

# Decision Trees and Random Forests

## Data science Certificate

Aymeric DIEULEVEUT

May 2022

- Supervised Learning framework
- why
- how we can use a DT
- how to build it!
  - \* what do we need to choose
  - \* what needs to be learned, how?
- how to implement a DT in practice (Python)
- + / -

## 1 Decision Trees

## 2 Random Forests

# Outline

1 Decision Trees

2 Random Forests

# A Reminder on Supervised Learning

- Observations:  $(X_i)_{i=1,\dots,n}$ . Each  $X_i$  has  $d$  components.
- Outputs:  $(Y_i)_{i=1,\dots,n} \in \mathcal{Y}$

3 Examples:

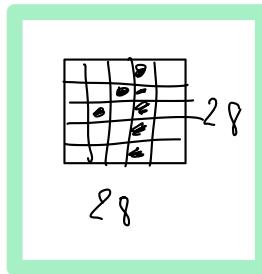
- Iris Dataset: predict type of Iris from

Classif  $\mathcal{Y}$  points  
handwritten

$\mathcal{Y}$   
 $d$  features  
characteristics of the flower!

- MNIST: digit predict.

Classif



- House price prediction

Regress

$$Y \rightarrow (\text{price})$$

$$X \rightarrow (\text{size}, \text{location}, \text{etc})$$

$$\mathcal{Y} = 1$$

$$\mathcal{Y} = \{0, \dots, 3\}$$

## Example: Iris Dataset



Figure: Iris: setosa, versicolor, virginica

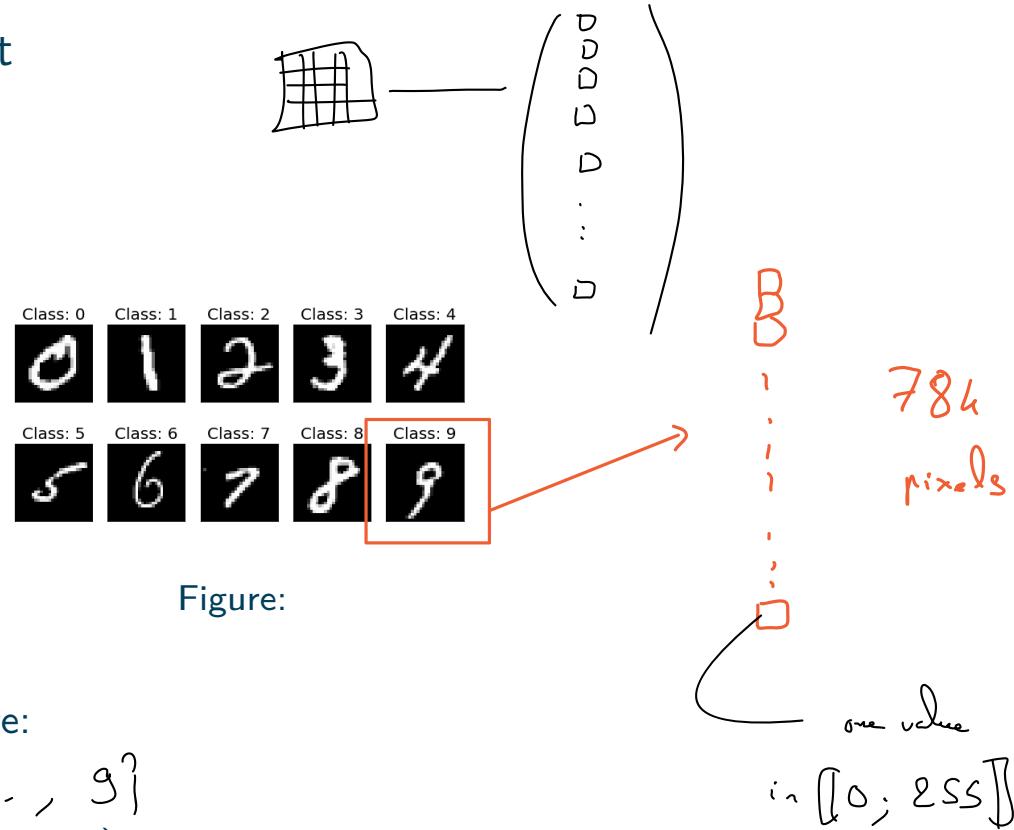
One of the most widely used toy dataset in classification:

- 3 classes : `['setosa', 'versicolor', 'virginica']` = 3 .
  - 50 points per class.  $\rightarrow n = 150$  number of training points
  - 4 features: `['sepal length', 'sepal width', 'petal length', 'petal width']`

*input*       $d = 4$

*4 continuous variables .*

## Example: MNIST Dataset



Digit recognition from an image:

- 10 classes :  $Y_i \in \{0, \dots, 9\}$
- 60k images, (50k train, 10k test), balanced
- 784 features =  $d$

↪  $28 \times 28$

$$\text{Input} \in \mathbb{R}^d$$

# House price Prediction

- Output : house value  $Y \in \mathbb{R}$   regression task!
- 1500 examples
- 81 features: Construction year, Neighborhood, Surface, Floor, Balcony: both quantitative and qualitative, some ordinal variables.

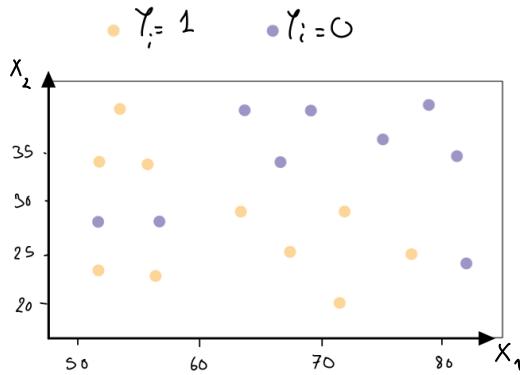
quantitative

Simplified version:

3 features:	$\text{surface} \in \mathbb{R}$
	$\text{neighborhood} \in \{0, \dots, 20\}$
	$\text{has garden} \in \{0, 1\}$

$$Y = \begin{cases} 1 & \text{if } \text{price} > \$1M \\ 0 & \text{else.} \end{cases}$$

