

Decision Trees and Random Forests Data science Certificate

Aymeric DIEULEVEUT

May 2021

- why ?
- what they are used for?
- how is it built → training / algorithm.
- how to implement a decision tree in python.
- pros / cons / tips.

1 Decision Trees

2 Random Forests

Outline

1 Decision Trees

2 Random Forests

A Reminder on Supervised Learning

we have access to :

$$(x_i, y_i) \text{ on the training set}$$

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

goal: find a fund^o f.

$$f: X \longrightarrow Y$$

$$x_{\text{new}} \longrightarrow f(x_{\text{new}})$$

3 Examples:

- Iris Dataset: Y is 3 possible classes

- MNIST: Y is $\{0, \dots, 9\}$

- House price prediction $Y \in \mathbb{R}_+$ $y_i \geq 0$

Example: Iris Dataset

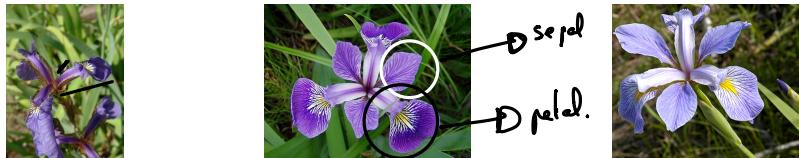


Figure: Iris: setosa, versicolor, virginica

One of the most widely used toy dataset in classification:

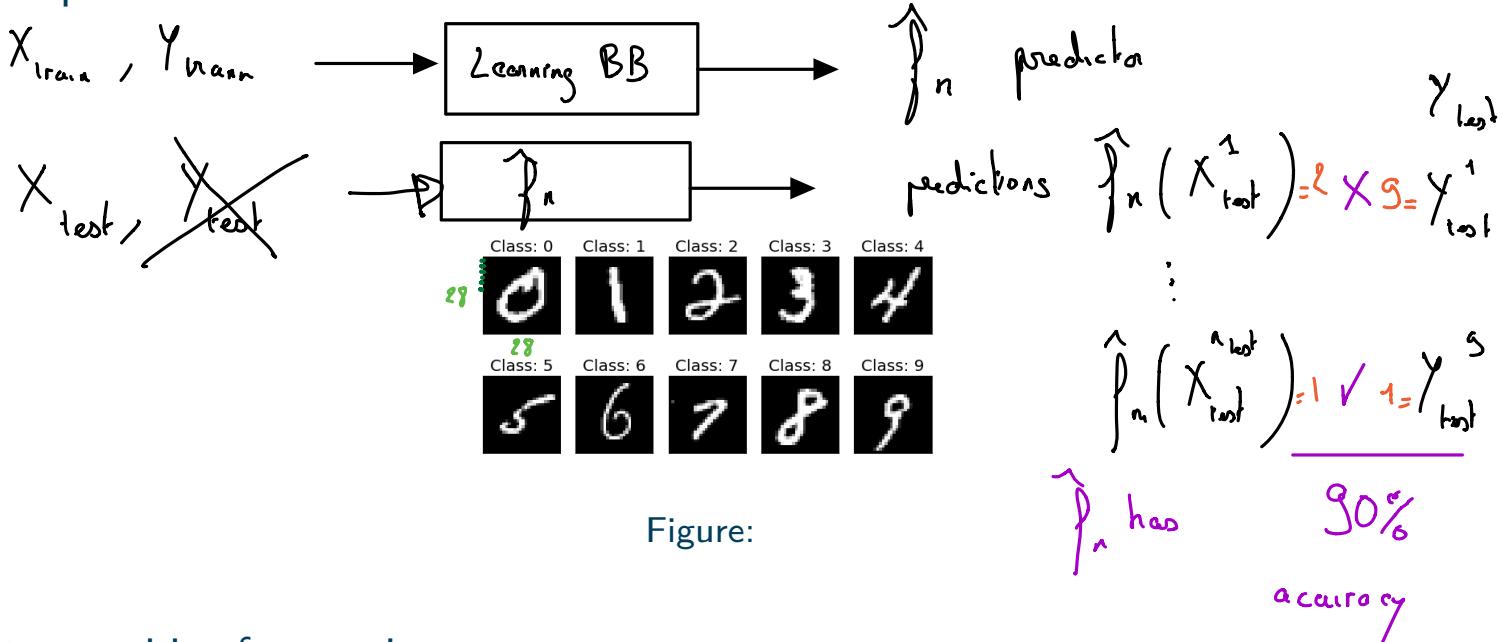
- 3 classes : ['setosa', 'versicolor', 'virginica'] $\longrightarrow \gamma_i \in \{0, 1, 2\}$

• 50 points per class. $n = 3 \times 50 = 150 \text{ obs}$

- 4 features: ['sepal length', 'sepal width', 'petal length', 'petal width']

$$\chi_i = \left[\begin{array}{c} \end{array} \right]$$

Example: MNIST Dataset



Digit recognition from an image:

- 10 classes : $Y_i \in \{0, \dots, 9\}$
- 60k images, (50k train, 10k test), balanced as many 0 than 1, 2 etc
- 784 features.

28×28 .

grey intensity per pixel.

$$X = (0, 0, \dots, 120, 256, 256, 0, \dots,) \in \mathbb{R}^{784}$$

House price Prediction

$X_1(\underline{\text{Paris}}, \underline{\text{Pd}}, \underline{\text{Vas}})$...
 $X_2(\underline{\text{Pd}}, \underline{\text{...}})$...
 $X_3(\underline{\text{Vas}}, \underline{1})$

Per	Pd	Vas	...
1	0	0	...
0	1	0	...
0	0	1	...

trick: one-hot-encoding.
get-dummies (random)
drop first = True.

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

quantitative
{ Paris, Palaisse, Vers. }
 \mathbb{R}_+

Iris: } all features were quantitative

$$d^+ \text{ columns} = 4 + 3$$

Mnist:

{0, 1}

linear link between features

$$X = \left(\begin{array}{|c|} \hline \text{---} \\ \hline \end{array} \right) \quad \text{n Rows}$$

drop 1 to avoid

linearly dep features!!

Example about linearly dependent features:

$$Y = \text{Salary} . = h \in$$

$$X = \text{height}$$

$$Y \approx 20 + \underbrace{(X - 1m)}_{\text{one feature}} \times \frac{10}{\alpha_1}$$

$$10 \times h_m + 0 \cdot h_{\text{inches}} = 0 \cdot h_m + 300 \cdot h_{\text{inches}}$$

Second character. / feature

$$10 \times (\text{height m})$$

$$\Leftrightarrow 300 (\text{height inches})$$

$$X = (\text{height in m}, \text{height in inches})$$

$$\alpha_1 (\text{height m})$$

$$\hat{\alpha} = (X^T X)^{-1} X Y$$

does not exist

$$+ \alpha_2 (\text{height inches})$$

Summary of those three datasets

	Iris	MNIST	House price prediction
Number of examples n	156	66k	1500
Number of the features d	4	78h	81
Type of the features	quant	quant	quant + qual.
Space \mathcal{Y}	$\{0, 1, 2\}$	$\{0, \dots, 9\}$	\mathbb{R}_+
Task	Classification.		Regression task.

$$\hookrightarrow \mathcal{Y} = \{0, \dots, k\}$$

Decision Tree

| also trees in unsupervised
| but they are ≠ !

