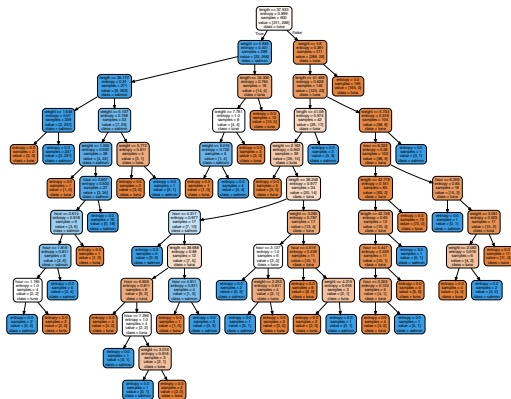


Fondamentaux théoriques du machine learning



Risks and risk decompositions

Probabilistic modelling

Optimization of neural networks

- Difficulties of optimizing neural networks

- Specific methods for neural networks

Classification and regression trees

- Decision trees

- Construction of a decision tree

- Tree pruning

Risks and risk decompositions

Probabilistic modelling

Optimization of neural networks

Difficulties of optimizing neural networks

Specific methods for neural networks

Classification and regression trees

Decision trees

Construction of a decision tree

Tree pruning

Notion of risk decomposition

Context : usual supervised learning.

- ▶ empirical risk of estimator $f : R_n(f)$ (random variable)
- ▶ risk (real risk, generalization error) of estimator $f : R(f)$ (real number).
- ▶ dataset D_n .

We would like to interpret the risk of any estimator.

Convergence of empirical risk

We fix $f \in H$ (hypothesis space). We assume that the samples (x_i, y_i) are i.i.d, with the distribution of (X, Y) , noted ρ . Then, according to the law of large numbers, under some assumptions (for instance, if the empirical risks are bounded), we have that in probability :

$$\lim_{n \rightarrow +\infty} R_n(f) = R(f) \quad (1)$$

The empirical risk of a fixed f converges to its real risk.

Notion of risk decomposition

We would like to compare $R(g)$, for some estimator g , to $R^* = R(f^*)$, f^* being the Bayes estimator.

First decomposition

- ▶ f^* : Bayes predictor
- ▶ F : Hypothesis space
- ▶ g : some predictor ($\in F$).

$$R(g) - R^* = \left(R(g) - \inf_{f \in F} R(f) \right) + \left(\inf_{f \in F} R(f) - R^* \right) \quad (2)$$

First decomposition

- ▶ f^* : Bayes predictor
- ▶ F : Hypothesis space
- ▶ g : some predictor
- ▶ f_n : empirical risk minimizer

Most of the time, we will be interested in such a decomposition for the learned estimator, e.g. the empirical risk minimizer f_n . Now, $R(f_n)$ is a **random variable** !

$$E[R(f_n)] - R^* = \left(E[R(f_n)] - \inf_{f \in F} R(f) \right) + \left(\inf_{f \in F} R(f) - R^* \right) \quad (3)$$

Underfitting and overfitting

Approximation error (bias) : depends on f^* and F , not on f_n , D_n .

$$\inf_{f \in F} R(f) - R^* \geq 0$$

Estimation error (variance, fluctuation error, stochastic error) : depends on D_n , F , f_n .

$$E(R(f_n)) - \inf_{f \in F} R(f) \geq 0$$

- ▶ too small F (compared to f^*) : underfitting (large bias, small variance)
- ▶ too large F (compared to n) : overfitting (small bias, large variance)

Deterministic bound on the estimation error

We consider the best estimator in the hypothesis space F .

$$f_a = \arg \min_{h \in F} R(h)$$

Exercise 1: Show that

$$R(f_n) - R(f_a) \leq 2 \sup_{h \in F} |R(h) - R_n(h)| \quad (4)$$

Deterministic bound on the estimation error

$$R(f_n) - R(f_a) \leq 2 \sup_{h \in F} |R(h) - R_n(h)| \quad (5)$$

Later in the course, based on **concentration inequalities** we will further build on this result and prove a probabilistic bound of the form

$$R(f_n) - R(f_a) \leq \frac{C}{\sqrt{n}} \quad (6)$$

(remember that by definition $0 \leq R(f_n) - R(f_a)$)

Approximate solution

- ▶ In machine learning, it is often not necessary to find the actual minimizer of the empirical risk , as there is an estimation error of $\mathcal{O}(\frac{1}{\sqrt{n}})$. [Bottou and Bousquet, 2009,]
- ▶ We can use an approximate solution \hat{f}_n , such that

$$R_n(\hat{f}_n) \leq R_n(f_n) + \rho \quad (7)$$

with ρ a predefined tolerance.