# T Gy dataset

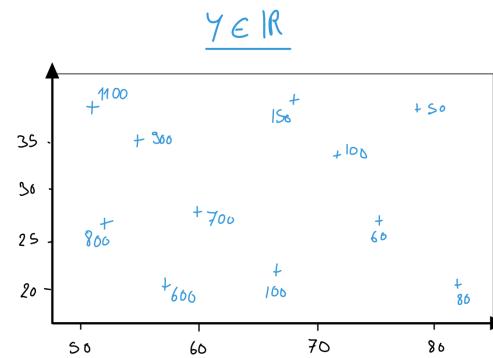


$$\mathcal{Y} = \{0, 1\}$$

$$d = 2$$

classif ~~reg.~~

$$n = 19$$



$$\mathcal{Y} = \{0, 1\}$$

$$d = 2$$

reg.

$$n = 11$$

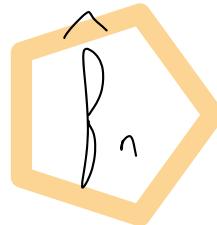
## Summary of those three datasets

	Iris	MNIST	House price prediction	Toy
Number of examples $n$	150	60k	1500	19
Number of the features $d$	4	78k-28 <sup>2</sup>	81	3
Type of the features	quantit.	quant.	quant + qual.	quant
Space $\mathcal{Y}$	{0,1,2}	{0,...9}	$\mathbb{R}_+$	{0,1}
Task	classif	classif	regression	classif

Pipeline:

## Training phase

$$\left( \begin{array}{l} \text{train set} \\ X_i^{\text{train}}, Y_i^{\text{train}} \end{array} \right)_{i=1..n}$$



such that

$$\hat{f}_n(X_i^{\text{train}}) \approx Y_i^{\text{train}}$$

\* type of fundo?

\* how to train? [need hyperparam]

Validation phase: find good hyperparameters.

Test phase

$$X_i^{\text{test}}$$



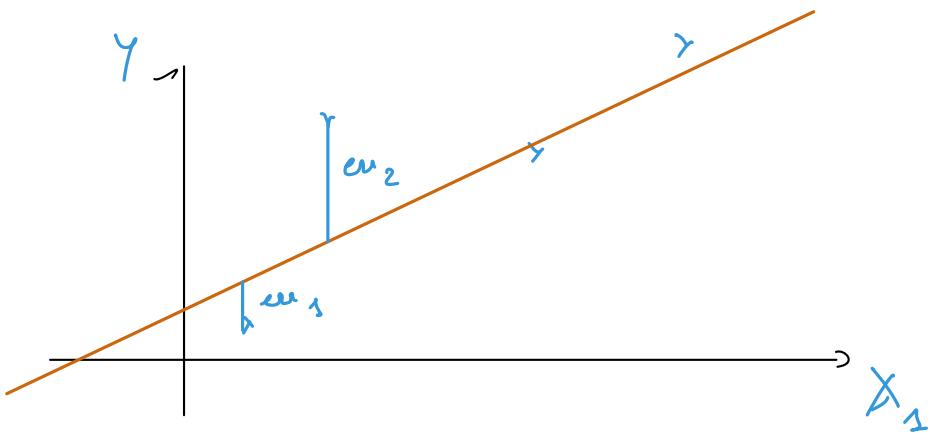
$$\hat{f}_n(X_i^{\text{test}})$$

compare to  $Y_i^{\text{test}}$

Use it in practice ~~X~~

↳  $\text{error}(\text{test set})$

Linear regression:



\* type of pred = ? Linear

\* how to train ? Find the line that minimizes a criterion  
 $\sum (\text{err})^2$

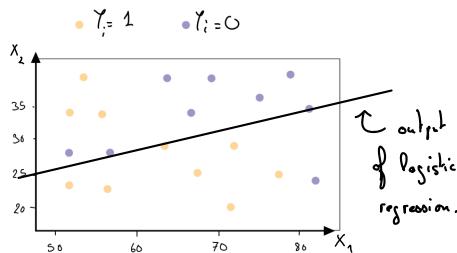
Logistic regression

vector  $\hat{\theta}$  /

$$\sigma(\langle \hat{\theta}, x_i \rangle) \geq \text{prob that } Y_i = 1$$

predict with

$$\text{if } \sigma(\langle \hat{\theta}, x \rangle) \geq \frac{1}{2} = \begin{cases} 1 & \text{if } \sigma(\hat{\theta}, x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$



## Limitations of linear models:

- ① \* data is not always linearly separable!
  - ② \* interpret. is ok-ish but still not easy for a human.
- \* Regression: House - price dataset.

Decision  
trees

$$Y = \underbrace{4000}_{\text{Learned!}} \times \text{surf} + \underbrace{16000}_{\text{Learned!}} \times \text{Floor} - \underbrace{1000}_{\text{Learned!}} \times \text{neighborhood}$$
$$+ 6000 \text{ surf}(2m_h) + \dots +$$

81 terms in the sum.

$$\begin{array}{c} 100 \text{ m}^2 (\rightarrow 2m_h) \\ | \\ 110 \text{ m}^2 \text{ overall} \end{array} \quad \begin{array}{c} 6000 \times 160 \\ | \\ 100 \times 110 \end{array}$$

# Decision Tree

## Goals of Decision Trees

- ① Supervised Learning: solve the classification or regression task
- ② Provide some understanding: what on this example allows me to classify it? What are the important features ?

↳ have an easy to interpret method

3. Go beyond linear models

# Decision Tree

## Goals of Decision Trees

- ① Supervised Learning: solve the classification or regression task
- ② Provide some understanding: what on this example allows me to classify it? What are the important features ?

