K\}$ (classification with K classes) or $\mathcal{Y} = \mathbb{R}$ (regression)

P : joint probability distribution of (X, Y) , fixed but unknown

$\mathcal{D}_n = \{(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}$: i.i.d. samples drawn from P

\mathcal{F} : collection of classifiers $f \in \mathcal{F}$

\mathcal{L} : loss function which measures the error of the classifier/model

- Examples (classification) : $\mathcal{L}(\mathbf{x}, y, f(\mathbf{x})) = \begin{cases} 1, & \text{si } f(\mathbf{x}) \neq y, \\ 0, & \text{sinon.} \end{cases}$
- Example (regression): $\mathcal{L}(\mathbf{x}, y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$

Goal: estimate from \mathcal{D}_n the function $f \in \mathcal{F}$ which minimizes the risk (cost) function $R(f) = \mathbb{E}_P[\mathcal{L}(X, Y, f(X))]$

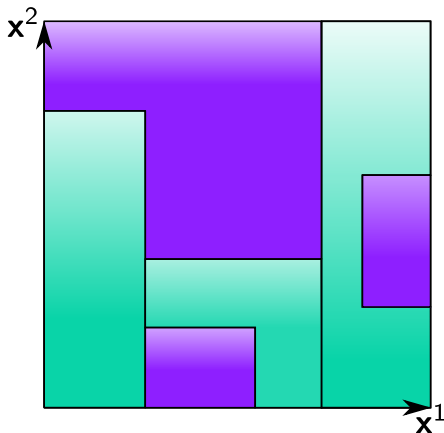
Estimate a classifier

We need to define:

- **input and output data space** ($\mathcal{X} \mathcal{Y}$)
- **type of classifier** (\mathcal{F})
- **cost function** (\mathcal{L}) to minimize for finding the best f
- **minimization algorithm** for \mathcal{L}
- **method for model selection** to estimate hyper-parameters
- **method to evaluate performance**

Constant piecewise model/classifier

\mathcal{F} belongs to the set of **constant piecewise functions**. We divide the input space \mathcal{X} in M disjoint partitions \mathcal{C}_m : $\mathcal{X} = \mathcal{C}_1 \sqcup \dots \sqcup \mathcal{C}_M$. To simplify things, we assume that the separation lines are parallel to the coordinate axes:



Constant piecewise model/classifier

Let $\mathbf{x} = (\mathbf{x}^1, \dots, \mathbf{x}^p)$, we model the function f with a different constant value α_m in every partition \mathcal{C}_m : $f(\mathbf{x}) = \sum_{m=1}^M \alpha_m \mathbb{1}_{\mathbf{x} \in \mathcal{C}_m}$, where $\mathbb{1}_{\mathbf{x} \in \mathcal{C}_m}$ is equal to 1 if $\mathbf{x} \in \mathcal{C}_m$ and 0 otherwise.

- **Regression:** If we use the L2-norm ($\sum_i (y_i - f(\mathbf{x}_i))^2$), then the best α_m^* is simply the average of the y_i within the region \mathcal{C}_m :
$$\alpha_m^* = \frac{1}{|\mathcal{C}_m|} \sum_{\mathbf{x}_i \in \mathcal{C}_m} y_i, \text{ where } |\mathcal{C}_m| \text{ is the number of elements within } \mathcal{C}_m$$
- **Classification:** We define the proportion of observations belonging to class k within region \mathcal{C}_m as $\rho_{mk} = \frac{1}{|\mathcal{C}_m|} \sum_{\mathbf{x}_i \in \mathcal{C}_m} \mathbb{1}_{y_i=k}$. The best α_m^* is then equal to the majority class: $\alpha_m^* = \arg \max_k \rho_{mk}$

- Motivation: easy to interpret
- Limitations:
 - regions are difficult to describe
 - If the partition is fixed beforehand, many regions might end up being empty
- Possible solution: learn the partitions from the data ! How to avoid curse of dimensionality ?

1 Introduction

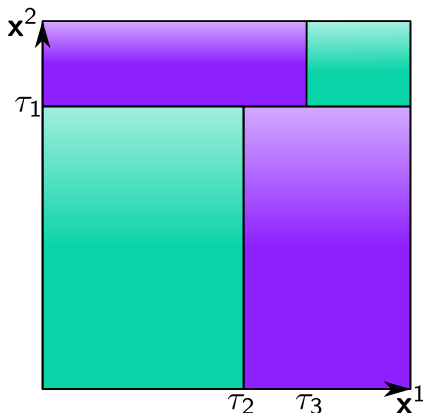
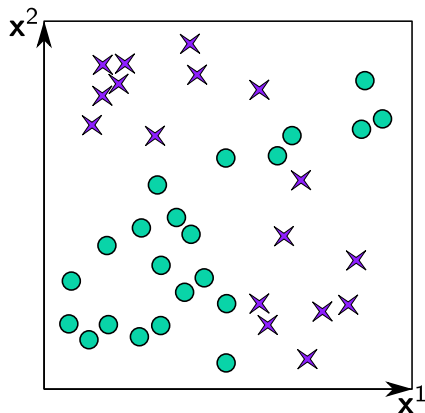
- Reminder about classification
- Constant piecewise model/classifier