- ▶ This important because in large-scale ML, the **computation time** needs to be optimized.

Approximate solution

This gives a new risk decomposition :

$$\begin{aligned} E[R(\hat{f}_n)] - R^* &= \left(E[R(\hat{f}_n)] - E[R(f_n)] \right) \\ &\quad + \left(E[R(f_n)] - \inf_{f \in F} R(f) \right) \\ &\quad + \left(\inf_{f \in F} R(f) - R^* \right) \end{aligned} \tag{8}$$

$$\begin{aligned} E[R(\hat{f}_n)] - R^* &= \left(E[R(\hat{f}_n)] - E[R(f_n)] \right) \\ &\quad + \left(E[R(f_n)] - \inf_{f \in F} R(f) \right) \\ &\quad + \left(\inf_{f \in F} R(f) - R^* \right) \end{aligned} \tag{9}$$

$E[R(\hat{f}_n)] - E[R(\tilde{f}_n)]$ is the **optimization error**.

To conclude, we have :

- ▶ an approximation error
- ▶ an estimation error
- ▶ an optimization error

Risks and risk decompositions

Probabilistic modelling

Optimization of neural networks

Difficulties of optimizing neural networks

Specific methods for neural networks

Classification and regression trees

Decision trees

Construction of a decision tree

Tree pruning

Context

We are given a set of observations $\{y_1, \dots, y_n\} \in \mathcal{Y}$ that we assume are generated i.i.d from an unknown distribution. We look for a **probabilistic model** that explains well the data. We could for instance use this model to generate new data, that would be statistically similar to the observed ones.

Density estimation

We will consider **parametric models** for density estimation.

Definition

Parametric model

Let $\Theta \subset \mathbb{R}^p$. A parametric model \mathcal{P} is a set of probability distributions on \mathcal{Y} , indexed by elements of Θ .

$$\mathcal{P} = \{p_\theta | \theta \in \Theta\}$$

Examples :

- ▶ Bernoulli model (parameter p)
- ▶ Gaussian model (parameter (μ, σ))
- ▶ Binomial model (parameter (n, p))

Objective

If we assume that the data were generated from some $p_{\theta^*} \in \mathcal{P}$, with a unknown parameter θ^* , our goal is to find a good estimation of θ . If the data are indeed generated by a distribution in \mathcal{P} , the problem is said to be **well specified**. Otherwise, the problem is said to be **misspecified**.

Definition

Likelihood

Let $\mathcal{P} = \{p_\theta, \theta \in \Theta\}$ be a parametric model.

Given $y \in \mathcal{Y}$, the **likelihood** is the function :

$$\theta \mapsto L(\theta|y) = p_\theta(y) \quad (10)$$

Given $D_n = (y_1, \dots, y_n)$, the likelihood $L(.|D_n)$ is the function :

$$\theta \mapsto L(\theta|D_n) = \prod_{i=1}^n p_\theta(y_i) \quad (11)$$

The **maximum likelihood estimator** (MLE) is the parameter θ that maximises the likelihood :

$$\hat{\theta}_n \in \arg \max_{\theta \in \Theta} (L(\theta|D_n)) \quad (12)$$

Remarks

- ▶ Since the samples y_i are assumed to be independent, the likelihood corresponds to the probability (or probability density) of observing the dataset according to p_θ .
- ▶ We often maximise the log of the likelihood, as it is easier to differentiate a sum. Since log is an increasing function, the MLE is also the maximiser of the log of L .

Example 1

Exercise 2: We observe the data $(1, 0)$. We model these data with a Bernoulli distribution of parameter p .

