

Decision Trees and Random Forests

Data science Certificate

Aymeric DIEULEVEUT

May 2023

x task

1 A Reminder on Supervised Learning : x examples

x why we choose trees

2 Decision Trees

x how we build trees

Choose

hyperparameters

Learned
model

x how we predict with them

3 Random Forests

Outline

1 A Reminder on Supervised Learning

2 Decision Trees

3 Random Forests

A Reminder on Supervised Learning

Task: predict Y from X .

given $(X_i, Y_i)_{i=1 \dots n}$ observed

inputs / Features

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

3 Examples:

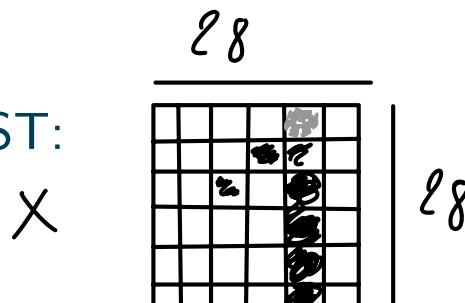
output

- Iris Dataset: Predict species

from 4 Pictures :

$\begin{cases} \text{petal length} \\ \text{width} \\ \text{sepal} \\ \text{"} \\ \text{"} \end{cases}$

- MNIST:



$$28^2 = 784 \text{ pixels}$$

- House price prediction :

predict the price of a house

$Y = \text{house price}$

$\rightarrow \in \mathbb{R}_+$: REGRESSION

$X = \{ \text{surface (m}^2\text{)}, \text{neighborhood}, \text{garden}, \text{floor, construct year, ...} \}$

$Y = \{ 0, \dots, 9 \}$

Finite number
of CLASSES
→ CLASSIFICATION

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. . $n = 150$
- 4 features: `['sepal length', 'sepal width', 'petal length', 'petal width']`

Example: MNIST Dataset

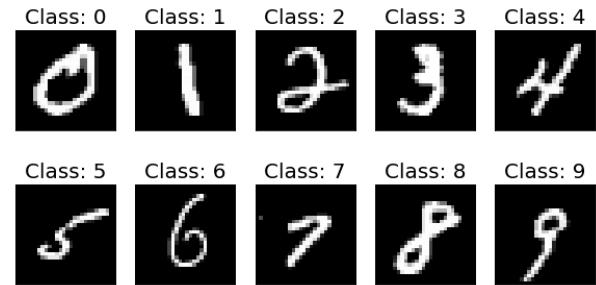


Figure:

Digit recognition from an image:

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

House price Prediction

Famous haggle.

X_1	100	, 1 st	, Yes	1M
X_2	10	, 13 th	, No	200k
X_3	200	, 5 th	, Yes	2M.

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

Simplified version when $X = \{ \text{Surface}, \text{Neigh}, \text{Garden} \}$.