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 ?

Age BMI

if age is $70+$ $\rightarrow Y=1$,

if BMI is $35+$ $\rightarrow Y=1$

Y is severe candidiasis

$$Y = 1 \quad 0,5 < \frac{1}{40} (\text{age} - 50) + \frac{1}{20} (\text{BMI} - 25)$$

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)
- How to predict with a given decision tree? → inference phase

CART .
↑ —

Inference Phase

$$Y = \{0, 1\} \quad \text{at risk COVID 19}$$

$$X = (\text{age}, \text{Bmi}, \text{"have-precondit"})$$

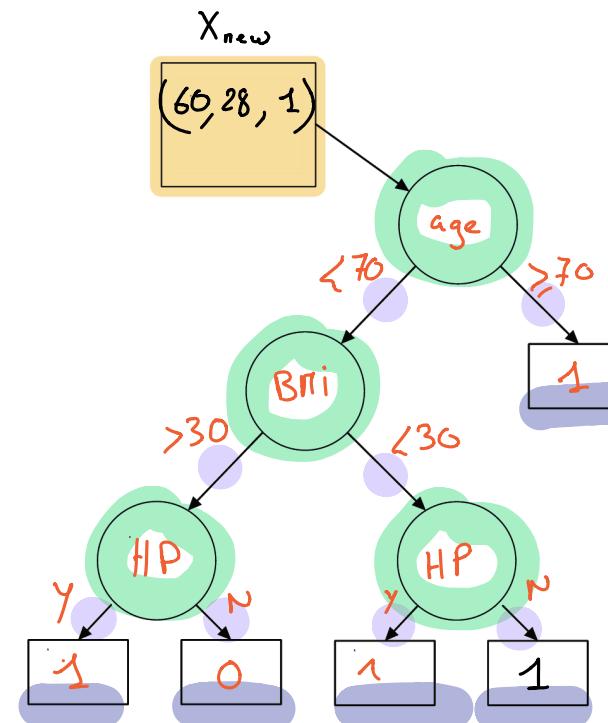
- Root = observation = input

- Nodes = tests

- Branches = possible outcomes

- Leaves = final decision = output

orange = result of the training phase
↳ latent



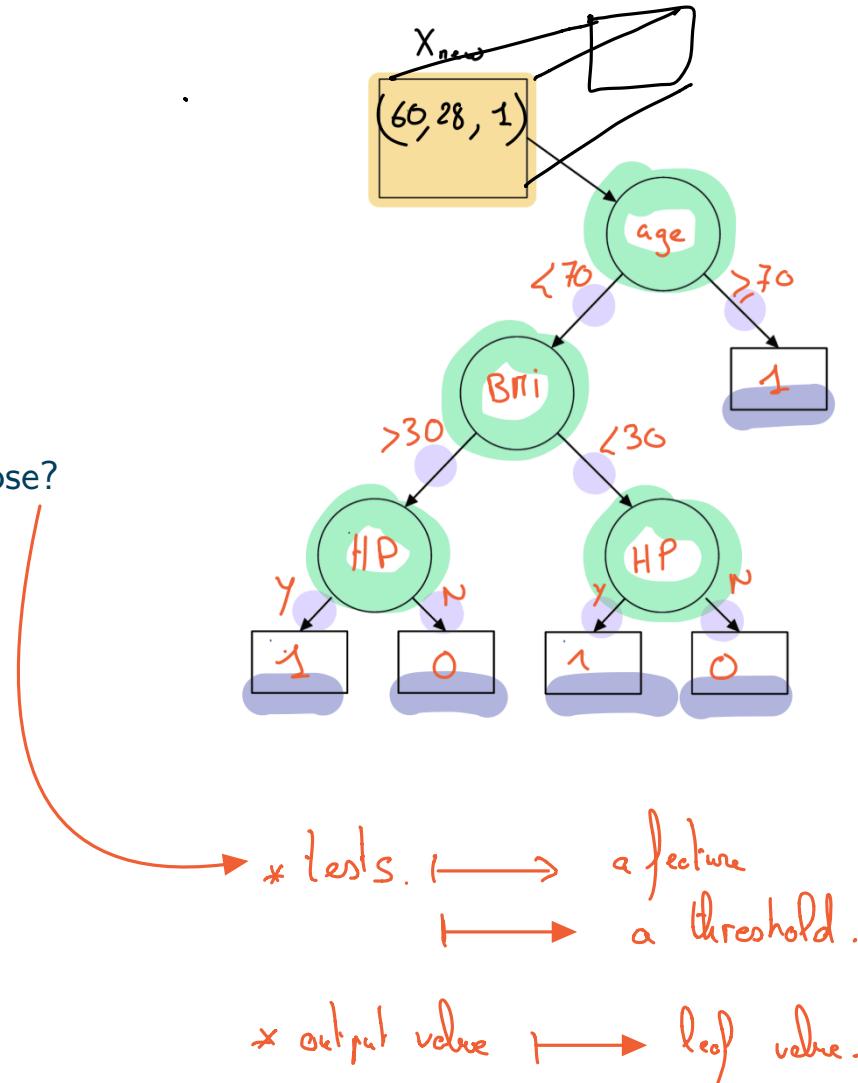
already trained

Training: how to build a tree.

Key question !