- ▶ What is the likelihood of these observations as a function of p ?
- ▶ What is the value \hat{p} that maximizes this likelihood?

Example 2

Exercise 3: We observe the data $(1, 0, 1)$ (same hypotheses)

- ▶ What is the likelihood of these observations as a function of p ?
- ▶ What is the value \hat{p} that maximizes this likelihood?

Example 3

We observe the data $(2.5, 3.5)$. We assume that these data come from a normal law of parameters μ and σ .

$$\begin{aligned} L &= p(2.5|\mu, \sigma)p(3.5|\mu, \sigma) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{2.5-\mu}{\sigma})^2} \times \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{3.5-\mu}{\sigma})^2} \end{aligned} \quad (13)$$

We want to show that the likelihood is maximum for :

- ▶ $\hat{\mu} = \frac{2.5+3.5}{2}$
- ▶ $\hat{\sigma}^2 = \frac{(2.5-\hat{\mu})^2 + (3.5-\hat{\mu})^2}{2}$

ERM

In the context of density estimation, we can define a loss function as the **negative log-likelihood**.

$$\Theta \times \mathcal{Y} \mapsto -\log(p_{\theta}(y))$$

Given this loss, the risk writes :

$$R(\theta) = E_Y[-\log(p_{\theta}(y))]$$

and the empirical risk :

$$R_n(\theta) = -\frac{1}{n} \sum_{i=1}^n \log(p_{\theta}(y_i))$$

The MLE is then also the empirical risk minimizer.

KL divergence

The Kullback-Leibler divergence is a quantity used to compare two probability distributions.

Definition

Kullback-Leibler divergence

Given two distributions p and q , the KL divergence from p to q is defined as :

$$KL(p||q) = E_{Y \sim p} \left[\log \frac{p(Y)}{q(Y)} \right]$$

Lemma

If the data are generated by p_{θ^} , then $KL(p_{\theta^*}||p_{\theta})$ is the excess risk of p_{θ} , with the negative log-likelihood loss.*

Link with supervised learning methods

Probabilistic modelling can provide an interesting interpretation of several supervised learning methods, such as :

- ▶ logistic regression
- ▶ ordinary least squares

In a supervised learning context, we replace the likelihood $p_{\theta}(y)$ by a **conditional** likelihood $p_{\theta}(y|x)$ (conditional modelling).

$$R_n(\theta) = -\frac{1}{n} \sum_{i=1}^n \log (p_{\theta}(y_i|x_i)) \quad (14)$$

Link with logistic regression

We consider a binary classification problem, with $\mathcal{Y} = \{0, 1\}$.

Let us now consider the probabilistic model such that

$$p_{\theta}(1|x) = \sigma(\theta^T x)$$

Equivalently, this model can be written (remember that $y = 0$ or $y = 1$)

$$p_{\theta}(y|x) = (\sigma(\theta^T x))^y (1 - \sigma(\theta^T x))^{1-y} \quad (15)$$

Exercise 4: Show that the parameter θ with maximum likelihood is the logistic regression estimator θ_{logit} (cross entropy version).

Risks and risk decompositions

Probabilistic modelling

Optimization of neural networks

- Difficulties of optimizing neural networks

- Specific methods for neural networks

Classification and regression trees

- Decision trees

- Construction of a decision tree

- Tree pruning

Neural networks

References / tools :

- ▶ <https://www.deeplearningbook.org/>
- ▶ <https://d2l.ai/>
- ▶ https://mlelarge.github.io/dataflowr-web/dldiy_ens.html
- ▶ <https://playground.tensorflow.org/>
- ▶ <http://www.jzliu.net/blog/simple-python-library-visualize-neural-network/>

Learning representations / features

► $\mathcal{X} = \mathbb{R}^d$.

► $\mathcal{Y} = \mathbb{R}$.

A neural network with scalar output learns a feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ **and** a linear regressor or classifier, $\theta \in \mathbb{R}^m$.

$$\forall x, f(x) = \langle \theta, \phi(x) \rangle \quad (16)$$

We can add a intercept (bias) by adding a dimension to θ and adding a component with a 1 to each $\phi(x)$.

Learning representations / features

► $\mathcal{X} = \mathbb{R}^d$.