2 Decision trees

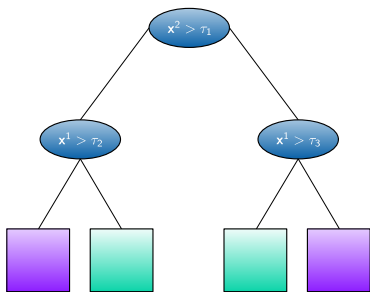
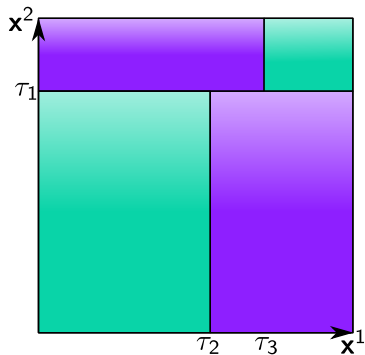
- Cost functions
- Impurity function
- Stopping criteria and variations
- Model selection

Decision trees

Presented by Breiman *et al.* in 1984 under the name of CART:
Classification and Regression Trees.



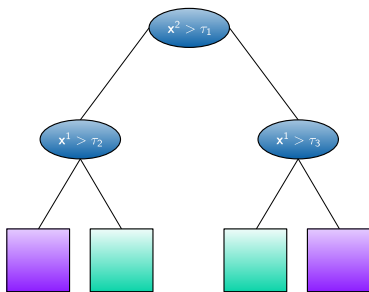
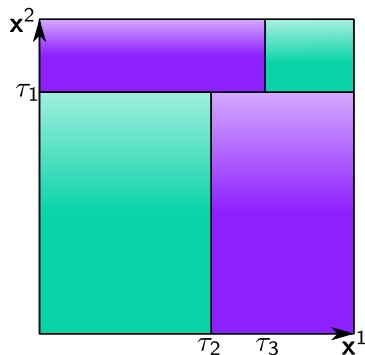
Decision trees



First idea:

Use several hyperplanes (and not only one) to build non linear decision boundaries.

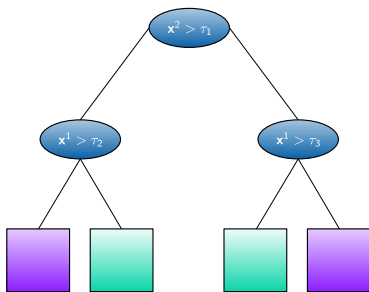
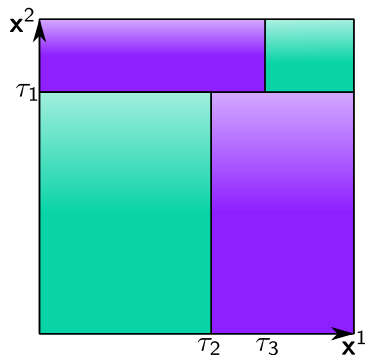
Decision trees



Secod idea:

Use separation lines orthogonal to the coordinate axis, *i.e.*, hyperplanes $\{\mathbf{x} \in \mathcal{X} : \mathbf{x}^j = \tau\}$ to ease interpretation.

Decision trees



Third idea:

Use recursive binary decision trees: The full data-set sits at the top of the tree. Every node (junction) is associated to a separating hyperplane $\{\mathbf{x} \in \mathcal{X} : \mathbf{x}^j = \tau\}$. The terminal nodes correspond to the regions.

Linear partition orthogonal to the axes

- We use the same model as before: $f(\mathbf{x}) = \sum_{m=1}^M \alpha_m \mathbb{1}_{\mathbf{x} \in \mathcal{C}_m}$, namely we assign a constant value at every region
- Let $\mathbf{x} = (\mathbf{x}^1, \dots, \mathbf{x}^p)$ with p variables. We define the split $t_{j,\tau}(\mathbf{x})$ along the direction \mathbf{x}^j with threshold τ as:

$$t_{j,\tau}(\mathbf{x}) = \text{sign}(\mathbf{x}^j - \tau) = \begin{cases} +1, & \text{si } \mathbf{x}^j > \tau \\ -1, & \text{si } \mathbf{x}^j < \tau \end{cases} \quad (1)$$

Efficient recursive algorithm

For a binary tree:

- 1 Given \mathcal{D}_n , the entire data-set is the root node
- 2 Look for the best separator $t_{j,\tau}$ on \mathcal{D}_n such that the local cost function $\mathcal{L}(t_{j,\tau}, \mathcal{D}_n)$ is minimal. This means looking for the “best” direction j and threshold τ that splits \mathcal{D}_n into \mathcal{D}_n^d and \mathcal{D}_n^g .
- 3 Be careful: the splitting values τ are not infinite ! They depend on the data \mathcal{D}_n . Hence, we can scan all inputs \mathbf{x} and quickly find the best j and τ
- 4 Split \mathcal{D}_n to \mathcal{D}_n^d and \mathcal{D}_n^g using the estimated separator hyperplane.
- 5 It results two nodes, a left (\mathcal{D}_n^g) and a right (\mathcal{D}_n^d) one
- 6 Evaluate the stopping criteria for the right node, if it is verified, the nodes becomes a terminal node, otherwise go to 3 using \mathcal{D}_n^d as input space
- 7 Evaluate the stopping criteria for the left node, if it is verified, the nodes becomes a terminal node, otherwise go to 3 using \mathcal{D}_n^g as input space

Examples

Given the input data-set \mathcal{D}_n and a binary separator $t_{j,\tau}$, we have

$$\mathcal{D}_n^d(j, \tau) = \{(\mathbf{x}, y) \in \mathcal{D}_n, t_{j,\tau}(\mathbf{x}) > 0\}$$

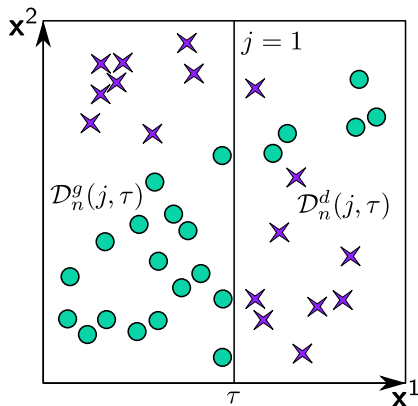
$$\mathcal{D}_n^g(j, \tau) = \{(\mathbf{x}, y) \in \mathcal{D}_n, t_{j,\tau}(\mathbf{x}) \leq 0\}$$

Examples

Given the input data-set \mathcal{D}_n and a binary separator $t_{j,\tau}$, we have

$$\mathcal{D}_n^d(j, \tau) = \{(\mathbf{x}, y) \in \mathcal{D}_n, t_{j,\tau}(\mathbf{x}) > 0\}$$

$$\mathcal{D}_n^g(j, \tau) = \{(\mathbf{x}, y) \in \mathcal{D}_n, t_{j,\tau}(\mathbf{x}) \leq 0\}$$

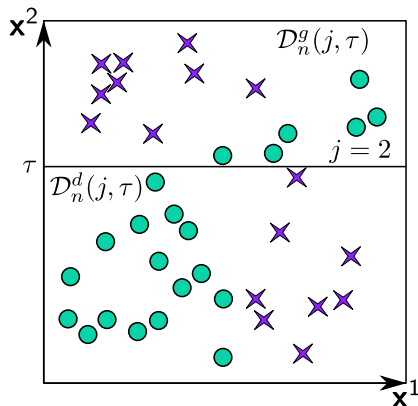
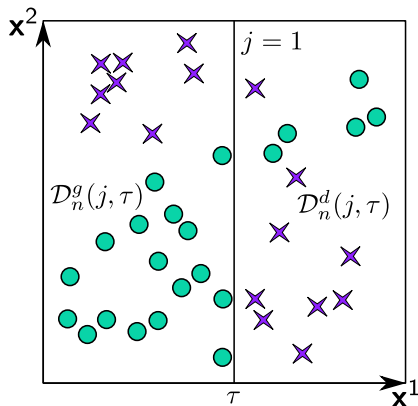


Examples

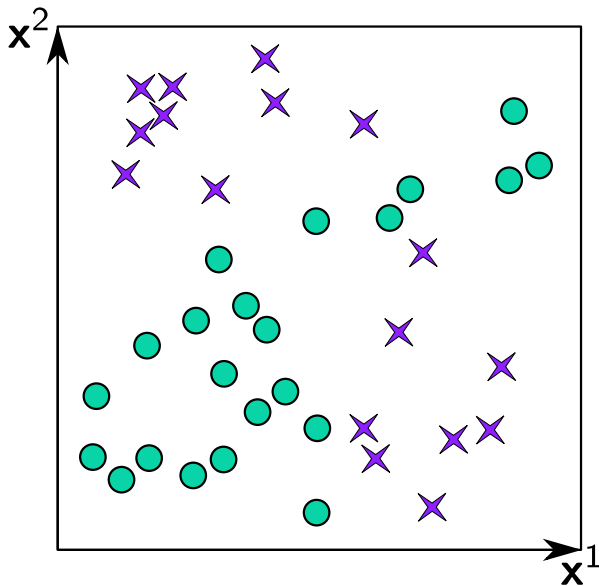
Given the input data-set \mathcal{D}_n and a binary separator $t_{j,\tau}$, we have

$$\mathcal{D}_n^d(j, \tau) = \{(\mathbf{x}, y) \in \mathcal{D}_n, t_{j,\tau}(\mathbf{x}) > 0\}$$