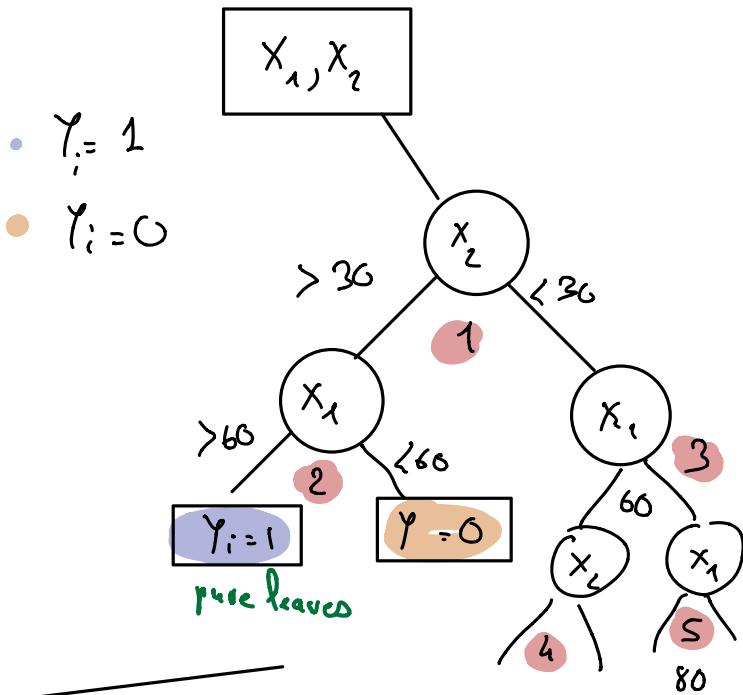
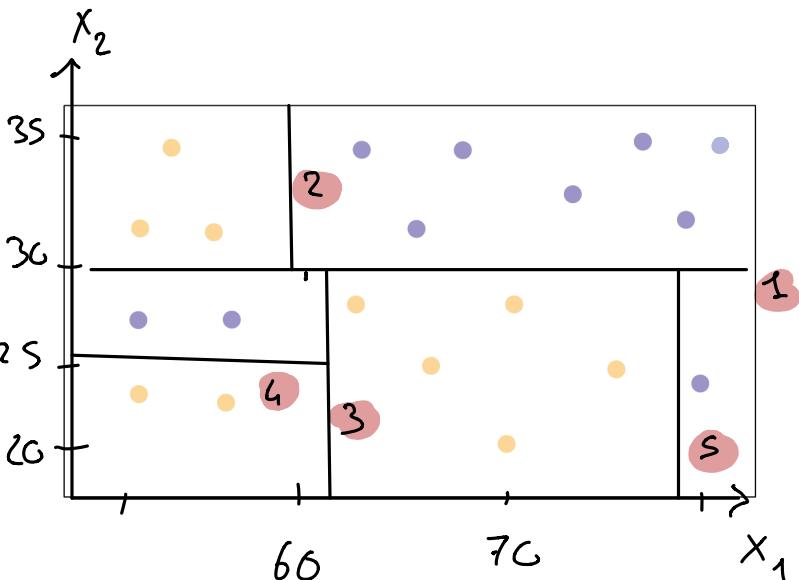
- What is a good tree?
- What do I need to choose?
- How to choose?

What do you think?



Building trees : CART

Restart at 11h30



Split 1 X_1 vs 60 $\rightarrow 45^\circ$ ↗
Split 1 bis X_1 vs 30 $\rightarrow 48\%$

To do: choose a quality criterion.

Gini's impurity

Classification.

group 1



$$\frac{5}{10} \times \frac{5}{10} + \dots = \frac{50}{100}$$

Gini : in Python, default choice!

group 2



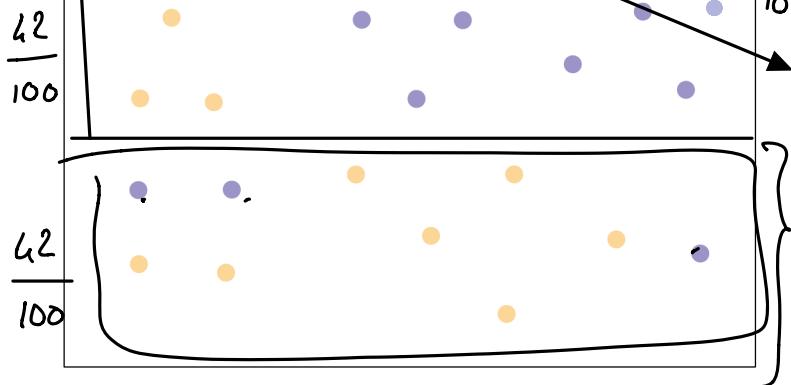
$$\frac{3}{10} \times \frac{7}{10} + \frac{7}{10} \times \frac{3}{10} \rightarrow \frac{42}{100}$$

group 3



$$\frac{1}{10} \times \frac{9}{10} + \frac{9}{10} \times \frac{1}{10} \rightarrow \frac{18}{100}$$

group 4



$$\sum_{i=1}^2 = \frac{3}{10} \times \left(\frac{7}{10} \right) = \frac{21}{100}$$

$$+ \frac{7}{10} \left(\frac{3}{10} \right) = \frac{21}{100}$$

Blue pts

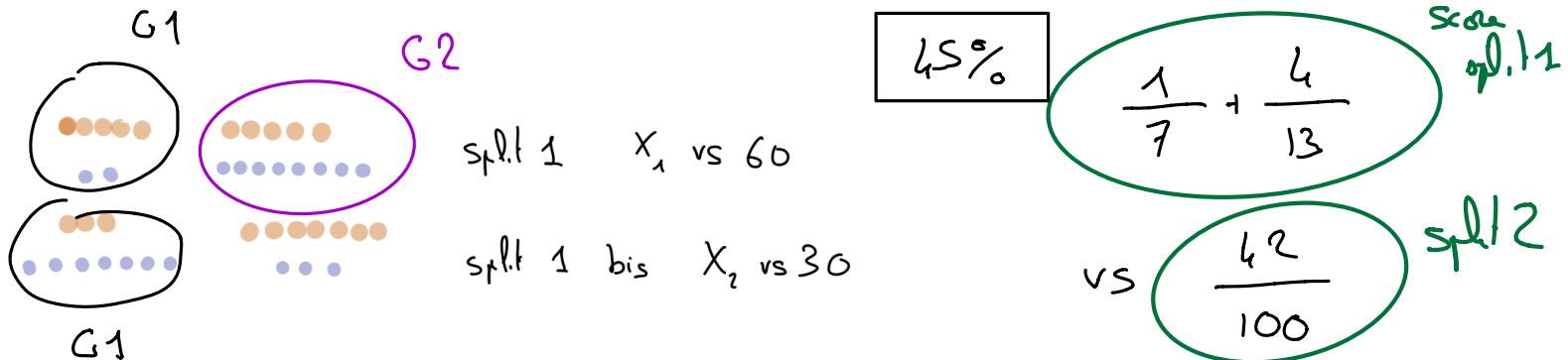
Gini impurity measures **how often a randomly chosen element from the set** would be **incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset.**

$$I_G(p) = \sum_{i=1}^J p_i (1 - p_i) \rightarrow \begin{array}{l} \text{# of classes.} \\ \text{prob of misclf.} \\ (\text{this point}) \end{array}$$

$$= 1 - \sum_{i=1}^J p_i^2$$

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

prob of picking class i



Best criterion of the end

Split 1

$G_1(2, S)$

$$2 \left(\frac{2}{7} \times \frac{5}{7} \right) = \frac{20}{49}$$

$G_2(8, 5)$

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

pls

7

13

$$\frac{7 \times 20}{20 \times 49} = \frac{1}{7}$$

$$\frac{13 \times 80}{20 \times 169} = \frac{4}{13}$$

Split 2:

$G_1 \rightarrow$

Gini's capacity

pls

10

Average

10

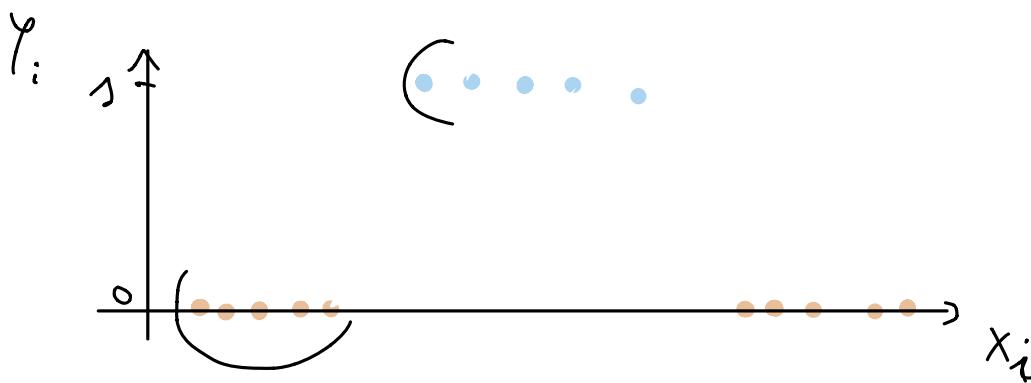
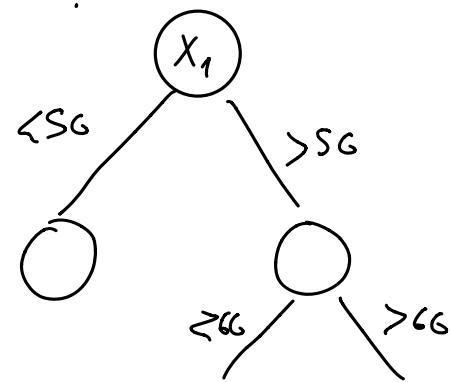
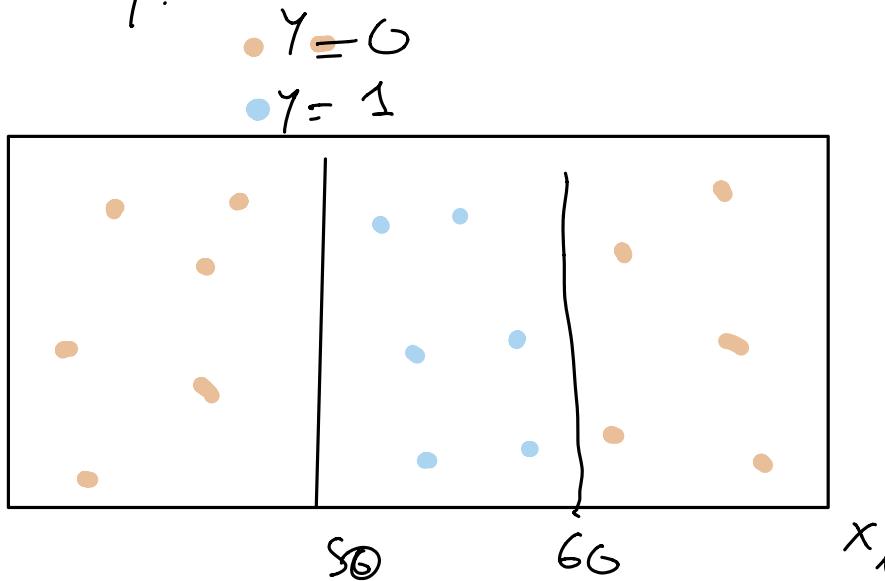
$G_2 \rightarrow$

$$\frac{62}{100}$$

10

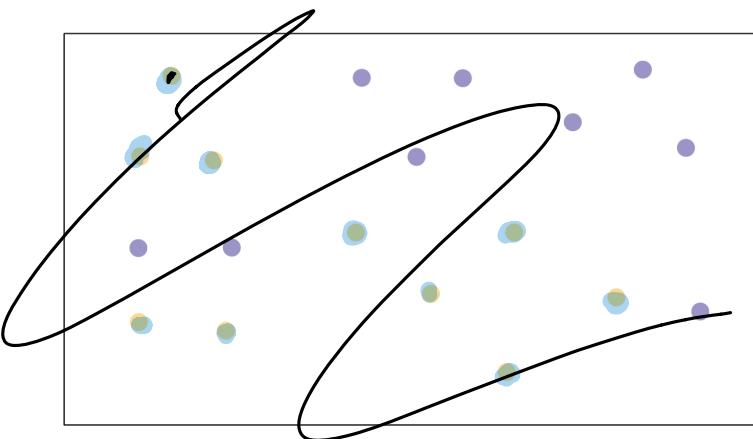
Score
P_a
spl.2.

Summary.

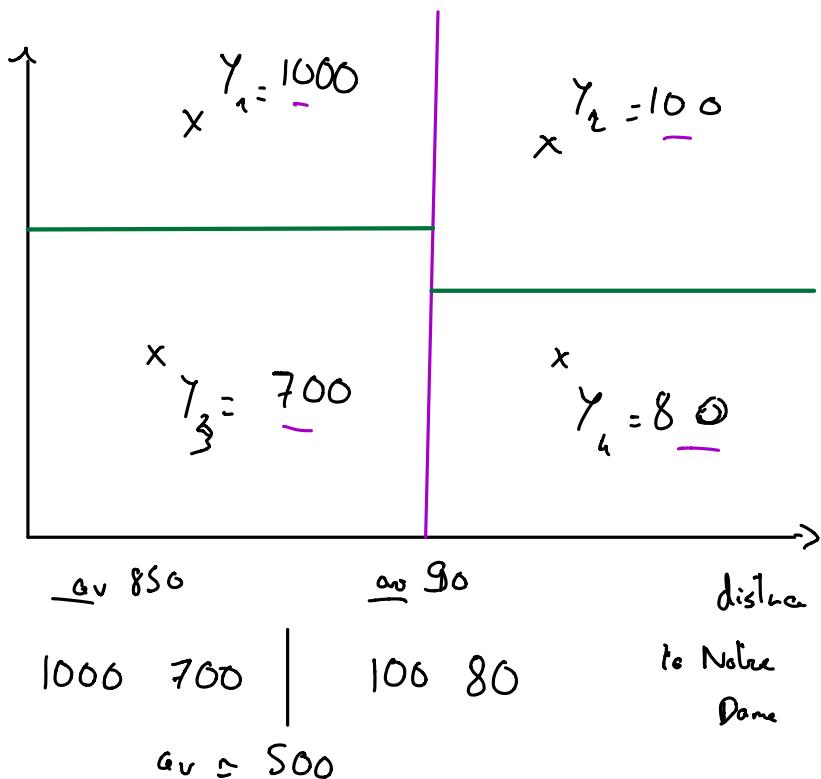


Yesterday : linear models

Building trees : CART



Surface



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

Learn

until "pure leaves" ($Gini = 0$). \rightarrow until reduction in $Gini$ is small
 until a maximal depth.
 until a minimal nb of pts per leaf.