► $\mathcal{Y} = \mathbb{R}$.

A neural network learns a feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ **and** a linear regressor or classifier, $\theta \in \mathbb{R}^m$.

$$\forall x, f(x) = \langle \theta, \phi(x) \rangle \quad (17)$$

We can add a bias by adding a dimension to θ and adding a component with a 1 to each $\phi(x)$.

Remark : kernel methods use hardcoded features ϕ , that can be **implicit** (kernel trick).

Multi output

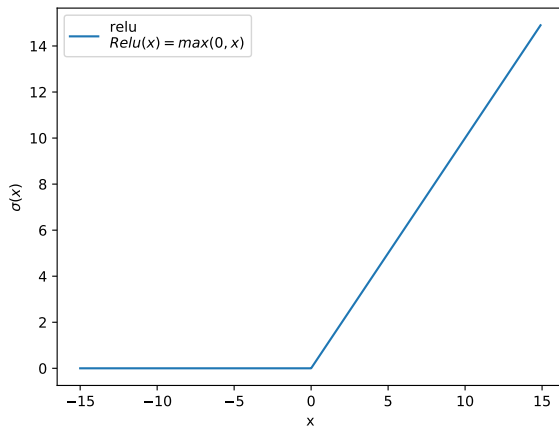
In order to learn a multidimensional output of dimension p ($\mathcal{Y} = \mathbb{R}^p$), it is sufficient to learn a matrix $\theta \in \mathbb{R}^{m,p}$.

Single layer neural network

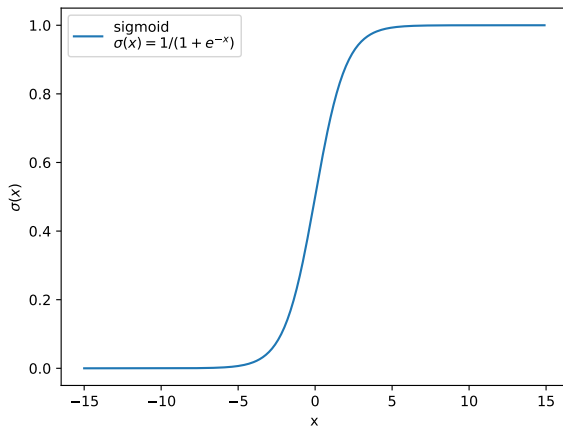
$$f(x) = \sum_{j=1}^m \theta_j \sigma(w_j^T x + b_j) \quad (18)$$

σ is an activation function (sigmoid, tanh, ReLu, etc).

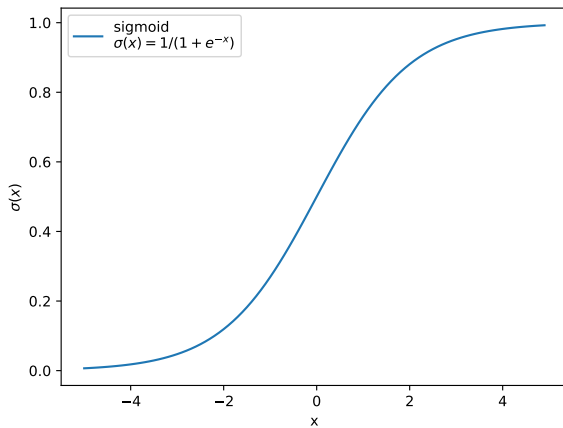
ReLU



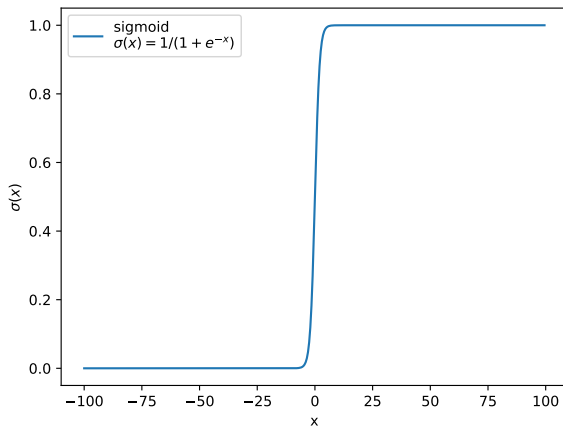
Sigmoid



Sigmoid



Sigmoid



Automatic differentiation

When working with neural networks, most used libraries implement automatic differentiation.