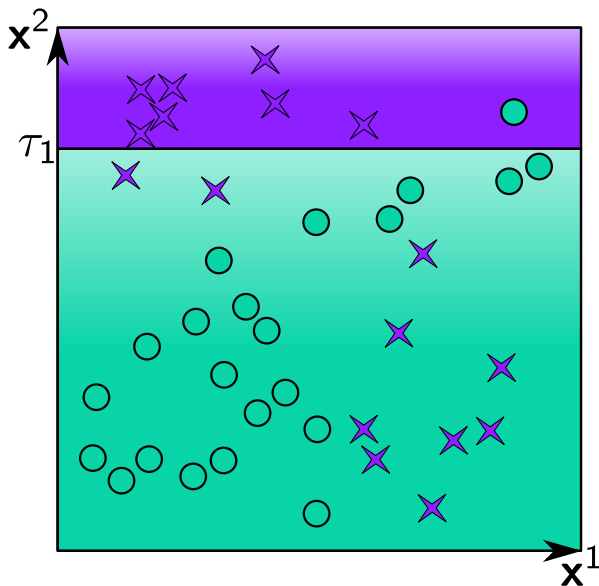
$$\mathcal{D}_n^g(j, \tau) = \{(\mathbf{x}, y) \in \mathcal{D}_n, t_{j,\tau}(\mathbf{x}) \leq 0\}$$



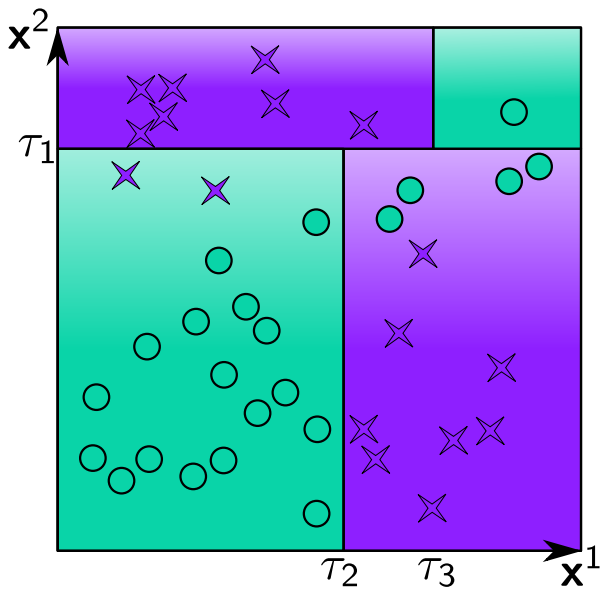
Example



Example



Example



- The presented algorithm (called CART) is greedy
- We do not optimize a global criteria. Instead, we locally look for an optimal separator (with respect to L), which means at every direction j independently.
- Why in your opinion ?

- We recall that for a classification tree, given the input data \mathcal{D}_n divided in K classes, we define the proportion of observations belonging to class k as: $\rho_k(\mathcal{D}_n) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i = k)$
- Note that $\rho(\mathcal{D}_n) = (\rho_1(\mathcal{D}_n), \dots, \rho_K(\mathcal{D}_n))^\top \in \Delta_{K-1}$ where
$$\Delta_{K-1} := \left\{ \rho_k \in \mathbb{R}^K : \sum_{k=1}^K \rho_k = 1 \text{ and } \forall k \in \llbracket 1, K \rrbracket, \rho_k \geq 0 \right\}$$
 is the $(K-1)$ -simplex. A simplex can be seen as the smallest convex set containing the given vertices, or the convex hull of the K vertices.
Example: a 2-simplex is a triangle.

Cost function - Classification

Among all parameters $(j, \tau) \in \{1, \dots, p\} \times \{\tau_1, \dots, \tau_m\}$, we look for j^* and τ^* which minimizes the following cost function:

$$L(t_{j,\tau}, \mathcal{D}_n) = \frac{n_g}{n} H(\rho(\mathcal{D}_n^g(j, \tau))) + \frac{n_d}{n} H(\rho(\mathcal{D}_n^d(j, \tau)))$$

avec $n_g = |\mathcal{D}_n^g(j, \tau)|$ et $n_d = |\mathcal{D}_n^d(j, \tau)|$

- H is an “impurity” function that evaluate the splitting. Pure means a node with observations from the same class. We want to **minimize** H in order to have pure nodes.
- The total cost is the sum of the impurity of each child node (\mathcal{D}_n^g and \mathcal{D}_n^d) weighted by the proportion of its observations (n_g and n_d)
- We evaluate a finite number of thresholds (max n)