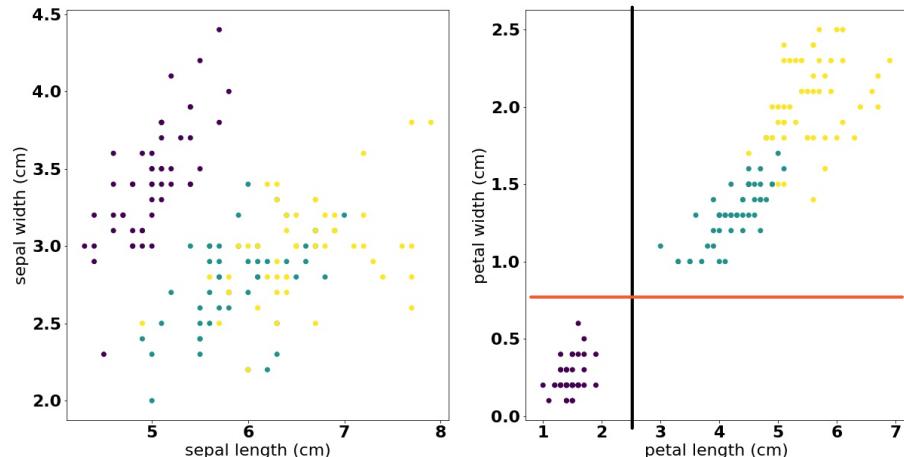
Example: Iris Dataset



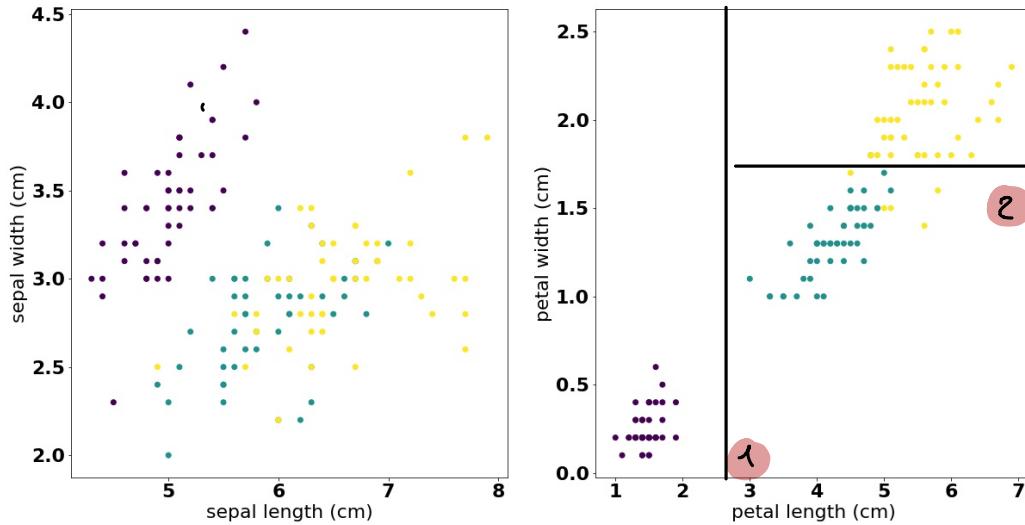
Figure: Iris: setosa, versicolor, virginica

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 ?

1: load data

2: choose parameters of the tree:

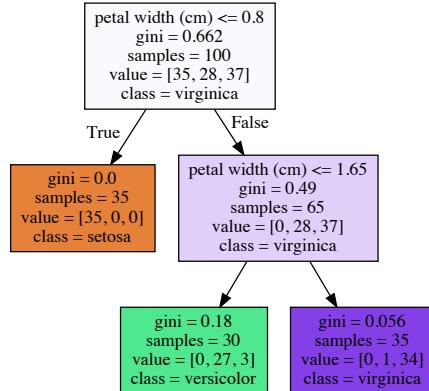
- * impurity : Gini is default

- * stopping criterion e.g. depth

3: Learn.

clf . fit (X_train , Y_train)