Two questions:

- How to train a decision tree? → **training phase (learning algorithm)** *moderate*
- How to predict with a given decision tree? → **inference phase** *easy*

## Inference Phase

House Price - simplified

- Root = observation = input

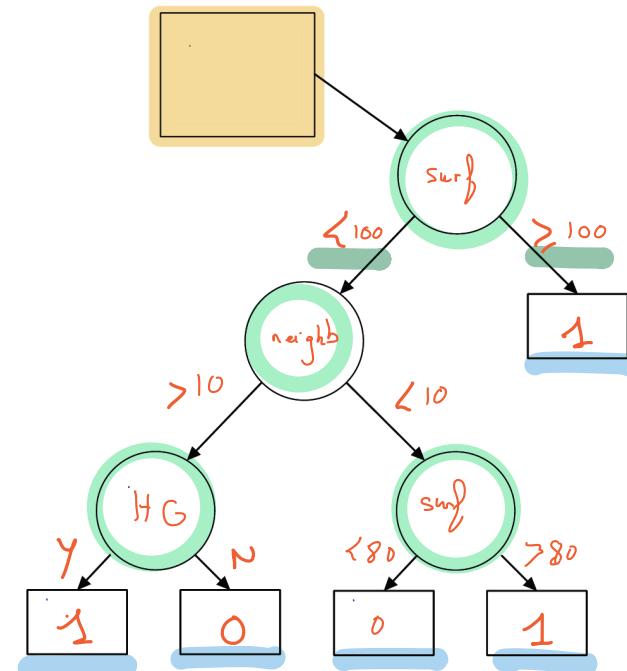
- Nodes = tests  
Looking at one variable

- Branches = possible outcomes  
only 2 ways

- Leaves = final decision = output

\* Inference is very easy, very fast and easy to interpret! → turn into a policy! (ex: Covid and vaccination calendar)  
\* Training: learn all nodes and splits!

$$x_{\text{new}} = (92, 6^*, 0)$$



# Training: how to build a tree.

## Data:

$X_1 = (62, 1, 0) \quad Y_1 = 1$   
 $X_2 = (72, 2, 1) \quad Y_2 = 1$   
 $X_3 = (50, 3, 1) \quad Y_3 = 0$   
⋮  
 $X_n = (45, 15, 1) \quad Y_n = 1$

ultimarily:

\* small

\* minimize error

xtrain.set : objective

xtest set : avoid overfitting  
(cross-validation - hyperpar)

## Key question !

- What is a good tree?
- What do I need to choose? / Learn?
- How to choose? / Learn?

## What do you think?

Choose

hyperparameters:  
\* depth  
\* nb of training pts  
(x stopping criteria)

Learn : splits / criterias

hyperparam : Cross validate

define a mathematical criteria to find the BEST SPLIT

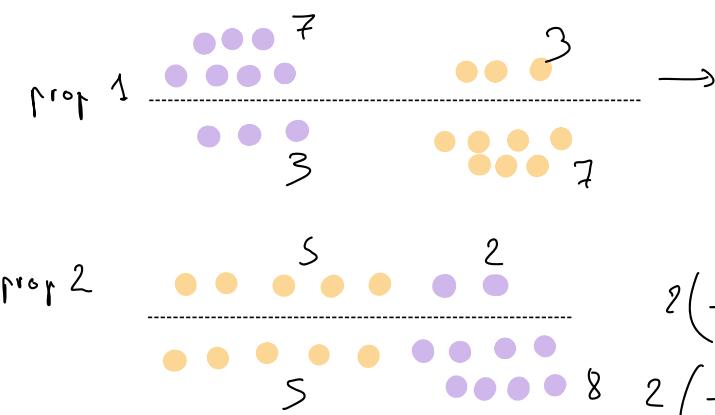
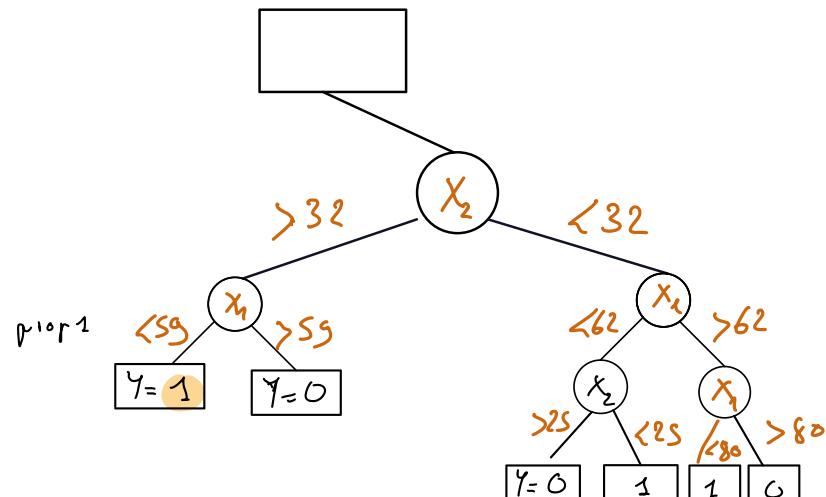
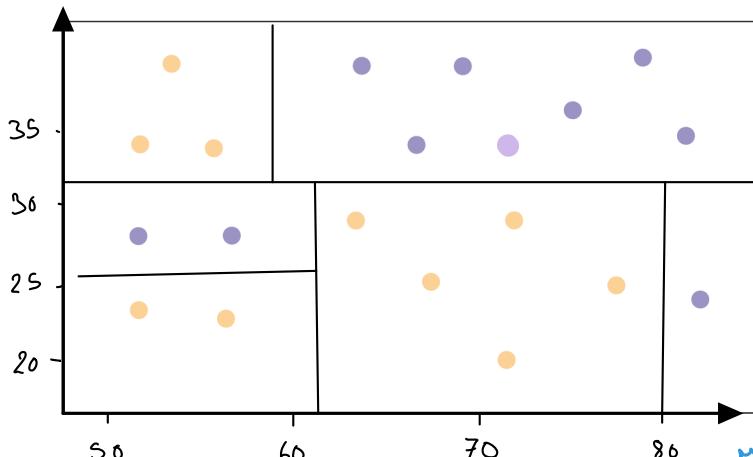
\* Rh: no given global criterion. ( $\neq$  to Lin reg, Logistic Reg, Deep Learning)

Instead : build a sequence of splits

# Building trees : CART

best split?

