**INSA**
ROUEN NORMANDIE

## Machine Learning

Ensemble learning, Random Forest, Boosting

Simon BERNARD
simon.bernard@univ-rouen.fr

Chapters :

# Introduction to ensemble learning

## Supervised learning

- Inputs : $x \in \mathcal{X} \subset \mathbb{R}^d$ (often called instances)
- Outputs [1] : $y \in \mathcal{Y} = \{\lambda_1, \lambda_2, \ldots, \lambda_c\}, c \geq 2$
- Core hypothesis : $\exists$ an unknown function $f : \mathcal{X} \to \mathcal{Y}$ that maps any input to its output
- Goal : find an approximation $h$ of $f$, s.t. $\hat{y} = h(x)$ is a good output prediction $\forall\, x \in \mathcal{X}$.
- Core idea : determine $h$ from a set of $n$ labeled instances :

$$\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)\}$$

where the $x_i \in \mathcal{D}$ are associated to their "true" output values $y_i$

1. For simplicity, we focus on classification but most of the concepts in this course are readily transposable to regression tasks.
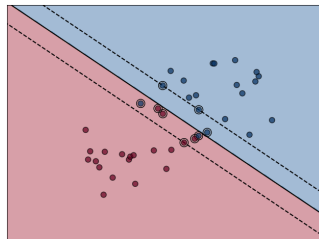
## Supervised learning

- The goal of machine learning is to find the best $h$ possible
- To do so, we rely on
    1. a hypothesis space $\mathcal{H}$
    2. a learning algorithm $\mathcal{A}$ that can output any $h \in \mathcal{H}$
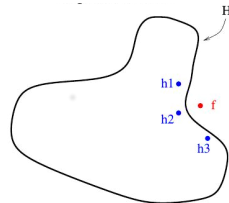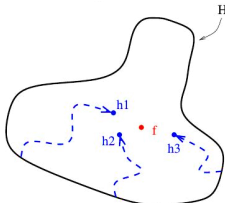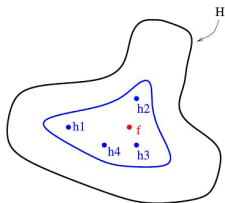
- $\mathcal{H}$ : linear models

    $$h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b, \quad \forall\, \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

- $\mathcal{A}$ : linear SVM for finding the best $h \in \mathcal{H}$

## Limitations of learning a single $h$

1. Statistical : if $\mathcal{D}$ is too small compared to $\mathcal{H}$, many $h_t$ that performs equally on $\mathcal{D}$
2. Computational : if $\mathcal{A}$ proceeds by local search, may lead to suboptimal solutions
3. Representational : in many situations, $f$ can not be represented by $\mathcal{H}$

**Limitations of learning a single $h$**

1. Statistical : if $\mathcal{D}$ is too small compared to $\mathcal{H}$, many $h_t$ that performs equally on $\mathcal{D}$
2. Computational : if $\mathcal{A}$ proceeds by local search, may lead to suboptimal solutions
3. Representational : in many situations, $f$ can not be represented by $\mathcal{H}$

Solution : Ensemble learning
1. Combining the $h_t$ reduces the risk of choosing the wrong one
2. Learning from different "starting points" and combining is better
3. A weighted combination of several $h_t$ expands $\mathcal{H}$

**Generic principle of ensemble learning**

- Produce several different $h_t, \forall t = 1, \ldots, T$
- Combine these $h_t$ with a combining operator $\Omega$ such that

$$\hat{y} = \Omega(h_1(\mathbf{x}), \ldots, h_T(\mathbf{x})), \forall \mathbf{x} \in \mathcal{X}$$

For example :

- $\mathcal{Y} = \{-1, 1\}$ and $h(\mathbf{x}) \in [-1, 1], \forall t = 1, \ldots, T, \forall \mathbf{x} \in \mathcal{X}$
- Split $\mathcal{D}$ into $T$ disjoints subsets $\mathcal{D}_t$ and learn a $h_t$ from each one
- $\Omega$ is typically the sign of the summed score :

$$\hat{y} = sign\left(\sum_{t=1}^{T} h_t(\mathbf{x})\right)$$

Key ingredients : accuracy and diversity

**Accuracy vs Diversity :**

1. A classifier is said to be accurate if it's better than random guessing, "at least"
2. The diversity between classifiers is high if their errors are mostly different