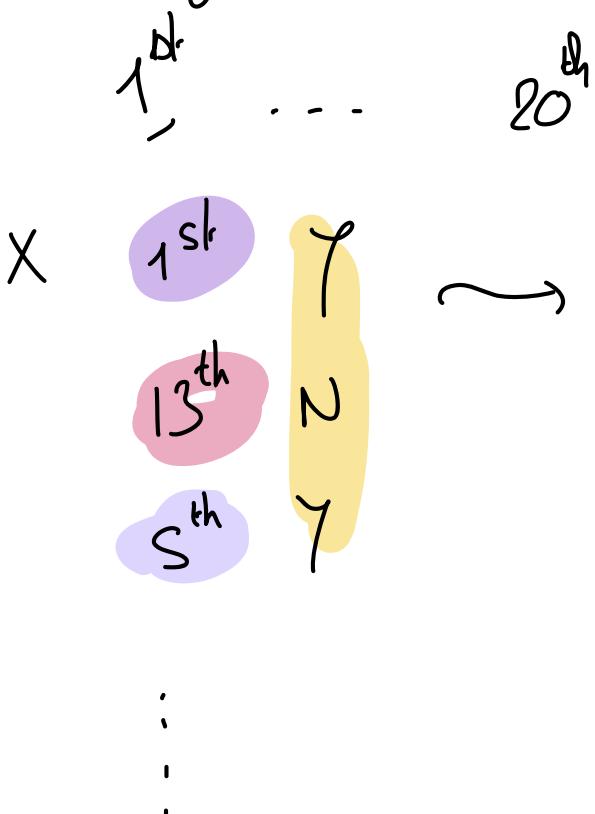
↓ ↓

type quantitative yes/no

~~Questions?~~

ONE HOT ENCODING

Categorical variable → encoding.

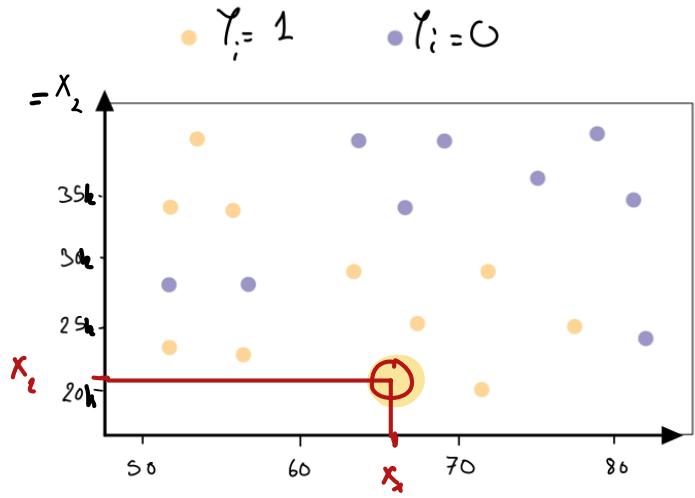


DUMM^Y ENCODING

Neighborhood	Surface	1 1 1 1 1 1	Garden
	1 0 0 0 0 0 0 0 0	1 0	1 0
	0 0 0 0 0 0 0 0 1 0	0 1	0 1
	0 0 0 0 1 0 0 0 0 0	↑ 13 th	1 0

Key idea : transform categorical variables with K categories into $K - 1$ { columns new features .

T Gy dataset

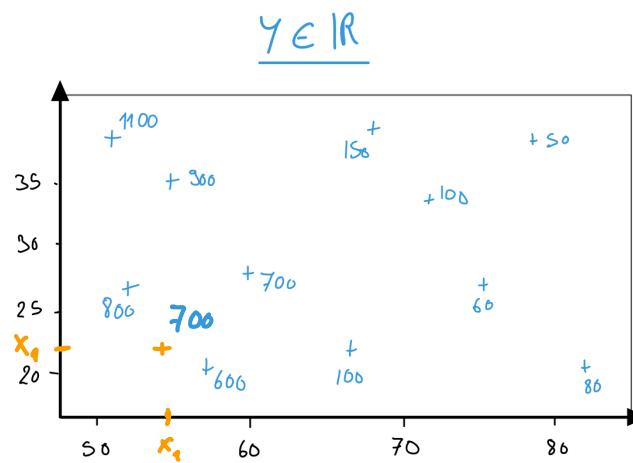


$$d = 2$$

$$n = 80$$

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

Classification



$$d = 2$$

$$n = 11$$

$$Y = \mathbb{R}$$

Regression



Summary of those three datasets

	Iris	MNIST	House price prediction	Toy
Number of examples n	150	50k	1500	80 - 11
Number of the features d	4	784	81 - 3	8
Type of the features	quantitative	quantitative	qualit + quant	quant.
Space \mathcal{Y}	{setosa, va, virg}	{0, ... 5}	\mathbb{R}_+	{0,1} / \mathbb{R}_+
Task	Classification		Regression	

Learning Pipeline

① Training phase

Training set: given a training set
on the train set

Goal: predict well! → most importantly: on the test set

Choose: method + its hyper-parameters

Learn: a predict^o function = model: \hat{P}

$$\hat{P}(X) \approx Y$$

② Validation phase

Choosing hyperparameters is

hard → try many, evaluate
the validate score for each.