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]]
```

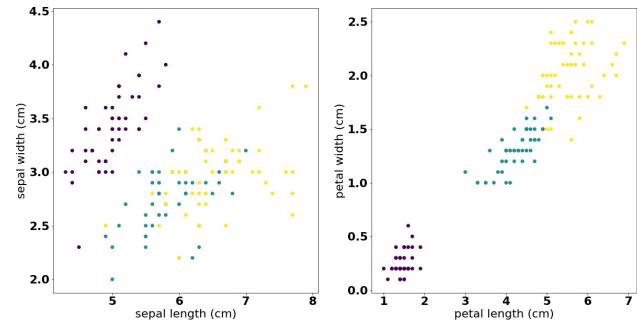
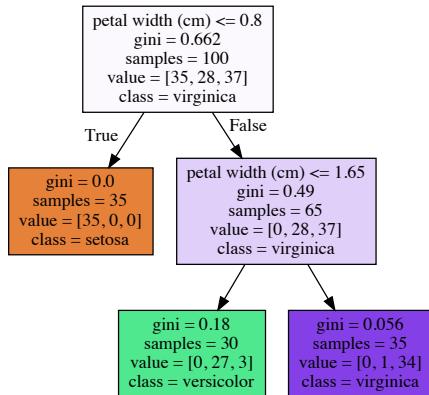
Returns:

(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/1dXYIjacNAeASgl54vIWjMQlfgigSVA?usp=sharing#scrollTo=VeJHkJd7RsDw>

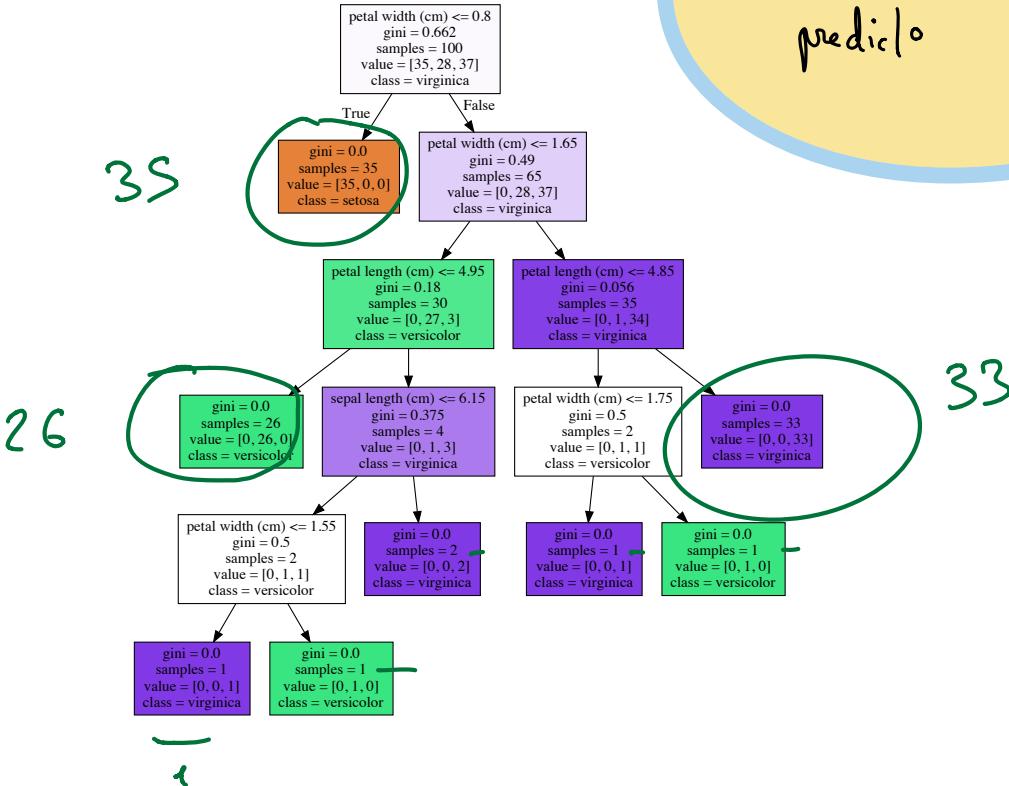
What if I use a deeper tree ? Or if I do not specify the depth in the model definition ?

Deeper trees

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

test - performance?
validation - predict

choose the
best depth.



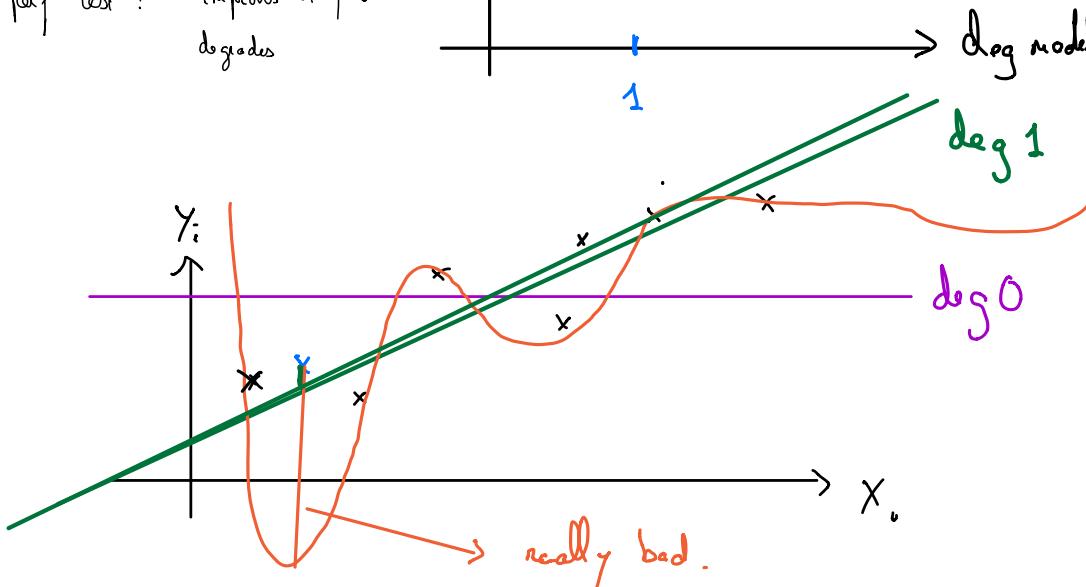
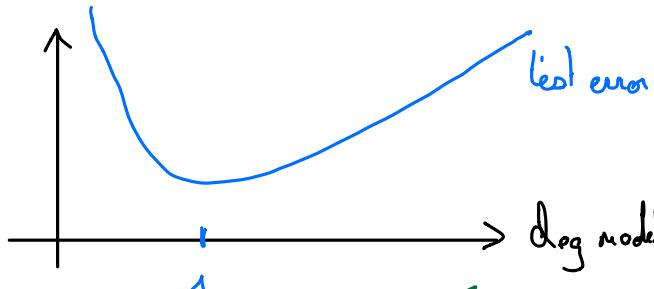
- When does the algorithm stop then ?
- What do you think about it ?

Overfitting: cross validatio
to choose parameters depth, etc.

When I increase the complexity of model.
I learn

perf train improves

perf test : improves at first
degrades



models depend on
the degree of the poly

deg 6 .

x train set .

x test point



easy.

sets

depth = 3

fit on train

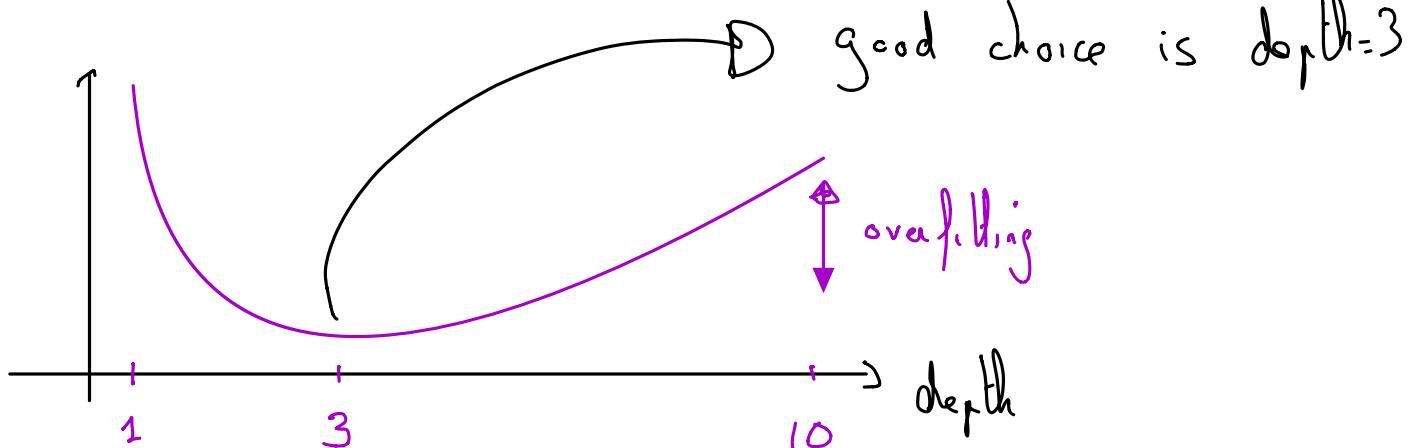
evaluate val/test error.

add validation



Several train/val splits.

depth = 10



Stopping / Pruning

on the train set !!

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

→ poor perf } on the test

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

① Stopping rules

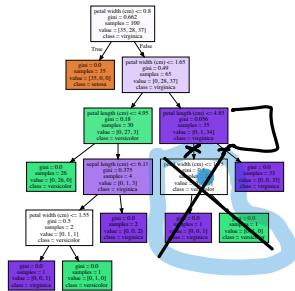
- ▶ Set a minimum number of samples inside each leaf.