- $\gamma_i = 1$
- $\gamma_i = 0$



$$\text{measure: } \begin{cases} 0,3 \times 0,7 \\ 0,7 \times 0,3 \end{cases} \left\{ \begin{array}{l} \text{misclassif.} \\ 0,21 + 0,21 = 0,42 \end{array} \right. \quad \rightarrow 42\%$$

$$42\% \oplus 42\% \rightarrow 84\%$$

$$2 \left( \frac{2}{7} \times \frac{5}{7} \right) = \frac{20}{49} \approx 40\% \quad \left. \begin{array}{l} + \\ 87\% \end{array} \right\} 87\%$$

$$2 \left( \frac{5}{13} \times \frac{8}{13} \right) = \frac{80}{169} \approx 47\%$$

Stopping criteria

-

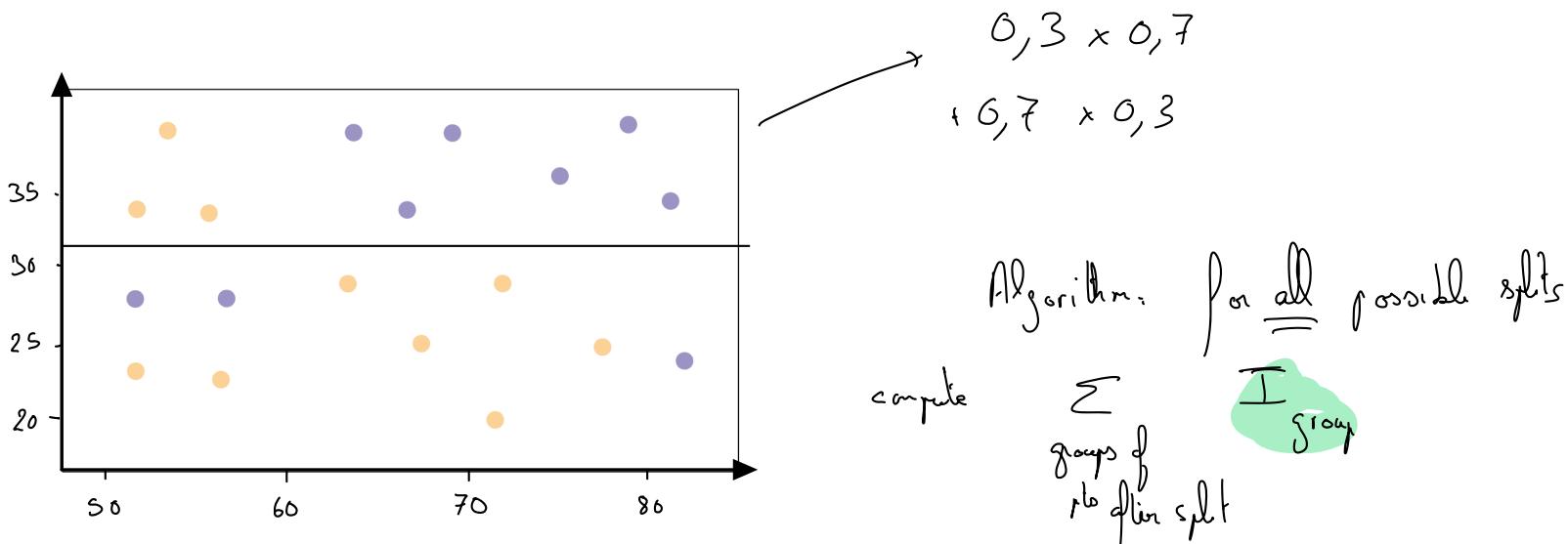
Complexity.

$\rightarrow d \times n$  inputs  
at each qt.

# splits  $\in n$

1 : perfect case

## Gini's impurity [classified : J classes]



Gini impurity measures **how often a randomly chosen element from the set** would be **incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset**.

$$I_G(p) = \sum_{i=1}^J p_i(1 - p_i) = 1 - \sum_{i=1}^J p_i^2$$

*J*  $\rightarrow$  *nb of classes*

For two classes, related to the variance if classified as Bernoulli r.v..