- × Class of functions?
- × How to find a good one? } ≠ methods

Take the best hyperparameter → use those!

③ Test phase

(X_{test}, Y_{test}) both observed

Compute your test accuracy on never-seen data. → only number trusted

④ Deployment

X_{new} — Y is really unknown

Reminder on Linear and Logistic regression

Linear method

Melhoun

input point X .

learn weights w (Beta before)

biases b

$\langle X, \beta \rangle$

$\langle X, w \rangle + b$

new notes

Take away!

Lin. reg & Logistic regress

Task

Regres^o

Classif

Result

Linear combin^o of features

Learn

Weights

+ single part
- not enough

Limits of Linear and logistic regression

$$\text{House price} \approx \beta_0 + \beta_1 \times \text{Surface} + 100k \times \text{Is_garden} + \dots$$

$$\beta_k >$$

$$+ 80k \times \text{nb_rooms} + 40k \text{ open-kitchen}$$

+ ...

Outline

1 A Reminder on Supervised Learning

2 Decision Trees

3 Random Forests

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 ?

③ go beyond linear functions → learn complex dependencies .

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**

Inference Phase

$$Y = 1_{\text{price} > 1m}$$

- Root = observation = input

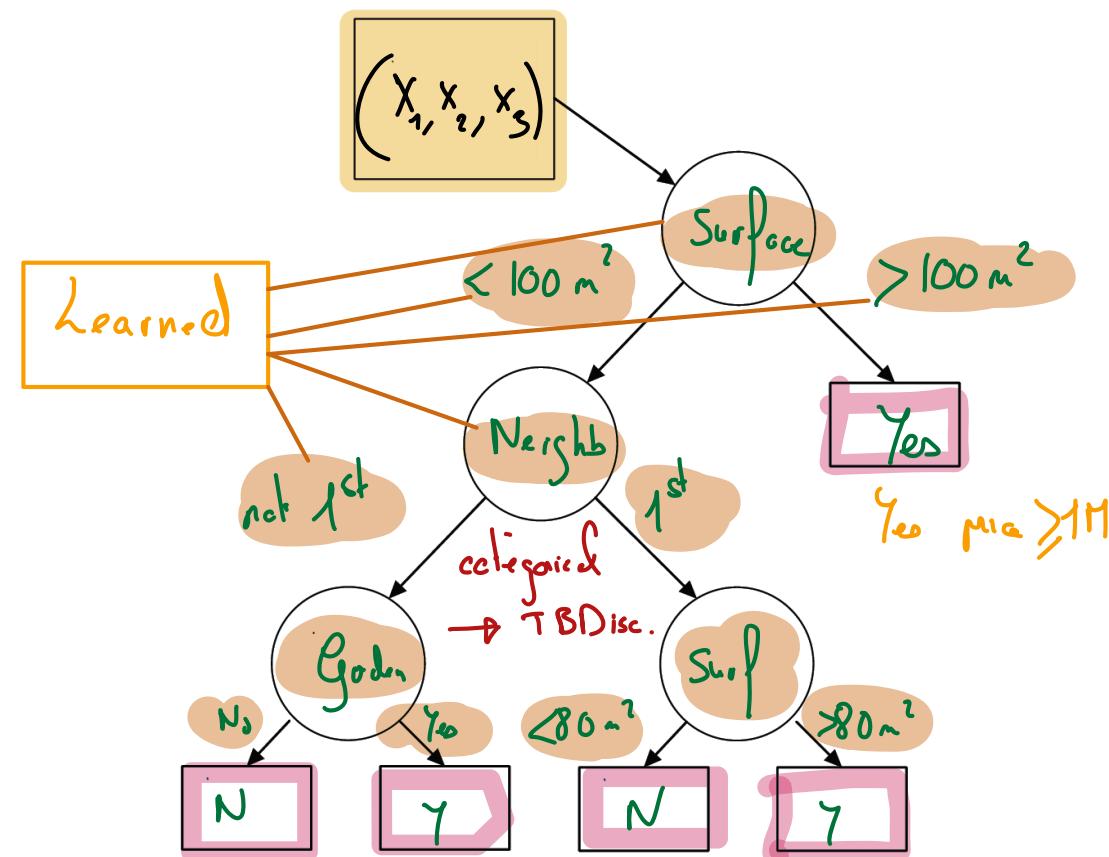
- Nodes = tests *binary answer*
- *only one feature per node*