- ▶ tensorflow
- ▶ pytorch (autograd)

Optimizing neural networks

Optimizing neural networks comes with specific difficulties.

- ▶ the problem is non-convex
- ▶ there is often a large number of parameters (optimization in a high dimensional space)
- ▶ specific problems due to depth (vanishing gradients, see below)

<http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>

Non-convexity

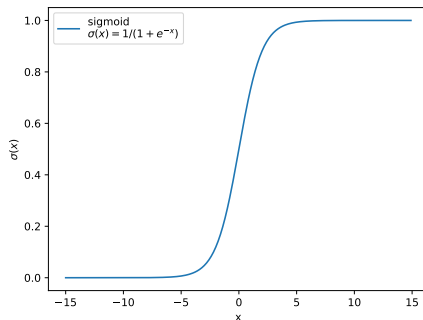
We know that

- ▶ If f is increasing and convex and g is convex, then $f \circ g$ is convex.
- ▶ If f is convex and g is linear, then $f \circ g$ is convex.

With neural networks, we are in neither of these cases, as the activations σ are non linear.

Hence **the objective function is non-convex**, and it remains difficult to understand why gradient based methods often perform well in practice

Vanishing gradients



Exercise 5 :

What is the maximum value of $|\sigma'(z)|$?

Exponentially decreasing gradients

- ▶ At each layer, the gradients are multiplied by a term of the form $\sigma'(u)$. Using a large number of layers leads to gradient norms that decrease rapidly when we move away from the output layer.
- ▶ This slows training down and caused deep learning to plateau for some years.
- ▶ Several initializations were necessary in order to obtain convergence, the result was unstable.

- └ Optimization of neural networks
 - └ Difficulties of optimizing neural networks

ReLU

The usage of ReLU solved this problem.

Other activation functions :

[https://dashee87.github.io/deep%20learning/
visualising-activation-functions-in-neural-networks/](https://dashee87.github.io/deep%20learning/visualising-activation-functions-in-neural-networks/)

SGD variants for neural networks

Several specific variations of SGD are commonly used for deep learning.

<https://pytorch.org/docs/stable/optim.html>

Specific methods for deep learning

Architectures :

- ▶ Convotutional networks
- ▶ Residual neural network (ResNet)

Optimization / regularization :

- ▶ dropout
- ▶ batch normalisation

Risks and risk decompositions

Probabilistic modelling

Optimization of neural networks

Difficulties of optimizing neural networks

Specific methods for neural networks

Classification and regression trees

Decision trees

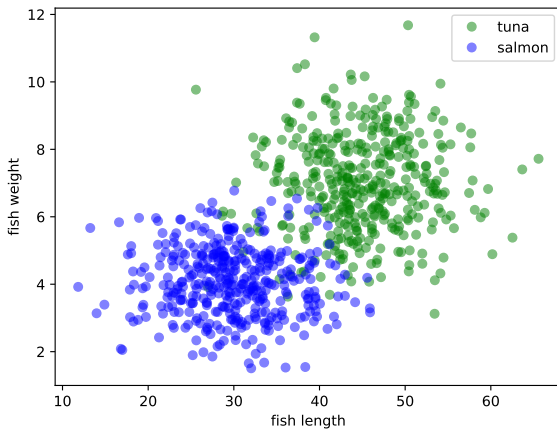
Construction of a decision tree

Tree pruning

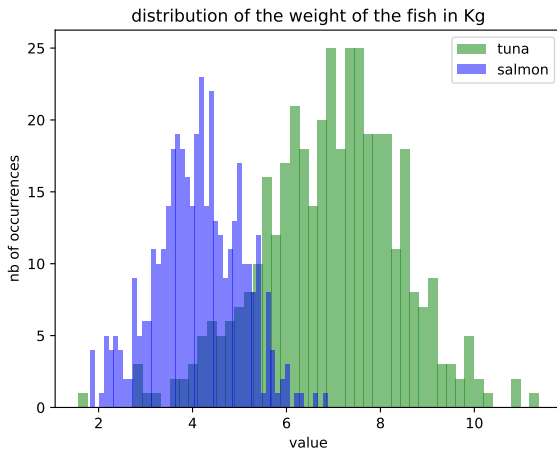
CART

- ▶ Segmentation method (partitionnement récursif), **binary splits**.
- ▶ First algorithm : **CART** (classification and regression tree, Breiman, 1984) [Breiman et al., 2017]