Group

$$\sum_{i=1}^3 p_i (z \cdot p_i)$$



50%



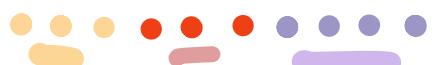
62%



18%



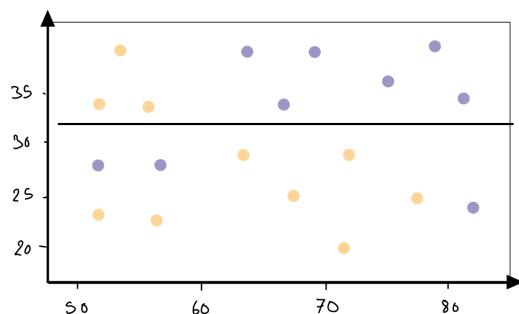
0%



$$\frac{3}{10} \times \frac{7}{10}$$

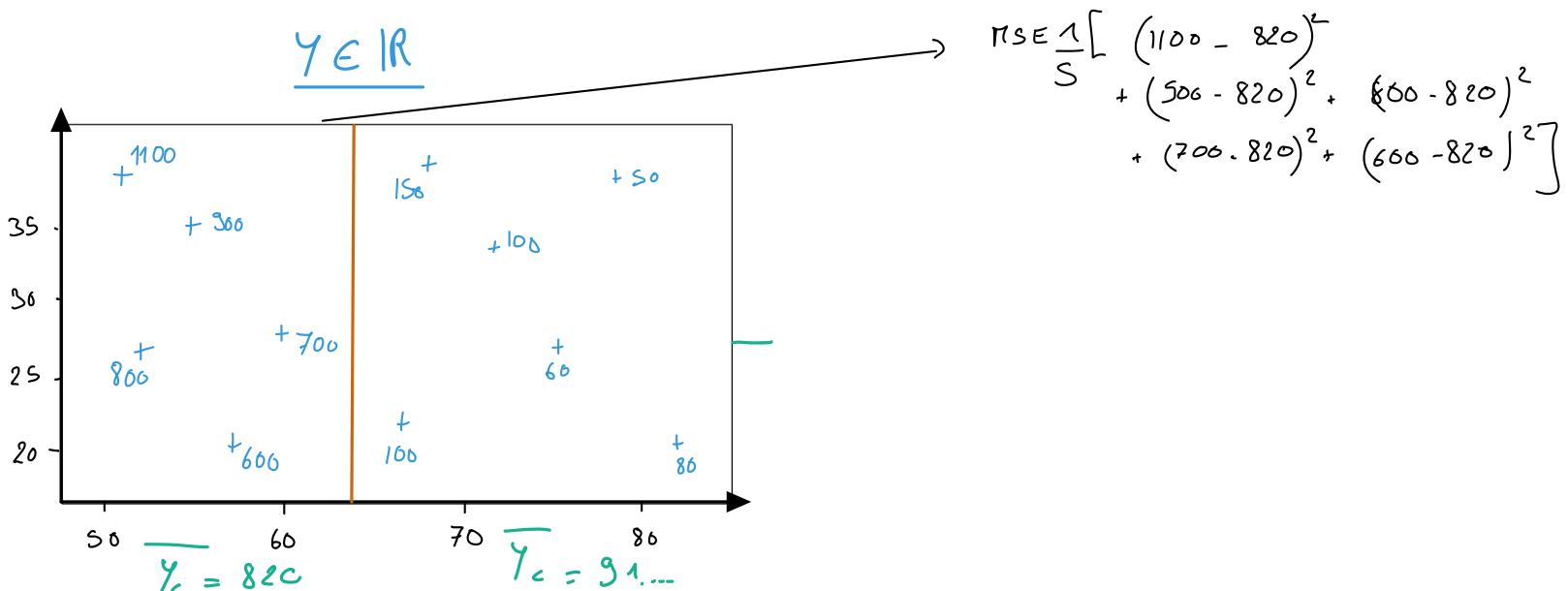
$$+ \frac{3}{10} \times \frac{7}{10}$$

$$+ \frac{5}{10} \times \frac{6}{10}$$



# Building trees : CART

## Regression case



How to chose the best cut ?

- e.g., Given a possible cut, measure the reduction of impurity :

- in regression: **mean-squared-error**  $\sum_{i \in C} (Y_i - \bar{Y}_C)^2$
- in classification: Gini's impurity.

$$\bar{Y}_C = \frac{\sum_{i \in C} Y_i}{\# C}$$

- Find dimension and threshold with optimal impurity reduction.
- Iterate until stopping criterion is met.

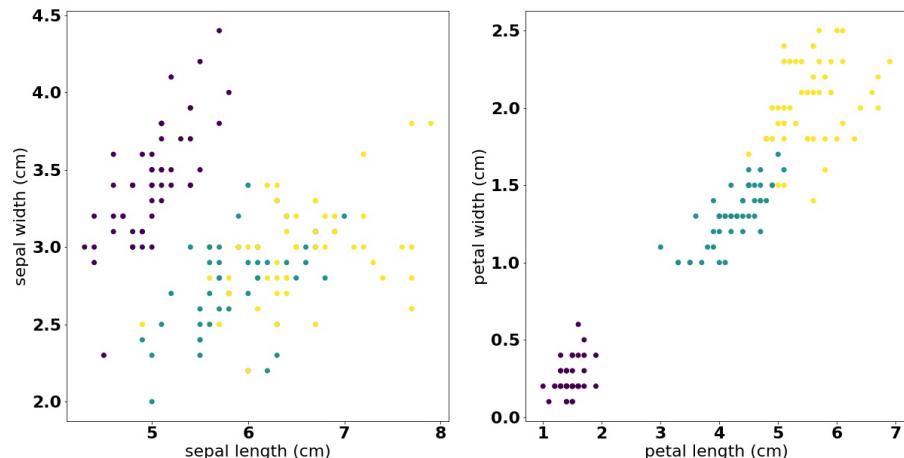
# Example: Iris Dataset



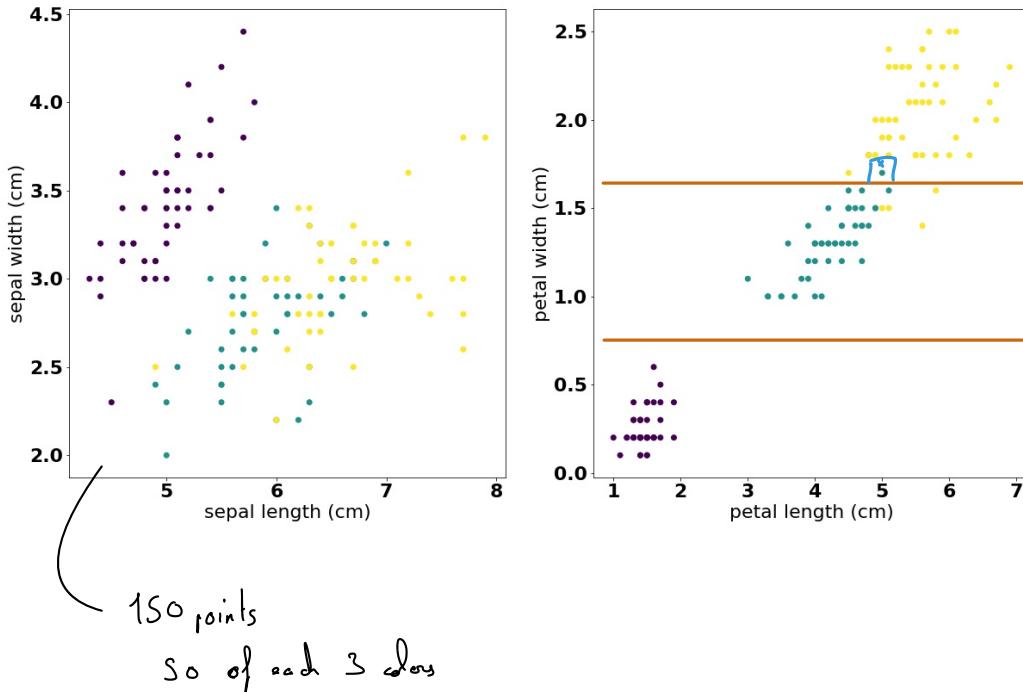
Figure: Iris: setosa, versicolor, virginica

One of the most widely used toy dataset in classification:

- 3 classes : `['setosa', 'versicolor', 'virginica']`
- 50 points per class.
- 4 features: `[ 'sepal length', 'sepal width', 'petal length', 'petal width' ]`



# Exercise: Guess Classification tree for Iris Dataset



# In practice : Scikit Learn

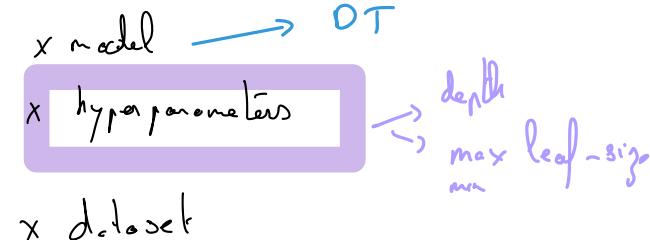
Universal Python library for Machine Learning:

- Very easy to use
- Very well documented
- Open Source

<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

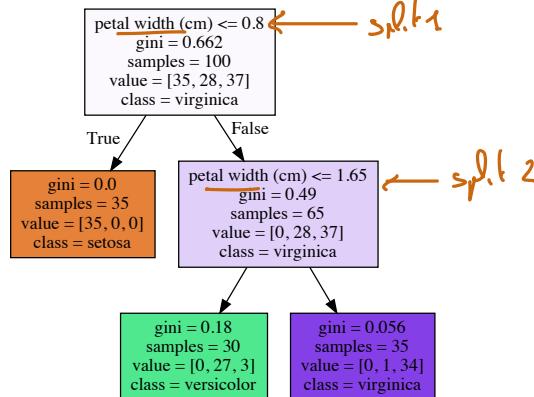
# Example: Iris Dataset

What do we need to specify ?



Output : Learned model  $\rightarrow$  DT .

Output:



```
from sklearn.datasets import load_iris
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.model_selection import train_test_split

# Choose the Learning algorithm
clf = DecisionTreeClassifier(max_depth = 2,
                             random_state=0)

# Load the dataset
iris = load_iris()
X_train, X_test, y_train, y_test = train_test_split(
    iris.data, iris.target, test_size=0.33, random_state=43)

# Check accuracy
cross_val_score(clf, iris.data, iris.target, cv=10)

# Fit the model
clf.fit(X_train,y_train)

# Plot tree
plot_tree(clf, feature_names = fn,
          class_names=cn,
          filled = True)
plt.show()
```

Example:

```
# Check
X_test[0,:], iris.target_names[y_test[0]]
```

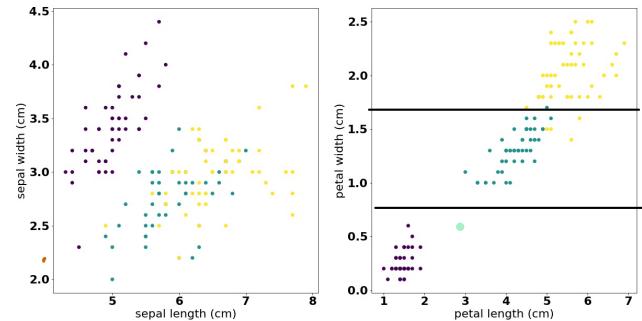
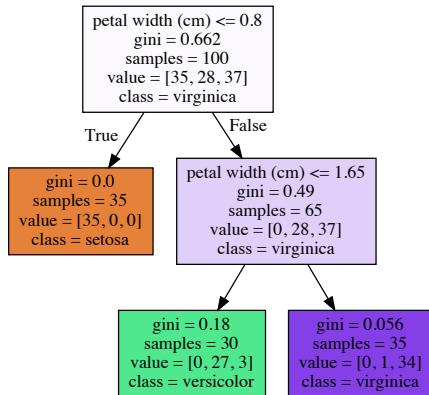
Returns:

```
(array([4.8, 3.1, 1.6, 0.2]), 'setosa')
```

Code: <https://colab.research.google.com/drive/1dXYIjacNAeASgl54vIWyWjMQlfgigSVA?usp=sharing#scrollTo=VeJHkJd7RsDw>

# Example: Iris Dataset

Output:



Example:

```
# Check  
X_test[0,:], iris.target_names[y_test[0]]
```

Returns:

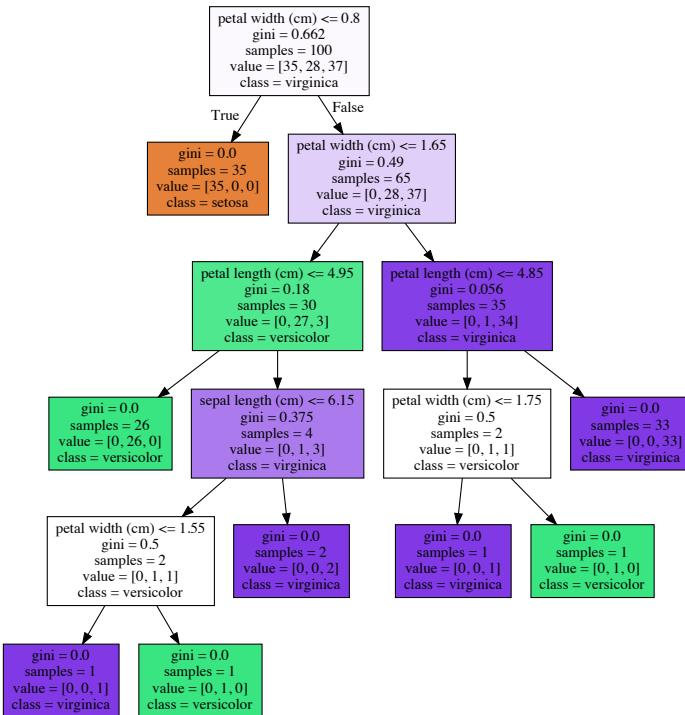
```
(array([4.8, 3.1, 1.6, 0.2]), 'setosa')
```

Code: <https://colab.research.google.com/drive/1dXYIjacNAeASgl54vIWyWjMQlfgigSVA?usp=sharing#scrollTo=VeJHkJd7RsDw>

What if I use a deeper tree ? Or if I do not specify the depth in the model definition ?

# Deeper trees

```
# Choose the Learning algorithm  
clf = DecisionTreeClassifier(random_state=0)
```



- When does the algorithm stop then ?
- What do you think about it ?

overfitting

# Stopping / Pruning

Trees that achieve perfect classification may suffer from **over-fitting**. (see example in the lab)

To avoid this problem, we can reduce the complexity of the trees:

## ① Stopping rules

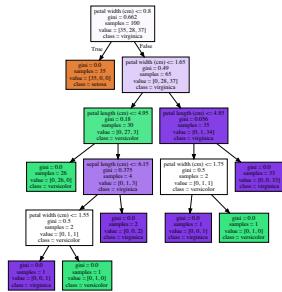
- ▶ Set a minimum number of samples inside each leaf.
- ▶ Set a maximum depth.