- Maximize both accuracy and diversity at the same time
- But these two objectives are generally opposed
- The challenge is to create diversity while maintaining sufficient individual accuracy

## How to create diversity ?

1. Randomization in the learning algorithm
   ex. : Random init. of parameters (e.g. weights of NN)
2. Manipulate the training instances
   ex. : Use training subsets
3. Manipulate the features
   ex. : Use description subspaces
4. Manipulate the output values
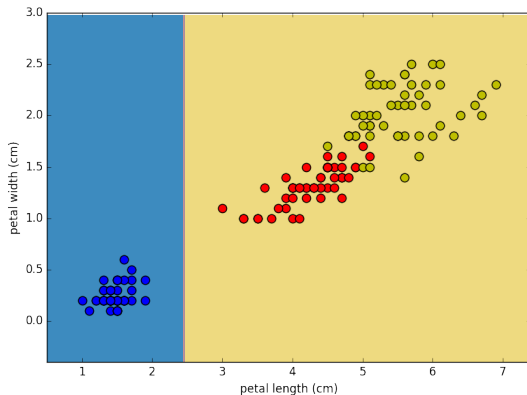   ex. : One-vs-One or One-vs-All strategies

# Decision tree classifier

Example : predict the iris species from the width and length of the petals



*Petal width* vs *Petal length* :

→ We build the classifier from these 2 features, with very simple and intuitive decision rules

Example : predict the iris species from the width and length of the petals



If *Petal length* $\leq 2.45$

$$h(\mathbf{x}) = setosa$$

Else

$$h(\mathbf{x}) = virginica$$

(or *versicolor*, difficult to decide)

Example : predict the iris species from the width and length of the petals



If *Petal length* $\leq 2.45$

$$h(\mathbf{x}) = setosa$$

Else if *Petal width* $\geq 1.75$

$$h(\mathbf{x}) = virginica$$

Else

$$h(\mathbf{x}) = versicolor$$

Example : predict the iris species from the width and length of the petals



If *Petal length* $\leq 2.45$

$$h(\mathbf{x}) = setosa$$

Else if *Petal width* $\geq 1.75$

$$h(\mathbf{x}) = virginica$$

Else if *Petal length* $\leq 4.95$

$$h(\mathbf{x}) = versicolor$$

Else

$$h(\mathbf{x}) = virginica$$

## Example : predict the iris species from the width and length of the petals

A tree is made up of nodes :

1. Root node or initial node
2. Splitting node or internal node
3. Leaf node or terminal node

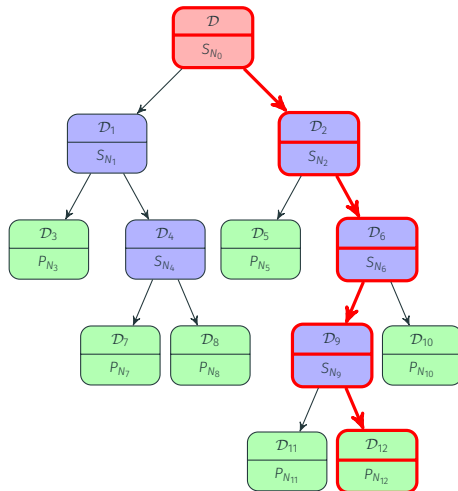The root node contains the whole training set $\mathcal{D}$, the other nodes contains subsets $\mathcal{D}_k$

Root and splitting nodes contains splitting rules :

$$S_{N_i} : \mathcal{X} \to \{N_{left}, N_{right}\}$$

Most often, it is a test on the value of one feature

$$S_{N_i}(\mathbf{x}) = \begin{cases} N_{left} & \text{if } x^{(i)} \leq \alpha_i \\ N_{right} & \text{otherwise} \end{cases}$$

Predict a new instance x :

1. Feed **x** to the root node
2. Evaluate $S_{N_0}$
3. Redirect to child node ($N_2$)
4. Evaluate $S_{N_2}$
5. and so on till reaching a leaf

Leaf node gives a prediction or a vector of class proba. :

$$P_{N_i} = \{\hat{P}(\lambda_1|\mathbf{x}), \hat{P}(\lambda_2|\mathbf{x}), \dots, \hat{P}(\lambda_C|\mathbf{x})\}$$