Definition: impurity function

An impurity function $H : [0 : 1]^K \rightarrow \mathbb{R}$ is a function defined on Δ_{K-1} for which the following properties hold:

- 1 H becomes maximum at points $(\frac{1}{K}, \dots, \frac{1}{K})^\top$, i.e., all ρ_k are equal
- 2 H becomes minimum at points $(1, 0, \dots, 0)^\top, (0, 1, 0, \dots, 0)^\top, \dots, (0, \dots, 0, 1)^\top$, i.e., the probability of being in a certain class is 1 and 0 for all other classes. These are the vertices of Δ_{K-1} .
- 3 H is symmetric with respect to its arguments ρ_1, \dots, ρ_K , i.e., even if we permute ρ_j , H does not change

Impurity function: binary case ($K = 2$)

When we have only two classes ($K = 2$):

- Δ_{K-1} is the line segment joining $(1, 0)$ and $(0, 1)$ in \mathbb{R}^2
- H becomes maximum at $(\frac{1}{2}, \frac{1}{2})$
- H becomes minimum at $(0, 1)$ or $(1, 0)$, which means when all observations in one region belong to the same class.

Misclassification error

Given the data of a node \mathcal{D}_n (it might be the root node or a child node), we assign the observations in \mathcal{D}_n to the majority class k^* :

$$k^* = \arg \max_{k=1,\dots,K} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_n} \mathbb{1}(y_i = k)$$

Then we define:

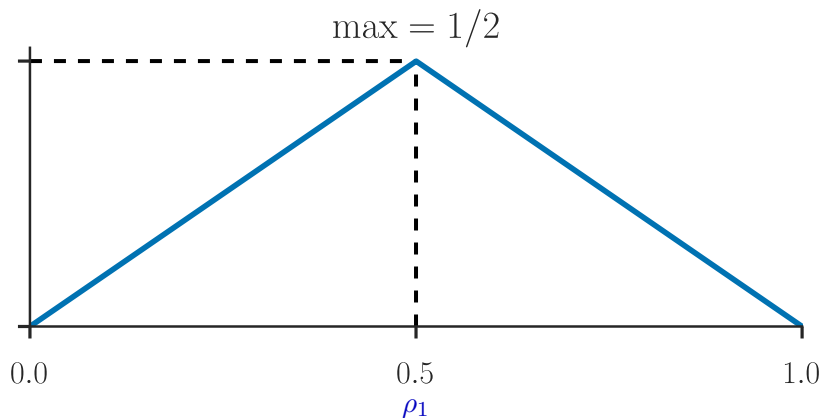
$$\text{Misclassification error: } H_{\text{mis}}(\mathcal{D}_n) = 1 - \rho_{k^*}(\mathcal{D}_n),$$

This is the error that we commit by assigning the observations to the class k^* . Remember that $\sum_{k=1}^K \rho_k = 1$.

Misclassification error

When the number of classes K is 2

$$H_{\text{mis}}(\mathcal{D}_n) = 1 - \max_{k=1,2} \rho_k(\mathcal{D}_n) = \min(\rho_1(\mathcal{D}_n), 1 - \rho_1(\mathcal{D}_n))$$

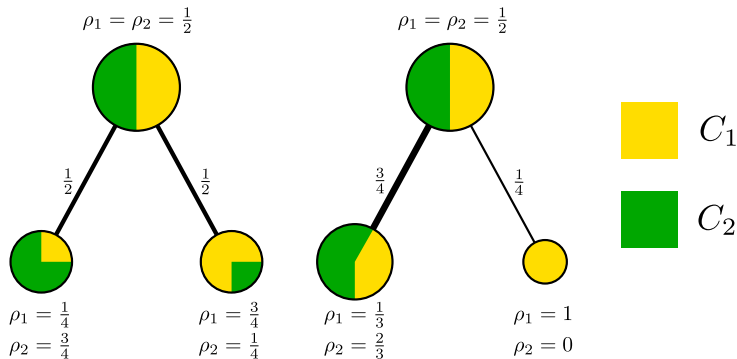


Limitations of the misclassification error

- Remember that the cost of a split is:

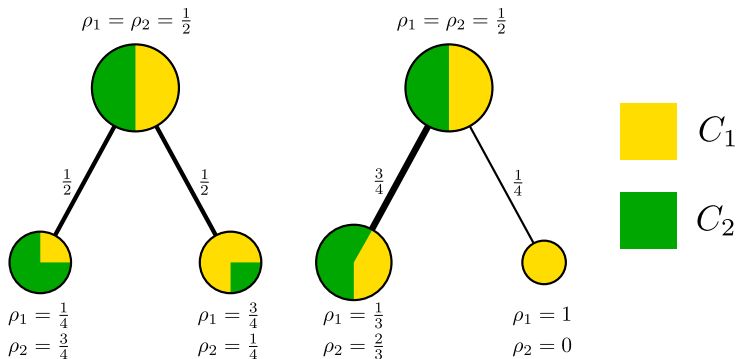
$$L(t_{j,\tau}, \mathcal{D}_n) = \frac{n_g}{n} H(\rho(\mathcal{D}_n^g(j, \tau))) + \frac{n_d}{n} H(\rho(\mathcal{D}_n^d(j, \tau)))$$

- For both splits we have $L_{\text{mis}} = \frac{1}{2} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{1}{4} = \frac{1}{4} = \frac{3}{4} \cdot \frac{1}{3} + \frac{1}{4} \cdot 0 = \frac{1}{4}$.
However it seems that the second one is definitely better !



Limitations of the misclassification error

- For a partition where a class has a clear majority, we might not find a split which reduces L
- The function is not differentiable (optimization is harder)
- It might underestimate pure nodes:



Strict impurity function

Definition: strict impurity function

Let $H : [0 : 1]^K \rightarrow \mathbb{R}$ be an impurity function, ρ, ρ' two distributions in Δ_{K-1} with $\rho \neq \rho'$ and $\alpha \in]0, 1[$. Then H is called strict, if it is strictly concave:

$$H(\alpha\rho + (1 - \alpha)\rho') > \alpha H(\rho) + (1 - \alpha)H(\rho')$$

If H is strict then it follows that

$$L(t_{j,\tau}, \mathcal{D}_n) = \frac{n_g}{n} H(\rho(\mathcal{D}_n^g(j, \tau))) + \frac{n_d}{n} H(\rho(\mathcal{D}_n^d(j, \tau))) \leq H(\rho(\mathcal{D}_n))$$
$$n_g = |\mathcal{D}_n^g(j, \tau)| \quad \text{et} \quad n_d = |\mathcal{D}_n^d(j, \tau)|$$

the equality is given iff $\rho_k(\mathcal{D}_n) = \rho_k(\mathcal{D}_n^g) = \rho_k(\mathcal{D}_n^d)$ for all k

Remark: The impurity function of the misclassification error is concave, but it is not strictly concave. L might be equal for all possible splittings.

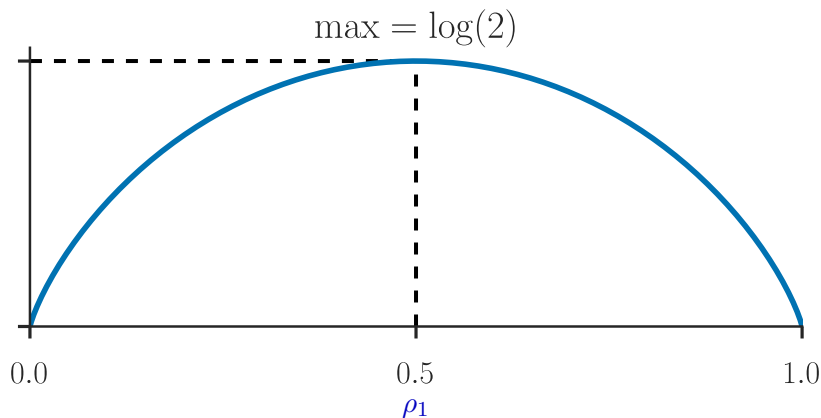
$$\text{Entropy: } H_{\text{ent}}(\mathcal{D}_n) = - \sum_{k=1}^K \rho_k(\mathcal{D}_n) \log \rho_k(\mathcal{D}_n)$$

- if we use \log_2 , it is called Shannon entropy
- $-\log \rho(\mathcal{D}_n)$ is the information content of \mathcal{D}_n
- Entropy is defined as the expected value of the information content (average amount of information). It measures the randomness
- when an event is certain, entropy is 0
- information gain is defined as reduction in entropy
- it is differentiable