- ▶ Set a maximum depth.

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.¹ Then the test error is evaluated along the sequence, to choose the optimal pruning.



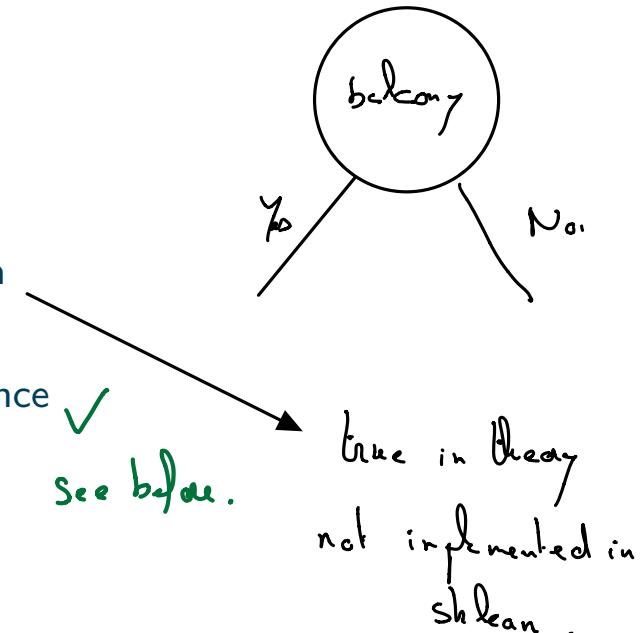
¹ [more details here](#)

Cart: pros and cons

8pm

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

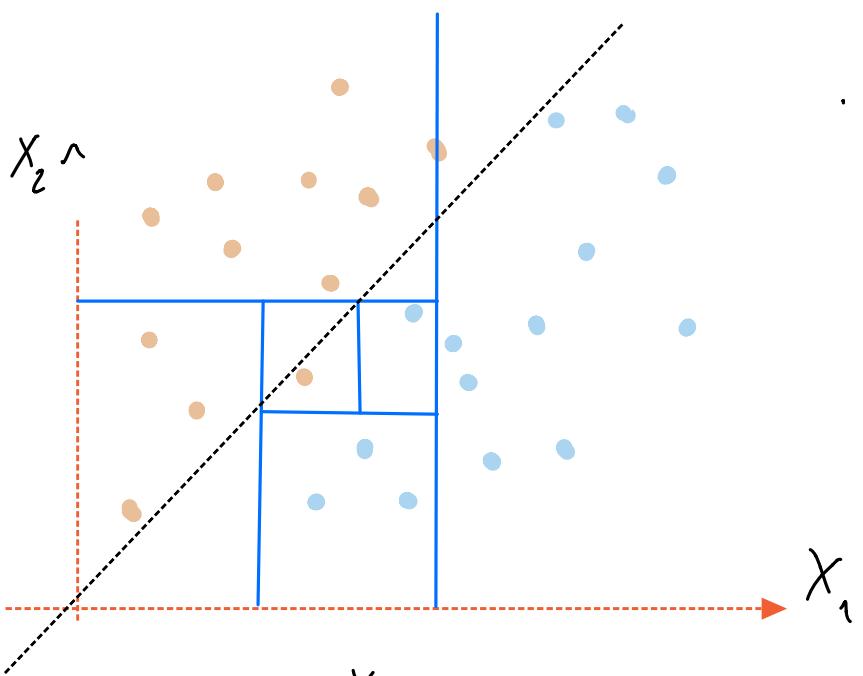
→ use one hot encoding!

first split.
d possible features
n possible thresholds.

Solution : Using Random Forests !

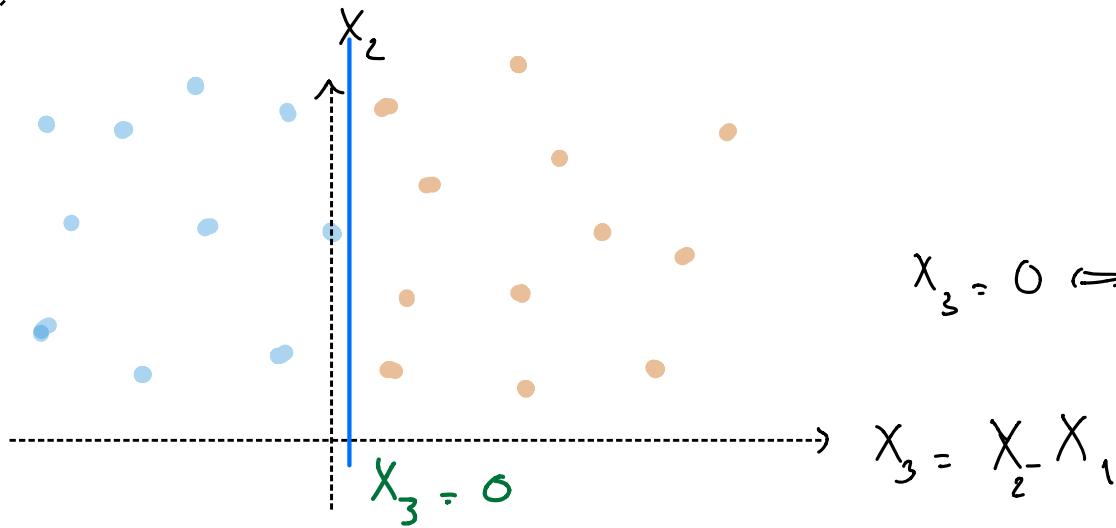
Creating new features can drastically simplify the tree!

Original tree



With
a new
feature

$$X_3 = X_2 - X_1$$



$$X_3 = 0 \Leftrightarrow X_1 = X_2$$

$$X_3 = X_2 - X_1$$

Outline

1 Decision Trees

2 Random Forests

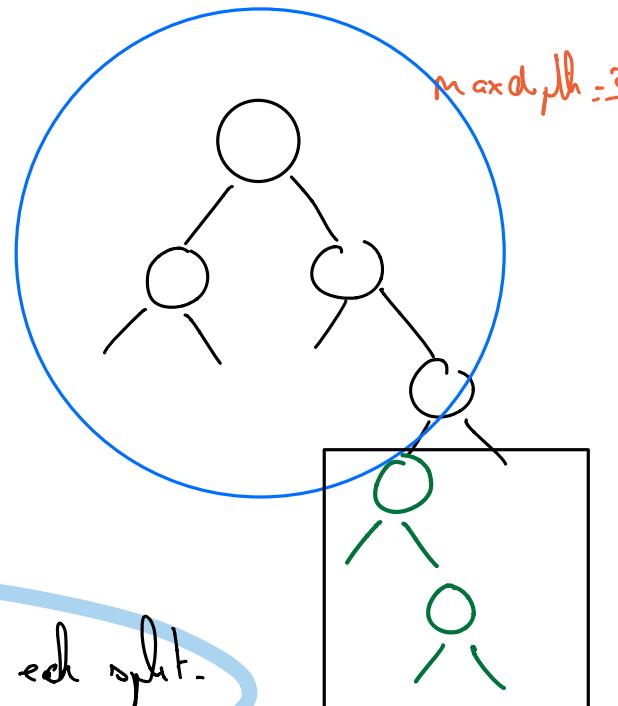
Random Forests

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 ?

- * different data !
- * change parameters
 - * max_depth
- * change purity criterion
- * allow only some features at each split.



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

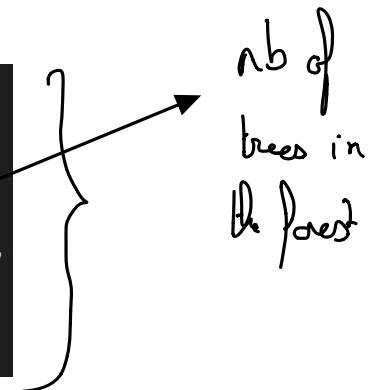
How to create variability ?

- Instead of the **most important features**, search the best feature among a random **subset (ntry)**. *ntry among d possible.*
- 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 = 0)

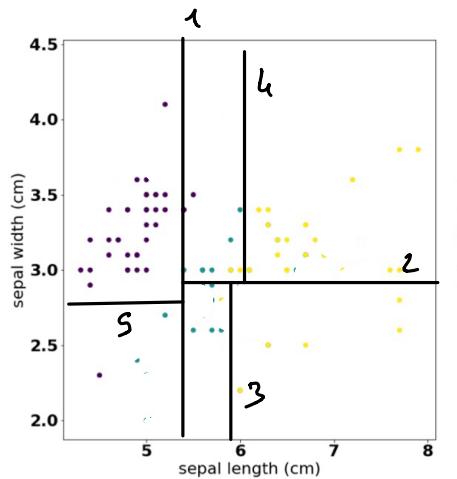
clf.fit(X_train,y_train)
```