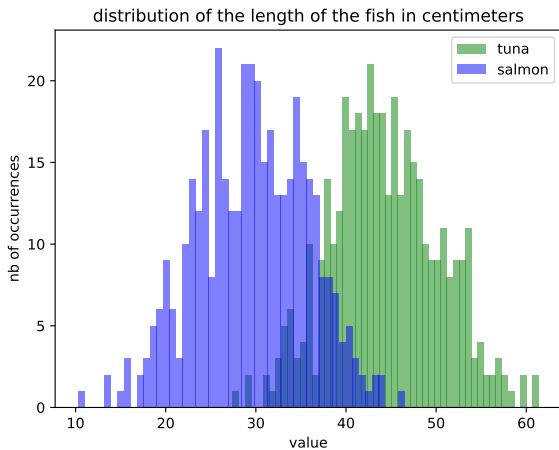
Example dataset : fishes



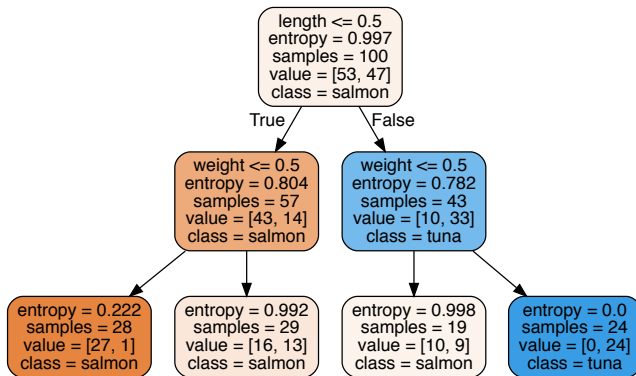
Fish dataset



Fish dataset

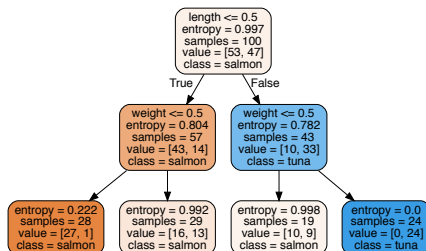


Example tree



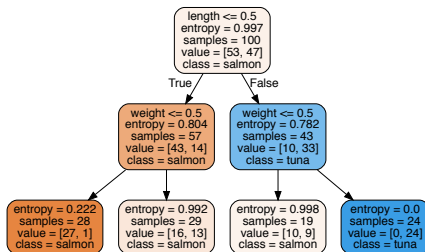
Splitting nodes

Each split should lead to more **homogeneous** nodes (in a sense that is to be formally defined)

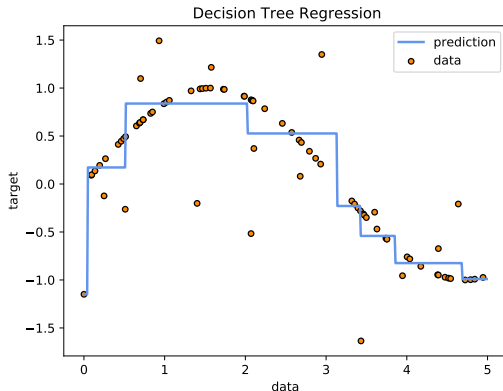


Splitting nodes

A leaf node corresponds to a predicted value.



Decision trees can be used for regression. In that case, the prediction is **piecewise constant**. It can be seen as a form of local averaging.



Construction of a tree

A split node corresponds to :

- ▶ a segmentation variable V
- ▶ a segmentation value / threshold if V is quantitative. If V is a class, it is a two-fold partitionning of the set of classes.