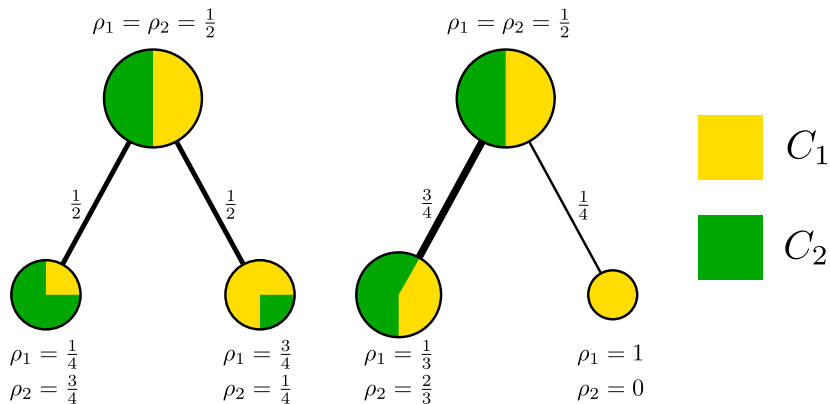
Entropy

When the number of classes K is 2

$$H_{\text{ent}}(\mathcal{D}_n) = -\rho_1(\mathcal{D}_n) \log(\rho_1(\mathcal{D}_n)) - (1 - \rho_1(\mathcal{D}_n)) \log(1 - \rho_1(\mathcal{D}_n))$$



Example



Question: Compute L_{ent} associated to H_{ent} . Which split is better ?

$$H_{\text{Gini}}(\mathcal{D}_n) = \sum_{k=1}^K \rho_k(\mathcal{D}_n)(1 - \rho_k(\mathcal{D}_n)) = \sum_{k=1}^K \sum_{\substack{k'=1 \\ k' \neq k}}^K \rho_k(\mathcal{D}_n)\rho_{k'}(\mathcal{D}_n)$$

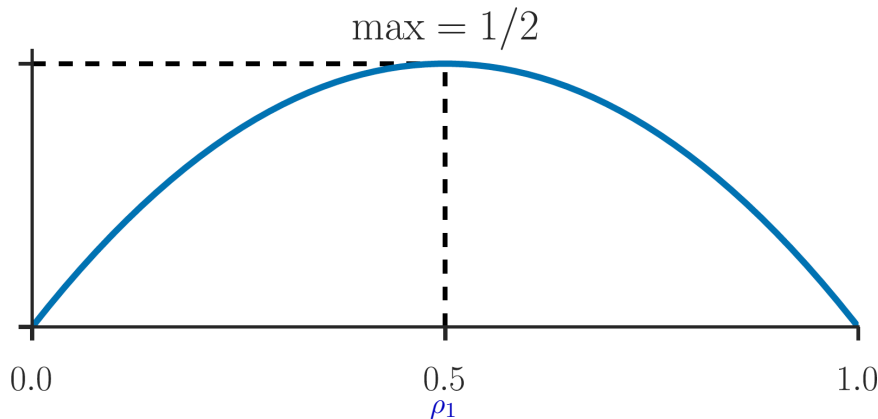
Two different interpretations:

- If we code each observation as 1 for class k and zero otherwise (bernoulli variable) the variance is $\rho_k(\mathcal{D}_n)(1 - \rho_k(\mathcal{D}_n))$. The Gini index is the sum of the variances of the “binarized” classes
- We do not assign observations to the majority class (as for H_{mis}) but we classify them to class k with probability $\rho_k(\mathcal{D}_n)$. The training error rate of this rule is $\sum_{\substack{k'=1 \\ k' \neq k}}^K \rho_k(\mathcal{D}_n)\rho_{k'}(\mathcal{D}_n)$. The Gini index is the sum over all classes.

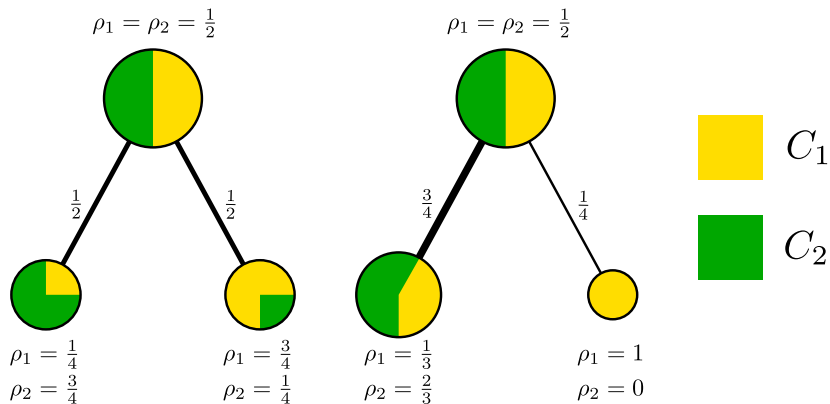
Gini index

When the number of classes K is 2

$$H_{\text{Gini}}(\mathcal{D}_n) = 2\rho_1(\mathcal{D}_n)(1 - \rho_1(\mathcal{D}_n))$$



Example



Question: Compute L_{Gini} associated to H_{Gini} . Which split is better ?