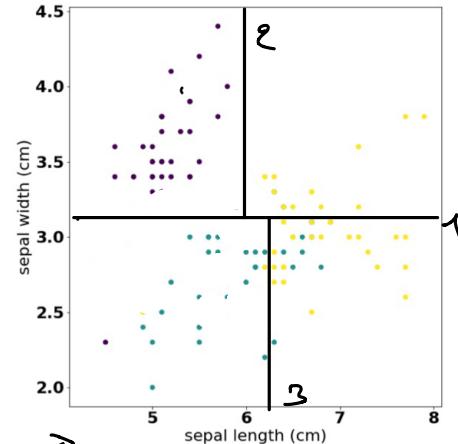
How to aggregate ? How to interpret a Random Forest ?

Bootstrap.

'Tree 1



Tree n° 2

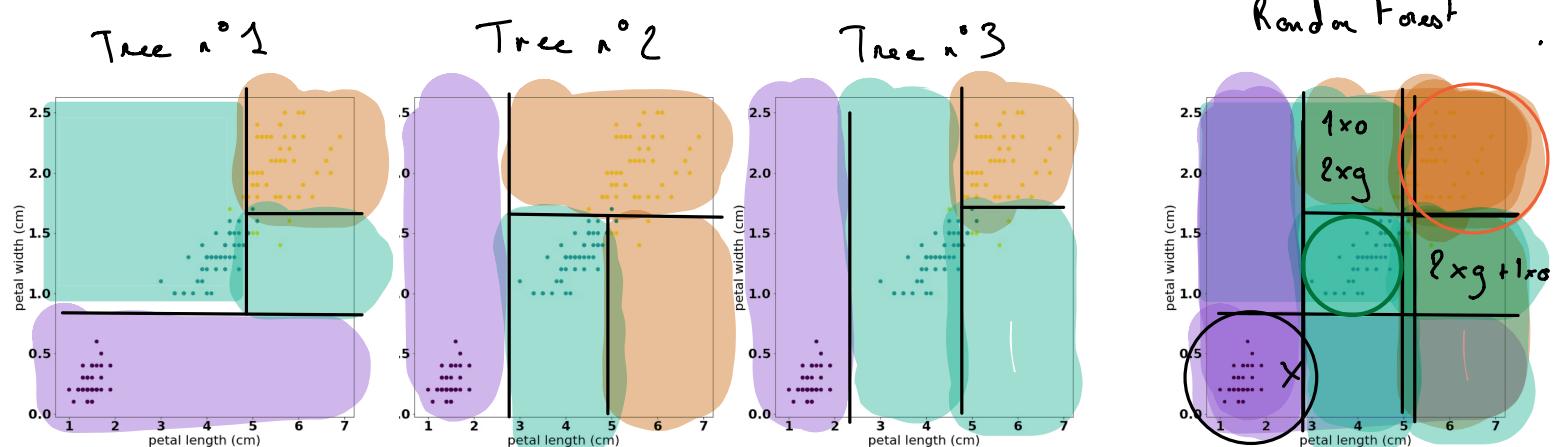


↑
subsampled points

→ I randomly delete points in the dataset
for each tree!

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.

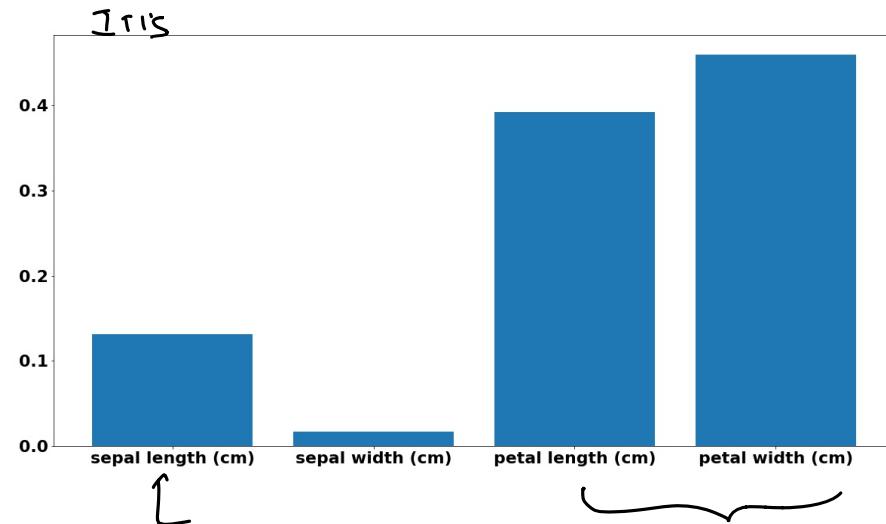


Figure: Feature importance output for a Random Forest Classifier in Python

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.
max_depth

arguments



RF pros and cons

$$n_{\text{pts.}} = 100 \text{ pts.}$$
$$\text{depth} : 7 \rightarrow 2^7 \approx 128$$



$$n_{\text{est}} = 100$$
$$\text{depth} = 3 \rightarrow \log n$$
$$\text{boosting} = \text{True}$$

RF Pros:

how many samples per leaf?

- 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. → s.o.t.a.
- ⑤ Bonus : supports multiple outputs!

$$10^6 \rightarrow \text{depth} \ll 20$$
$$\approx 10$$

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.

↳ inference time

have to compute n_{est} predictions

What about Regression ?

↳ aggregation is by averaging instead of voting !

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

XG_Boost

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.

- 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.