## ② Pruning

- ▶ Reduced error pruning.
- ▶ Cost complexity pruning.

# Pruning : example of Cost complexity pruning.

We create a sequence of Trees by changing one subtree at each step into a leaf, until we get only the root: the subtree is chosen to minimize the error made.<sup>1</sup> Then the test error is evaluated along the sequence, to choose the optimal pruning.



<sup>1</sup> [more details here](#)

# Cart: pros and cons

## Pros:

- Simple to understand, interpret, visualize
- Implicitly perform features/variable selection
- can handle both categorical and numerical data ✓
- little effort for data preparation ✓
- Non linear relationships do not affect performance ✓
- Can work with large number of observations. ✓

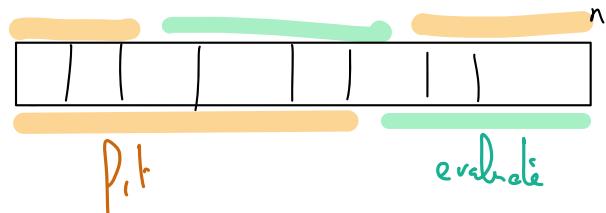
## Cons:

- ① Can create over-complex trees: overfitting ✗
  - ② Unstability: small variations of data can generate completely different trees ✗
  - ③ Cannot guarantee global optimal tree ✗
  - ④ Biased if some classes dominate ↘
- iterative algo - not  
a global solution

Solution : Using Random Forests !

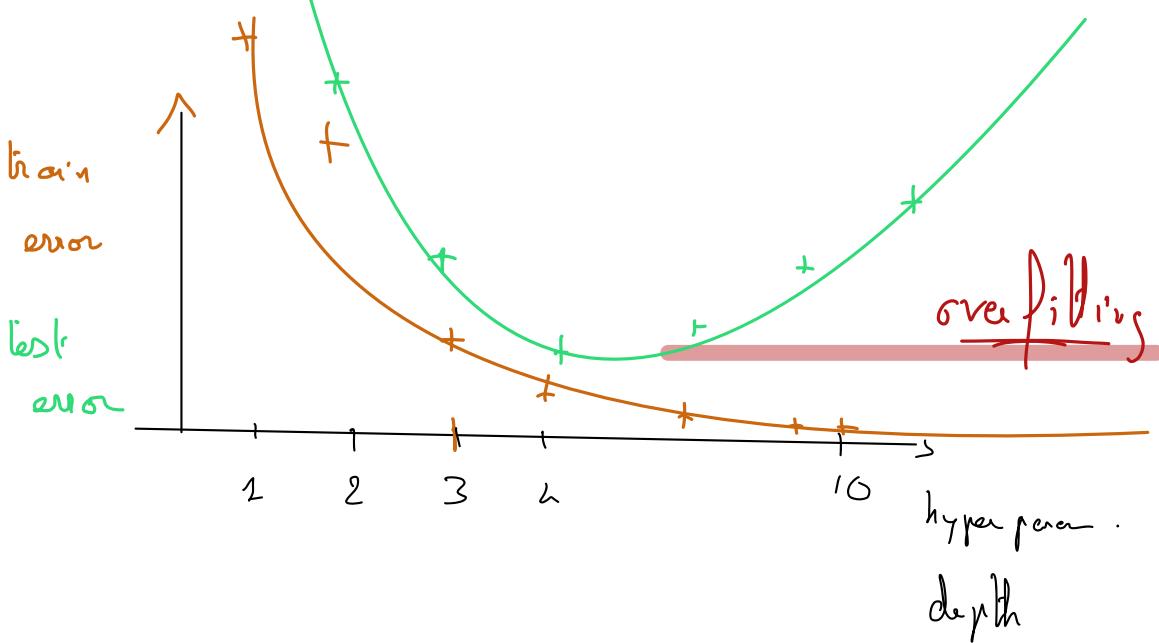
## Overfitting

Train  
validation



Test data

(train set)  
(validation set)



# Outline

1 Decision Trees

2 Random Forests

# Random Forests

To put it in simple words: random forests builds multiple decision trees and merges them together to get a more accurate and stable predictions.

→ key idea: create variable trees + aggregate them

How to create variability ?

- x sample the data
    - (x change hyperparam)
    - (x sample the classes ⌈ 3>>1)
  - x choose first condit° at random
- ← { sample at each step only some features  
use test }

# Random Forests

To put it in simple words: random forests builds multiple decision trees and merges them together to get a more accurate and stable predictions.  
→ key idea: create variable trees + aggregate them

How to create variability ?

- Instead of the most important features, search the best feature among a random subset (ntry)
- Work on subsets of data (bootstrap=True).
- More <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>

```
# Creating and fitting a Random Forest
from sklearn.ensemble import RandomForestClassifier

plt.figure(figsize=(20,10))
clf = RandomForestClassifier(n_estimators = 100,
                             max_depth=3, bootstrap = True,
                             random_state = 43)

clf.fit(X_train,y_train)
```

→ *Subsample features  
(-dired°)*

→ *nb  
trees*

→ *subsampled  
data*

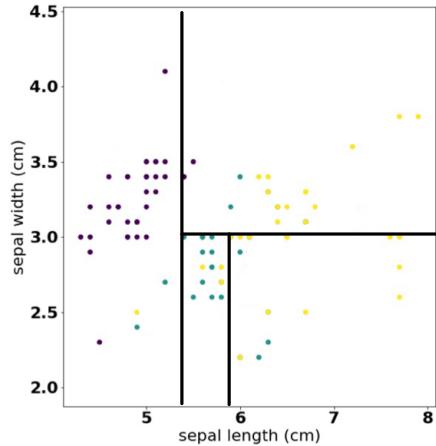
→ *hyperpar of the tree*

How to aggregate ? How to interpret a Random Forest ?