Without stopping criteria, we could grow a tree until a situation where each observation represents a terminal node. This would be computationally expensive, difficult to interpret and prone to over-fitting. Instead, we could use one or more of the following stopping criteria:

- maximal depth
- maximum number of terminal nodes
- a node becomes a terminal node when it reaches a maximum number of observations
- degree of purity of a node (*i.e.*, threshold on $\rho_k(\mathcal{D}_n)$)

- For a binary tree : if we have a categorical variable x which can take up to M values, we transform it into M binary variables
- Warning: The partitioning algorithm tends to favor categorical variables with many values since the number of possible partitions grows exponentially with M . This means that we have more choices to find a good partition. This can lead to over-fitting ! Try to avoid such variables.

Loss matrix

In some cases, the consequences of misclassifying observations can be very serious (*i.e.*, medicine). To account for that, we introduce a loss matrix $C \in \mathbb{R}^{K \times K}$, with $C_{k,k'}$ being the loss incurred for classifying observations of class k as belonging to class k'

$$C_{k,k'} = 0 \text{ si } k = k' \quad C_{k,k'} \geq 0 \text{ si } k \neq k'$$

We can then modify the Gini index as follows:

$$\text{Gini index: } \sum_{k=1}^K \sum_{\substack{k'=1 \\ k' \neq k}}^K C_{k,k'} \rho_k(\mathcal{D}_n) \rho_{k'}(\mathcal{D}_n)$$

Note: This works for $K > 2$ but it has no effect in the binary case (Why?). A different approach consists of weighting the observations of class k by $C_{k,k'}$. In a terminal node we classify the observations to $k' = \arg \min_k \sum_k C_{k,k'} \rho_k(\mathcal{D}_n)$

Regression trees

For regression the process is almost identical, we only change the impurity function. We use the squared error (or variance):

$$H(\mathcal{D}_n) = \frac{1}{|\mathcal{D}_n|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_n} (y_i - \bar{y}_n)^2$$

where

$$\bar{y}_n = \frac{1}{|\mathcal{D}_n|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_n} y_i$$

and we minimize as before

$$L(t_{j,\tau}, \mathcal{D}_n) = \frac{n_g}{n} H(\mathcal{D}_n^g(j, \tau)) + \frac{n_d}{n} H(\mathcal{D}_n^d(j, \tau))$$

Note: as before we want to maximize the homogeneity (purity) of the terminal nodes

Model selection (1)

We can compute one (or more) of the following hyper-parameters instead than fixing them as stopping criteria:

- maximal depth of the tree
- maximum number of terminal nodes
- maximum number of observations in a node to become a terminal node

→ we could use cross validation

Pruning (2)

What's the optimal size of a tree ? A large tree might overfit the data, while a small tree might not capture important structures.

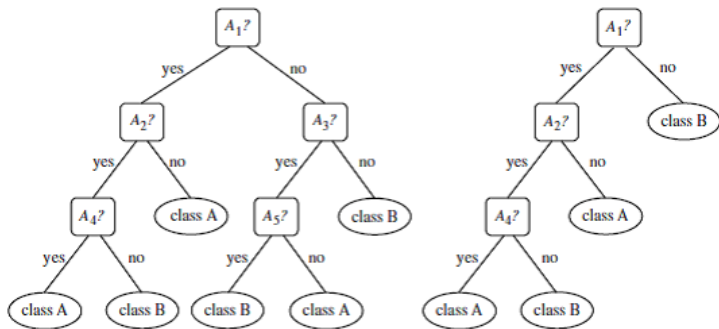
A possible solution is to grow a large tree in a training set stopping the splitting only when the terminal nodes have reached a minimum number of observations or a certain purity. Then, we produce several new trees by pruning the original tree at different nodes. When we want to prune with respect to node t , we delete all successor nodes of t in the original tree.

The new trees are then tested in a validation set. We select the tree that gives the best performance.

Trying all possible trees might be computationally unfeasible. Several greedy techniques exist. See Hastie *et al.* (2009) for more details.

Note: pruning is not currently supported in `sklearn` (use `rpart` in `R` if needed)

Example of pruning



Advantages

- Build a non-linear and interpretable decision function
- Invariant under scaling and other linear transformations of the input data X
- Robust to the inclusion of (few) irrelevant features x
- It works for multi-class
- Computationally efficient: $O(\log F)$, where F is the number of terminal nodes
- It works for continuous and categorical variables

Drawbacks

- Low bias but very high variance. A small change in (all) input data can bring to a completely different tree (noisy) ! This instability is due to the hierarchical nature of the process. → averaging trees reduces the variance (bagging, random forests)
- No global optimization
- Separation hyperplanes are aligned with the feature axes → it might entail a sub-optimal solution
- Splits are hard. This creates piecewise-constant predictions with discontinuities at the split boundaries → prediction function is not smooth !

- ① Hastie, Tibshirani and Friedman (2009). The Element of Statistical Learning. Springer.
- ② L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone (1984). Classification and regression trees. Wadsworth Statistics/Probability Series.