# Decision tree learning

## Generic procedure for DT learning

1. Init. the root node $N$ with $\mathcal{D}$
2. To split or not to split?
   2.1 To split :
       2.1.1 Select a splitting rule $S_N$
       2.1.2 Build the child nodes $N_{left}$ and $N_{right}$ according to $S_N$
       2.1.3 Go to step 2 with $N = N_{left}$ and with $N = N_{right}$
   2.2 Not to split : assign a class to predict (or $C$ class proba.) to $N$

## Questions to adress :

- How to select a splitting rule ? (step 2.1.1)
- How to decide if a node should be split or not ? (step 2)
- How to decide which class to predict at a given leaf ? (step 2.2)

Splitting rule based on a single feature $x^{(i)}$

- Find the $S_N = (x^{(i)}, \alpha_i)$ that best separates the instances according to their classes
- The classical procedure is an exhaustive search :



1. At node $N_k$, $S_{N_k}$ must be set from $\mathcal{D}_k$ only
2. Try all $x^{(i)}$ and for each of them, every possible $\alpha_i$ in $\mathcal{D}_k$
3. Retain the $(x^{(i)}, \alpha_i)$ that optimizes the splitting criterion

## Splitting criterion

· The splitting criterion is a measure of the quality of the split
· Based on impurity : a node is said to be pure if all the instances belong to the same class



· For a given impurity measure $I$ (here the Gini index)
· $I(\mathcal{D}_i) = 0.667$ when class distrib. is $\{50, 50, 50\}$
· $I(\mathcal{D}_i) = 0$ for left child since class distrib. is $\{50, 0, 0\}$

## Splitting criterion

- For each candidate rule $S_{N_i}$, the criterion is the impurity gain ($\nearrow$) :

$$\Delta(\mathcal{D}_i, S_{N_i}) = I(\mathcal{D}_i) - \left( \frac{|\mathcal{D}_{left}|}{|\mathcal{D}_i|} I(\mathcal{D}_{left}) + \frac{|\mathcal{D}_{right}|}{|\mathcal{D}_i|} I(\mathcal{D}_{right}) \right)$$

- Two popular impurity measures :
    - *Gini index*

$$G(\mathcal{D}_i) = 1 - \sum_{k=1}^{C} \left( \frac{n_{k.}}{|\mathcal{D}_i|} \right)^2$$

    - *Entropy (Information Gain)*

$$E(\mathcal{D}_i) = \sum_{k=1}^{C} -\frac{n_{k.}}{|\mathcal{D}_i|} log_2 \frac{n_{k.}}{|\mathcal{D}_i|}$$

Early stopping and prediction assignment

1. Stop a branch from growing :
   - Never : a node $N_k$ become a leaf node when it's pure
   - Based on threshold on depth, number of instances, impurity gain, etc.

2. which prediction to assign to the leaf node :
   - the majority class in $\mathcal{D}_k$
   - class counts in $\mathcal{D}_k$ + smoothing methods (recommended)
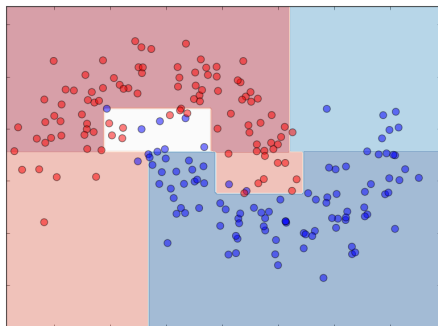     Example : Laplace smoothing

$$\hat{P}(\lambda_c|\mathbf{x}) = \frac{n_{c.} + 1}{n_{..} + C}$$

## Early stopping and prediction assignment



Max. depth tree always have 0% training error but very poor generalization performances

## Early stopping and prediction assignment



Early stopping can help but one can also cut some branches a posteriori (pruning)
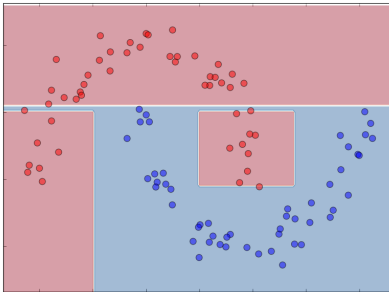
# Ensemble of decision trees