The construction of a tree requires a criterion to :

- ▶ choose the segmentation variable and value
- ▶ decide when a node is terminal (leaf node)
- ▶ associate a prediction value to all leaf nodes

Homogeneity

We want the homogeneity in the child nodes to be greater than in the parent node. Equivalently, we can define a non-negative heterogeneity function H that describes each node and that must be :

- ▶ 0 if the node is homogeneous (only one class or one output value is represented in this node)
- ▶ maximal if the values of y are uniformly distributed within the node.
- ▶ if $H(L)$ is the value of H for the left child node (and $H(R)$ for the right one), then the segmentation should minimize

$$H(L) + H(R) \tag{19}$$

Homogeneity criterion for regression

In the case of regression, $\mathcal{Y} = \mathbb{R}$ and the heterogeneity $H(n)$ of node n is its empirical variance :

$$H(n) = \frac{1}{|n|} \sum_{i \in n} (y_i - \bar{y}_n)^2 \quad (20)$$

Homogeneity criterion for classification

In the case of classification, $\mathcal{Y} = [1, \dots, L]$ and several criteria exist.

Homogeneity criterion for classification : entropy (information gain)

- ▶ We use the convention that $0 \log(0) = 0$
- ▶ p_n^l is the proportion of class l in node n .

The **entropy** writes :

$$H(n) = - \sum_{l=1}^L p_n^l \log(p_n^l) \quad (21)$$

Exercice 6 : What are the maximum and minimum values of the entropy ?

Homogeneity criterion for classification : Gini impurity

The Gini impurity writes :

$$H(n) = \sum_{l=1}^L p_n^l (1 - p_n^l) \quad (22)$$

Homogeneity criterion for classification : Gini impurity

The Gini impurity writes :

$$H(n) = \sum_{l=1}^L p_n^l (1 - p_n^l) \quad (23)$$

Exercise 7 : Interpret the meaning of the Gini impurity in terms of probabilities, assuming that we predict according to the proportions p_n^l .

Homogeneity criterion for classification : misclassification probability

$$H(n) = 1 - \max_l(p_n^l) \quad (24)$$

If we predict the most represented class in n , then this is the misclassification probability.

Prediction for a leaf node

- ▶ **regression** : predict the average value in the node
- ▶ **classification** : several possibilities
 - ▶ most represented class in the node
 - ▶ most probable class *a posteriori* if the *a priori* probabilities are known (Bayesian probabilities)
 - ▶ class that costs the less in case of missclassification

Pruning

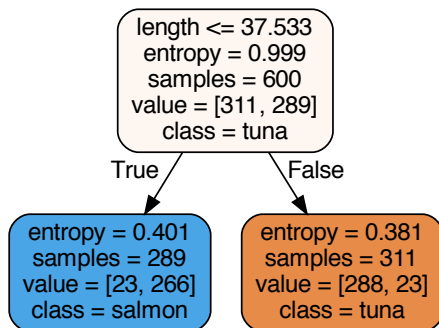
- ▶ If the segmentation is not restricted, the tree can fit the training dataset perfectly (overfitting).
- ▶ learning would then be highly dependent on the choice of the training dataset, leading to **instability**.
- ▶ To avoid it, **pruning** (élaguage) is used : it consists in restricting the size of the tree (and can be seen as an equivalent to regularization like in ridge regression or logistic regression).

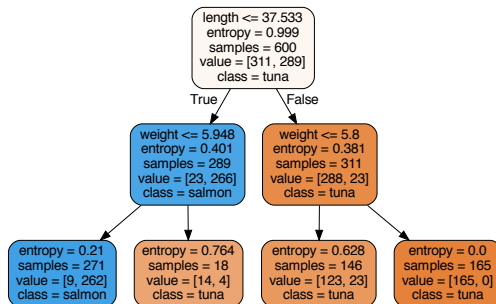
Pre-pruning by restricting divisions

We can impose

- ▶ minimum impurity decrease
- ▶ a minimum number of samples in a leaf node
- ▶ a minimum number of samples to authorize splitting a node