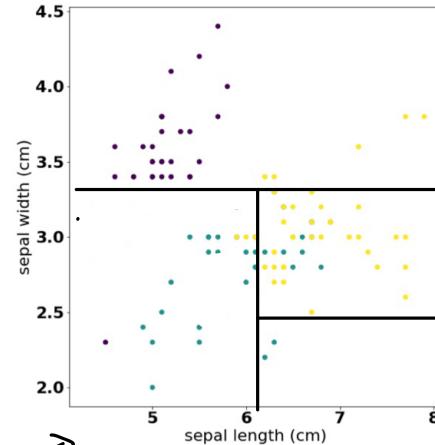
Bootstrap. = subsampling data

'Tree 1'

1



'Tree n° 2'

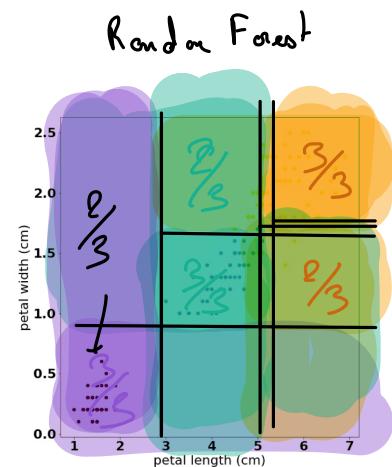
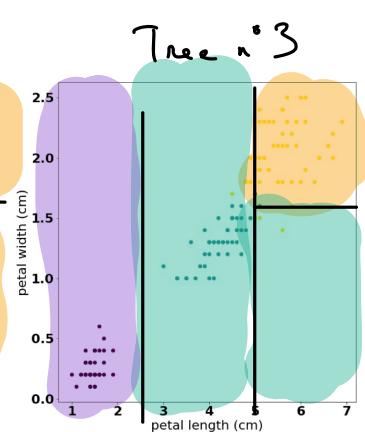
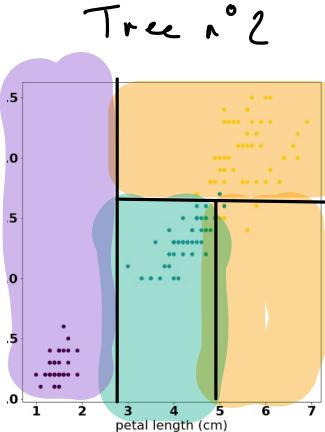
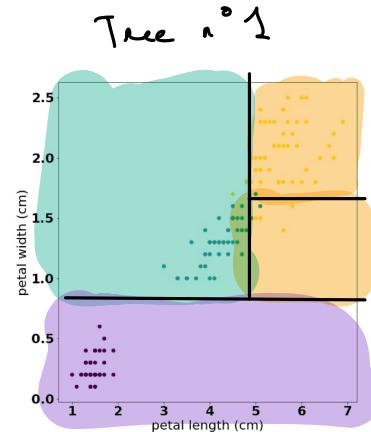


subsampled points



# Aggregation

We using a voting or averaging process !



# Aggregation

We using a voting or averaging process !

## Feature importance

- how much tree nodes using a particular feature reduce impurity along the trees.
- suggests feature selection rule: drop out features with low importance to avoid overfitting.
- Attribute **feature\_importance\_** in RandomForestRegressor of Sklearn.

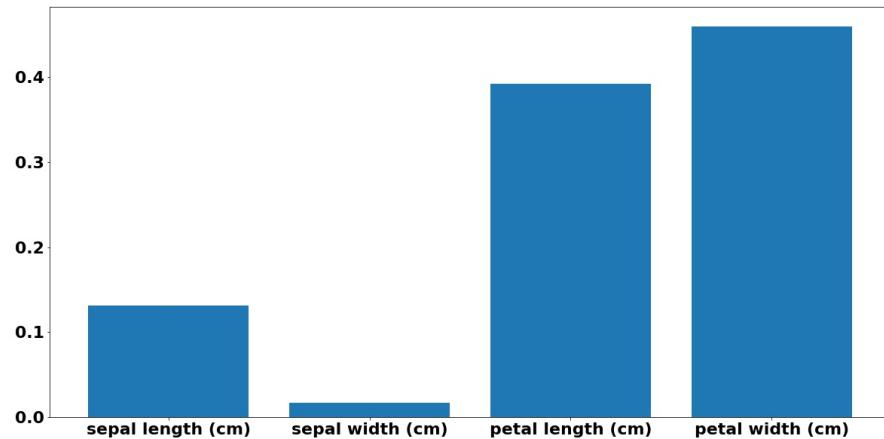


Figure: Feature importance output for a Random Forest Classifier in Python

# Comparison with decision trees

Cons:

- Less interpretable
- Slower to run

Pros:

- Better in higher dimension
- More stable outputs
- Feature selection

Improving predictions - more arguments.

- ① **n\_estimators** : number of trees in the forest
  - ▶ improves prediction / stability
  - ▶ slows down the algorithm
- ② **max\_features** : maximum number of features considered to split a node
- ③ **min\_sample\_leaf** : minimum number of samples that should remain in a leaf node.

# RF pros and cons

## RF Pros:

- ① works for classification and regression
- ② default hyperparameters often produce reasonable prediction -> easy to use.
- ③ avoid overfitting if some good trees in the forest and easy feature selection -> high dimensional pb.
- ④ hard to beat in performance.
- ⑤ Bonus : supports multiple outputs!

## RF Cons:

- ① ~~(Fast to train)~~ but slow to provide new predictions -> ineffective for real-time predictions
- ② Good for prediction but bad for description.

What about Regression ?

# Questions?

# Boosting

Another import technique (**very powerful !**)

Main idea:

- Give more importance to difficult point iteratively
- Incrementally building an ensemble by training each new instance to emphasize the training instances previously mis-modeled.

Example: AdaBoost

See: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostClassifier.html>

## Bonus: Multiple Outputs

- **Goal:** predict several outputs simultaneously
- **Solution:**
  - ▶ each leaf contains a value for each output
  - ▶ to chose the splits, we use a (weighted) average of the impurity of each output.

Face completion with multi-output estimators



## Bonus: complexity

**Total complexity for one decision tree:**

Worst case:

$$O(n^2 d)$$

Balanced case:

$$O(nd)$$

# Tips

- Decision trees tend to overfit: limit the depth or use `min_samples_leaf` (5 is a good initial pick)
- Visualize the tree with a small depth first
- Balance your dataset: trees tend to be biased towards dominating class.

# Conclusion

- ① Trees are one of the simplest method (the most intuitive / human like)
- ② Random Forests give excellent results in many applications.

## Lab :

- Example on a synthetic dataset of time series
- Application to inflation prediction in Brazil.