## Unstable classifier

A classifier is said to be unstable if small changes in the learning set involves significant changes in the resulting model.
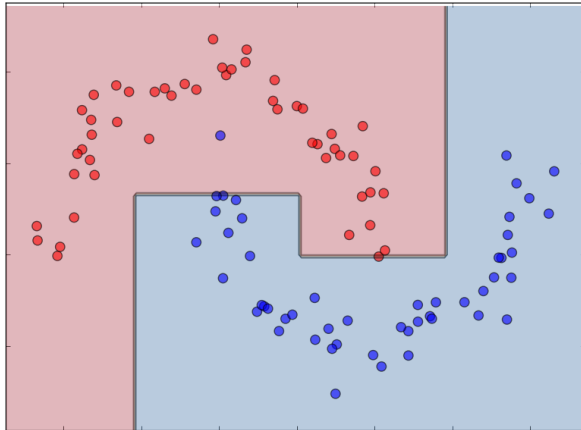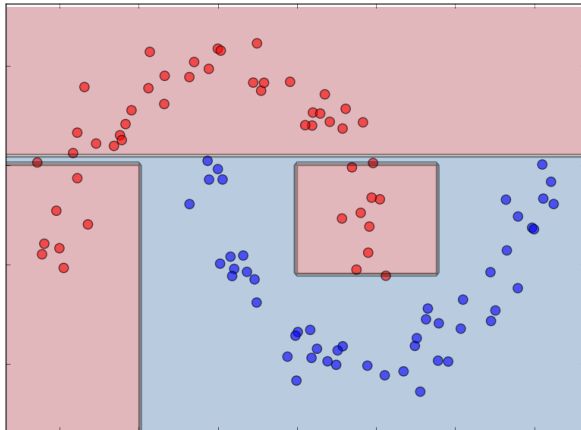
**Unstable classifier**

A classifier is said to be unstable if small changes in the learning set involves significant changes in the resulting model.

- Instability leads to a very high variance (and poor generalization performances)
- Pruning methods are generally not sufficient
- Better solution are ensembles of decision trees (decision forests)
- Why ? Because instability is a real asset to create diversity
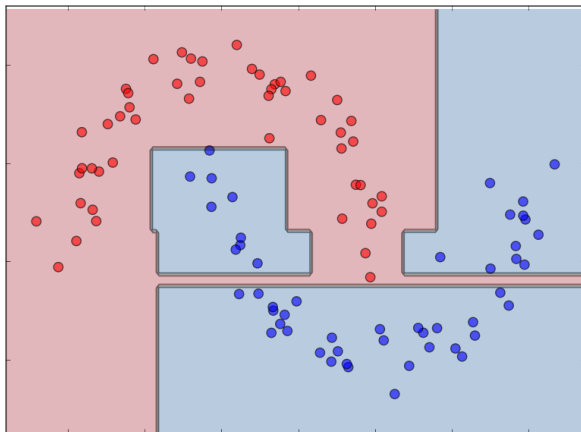
A first decision tree built (without pruning) on 100 training instances (50 in each class)