Simulation : fish dataset, max depth=1

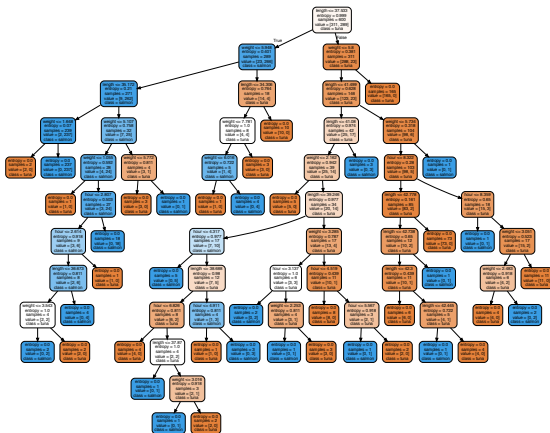




FTML

Classification and regression trees

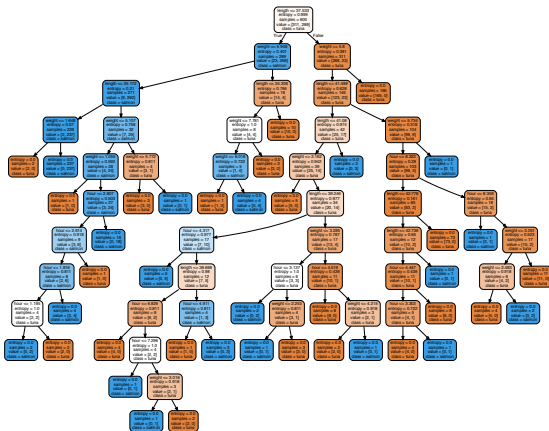
Tree pruning



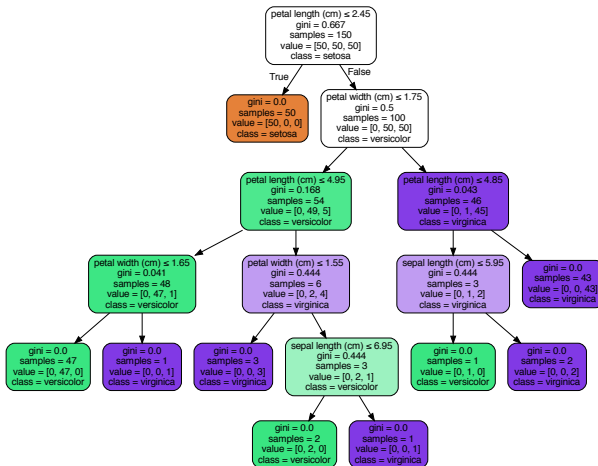
FTML

Classification and regression trees

Tree pruning



Overfitting : Iris dataset, max depth=34



Avoiding overfitting : Iris dataset

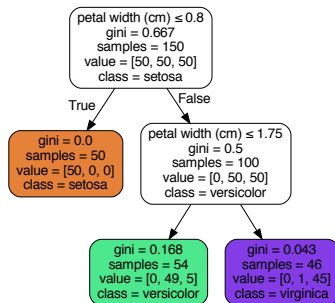


Figure – Avoid overfitting with pre-pruning. In this case, impose a minimum impurity decrease of 0.1

Post-pruning

It is also possible to prune after the construction of a large tree.

Remarks on decision trees I

- ▶ Decision trees do not require hypotheses on the distribution of de features.
- ▶ Feature selection is integrated in the algorithm, so they might be relevant in a situation with many features (variables explicatives)
- ▶ A segmentation is invariant to a monotonic change in a feature. It is thus robust to outliers or to very asymmetrical distributions.
- ▶ Decision trees allow a straightforward **interpretation** of the model.

Remarks on decision trees II

- ▶ Decision trees are greedy : they might find a local minimum
- ▶ they are hierarchical : a bad choice in a segmentation at the top of the tree propagates in the whole tree
- ▶ Hence their instability and sensibility to the choice of the training set.

References I



Bottou, L. and Bousquet, O. (2009).

The tradeoffs of large scale learning.

*Advances in Neural Information Processing Systems 20 -
Proceedings of the 2007 Conference*, (January 2007).



Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J.
(2017).

Classification And Regression Trees.

Routledge.