A second tree built on 100 additional training instances (from the same distribution)
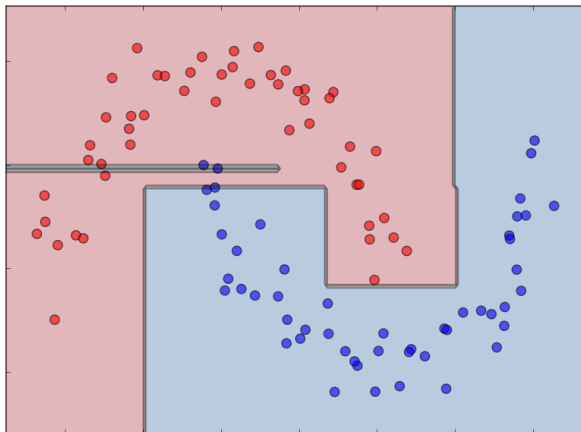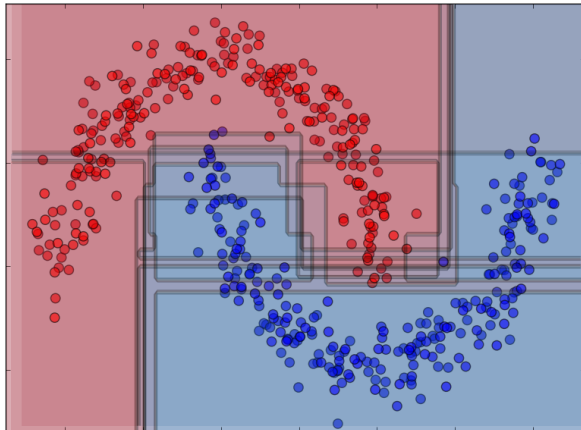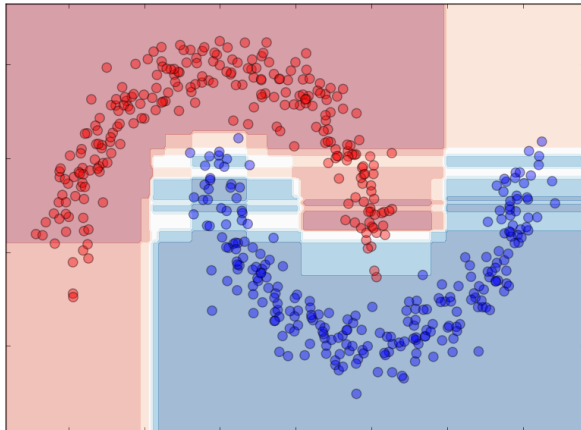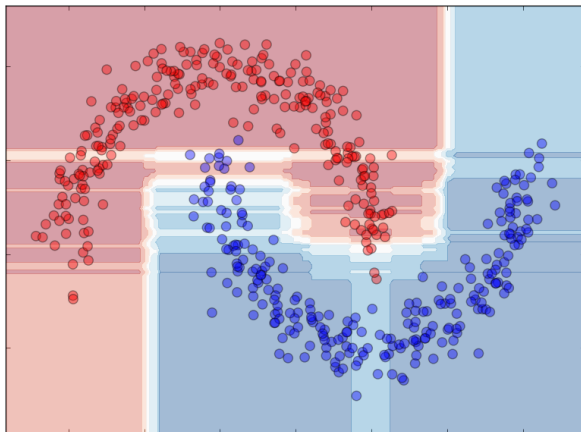
Third tree...

Fourth...

And fifth.

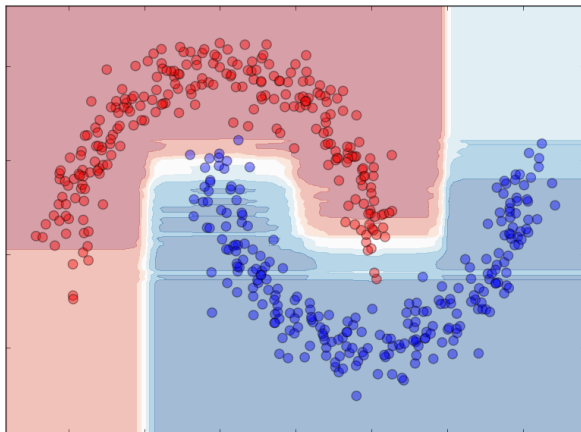All of them in the same plot (with their respective training instances)

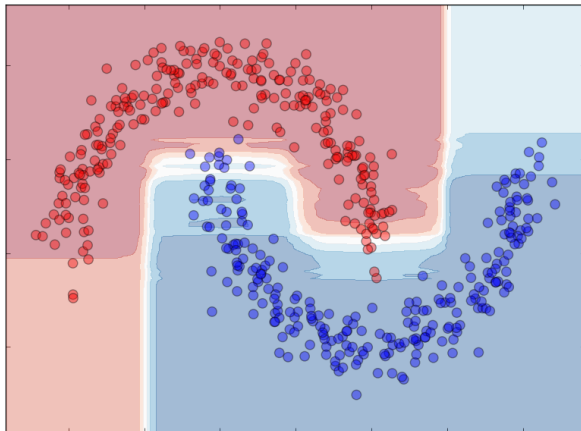And what about making them vote for the final prediction ?

And by combining 10 trees ?

And 100 trees ?

1000 ?!! (< 4 secondes for training)

Takeaways

DT are...

- easy to understand and interpret
- flexible (can handle many different types of data)
- reasonnably fast with reasonnably large datasets

But...

- do not scale up ("non reasonnably" large datasets)
- are unstable (high variance)
- are not very accurate in general

However, they are ideal for ensemble methods