- Branches = possible outcomes

- Leaves = final decision = output

house price. $Y = \mathbb{R}$

$X_1, X_2, X_3 = \text{Surface, Neighbor, Garden}$



For a given DT, predict inference | is | straight forward interpretable

Training: how to build a tree.

Key question !

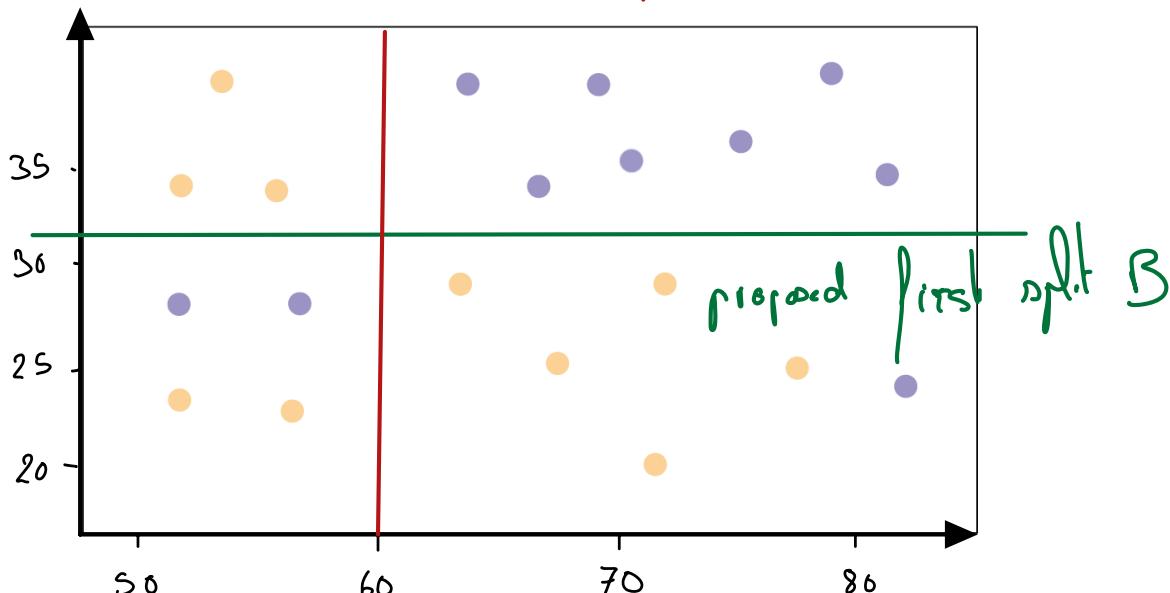
- What is a good tree?
- What do I need to choose?
- How to choose?

What do you think?

Building trees : CART



proposed first split A



Algorithm at each step :

- try all possible values for all possible features
- compute some heterogeneity of the resulting groups
- pick the split that minimizes heterog.

Heterogeneity?

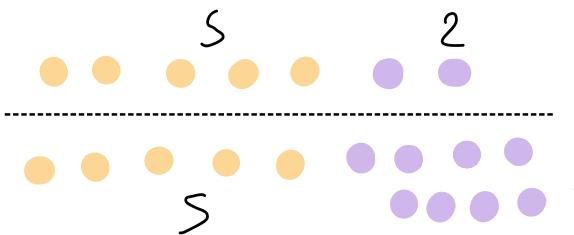
→ often several solns

Prop B



$$\frac{3}{10} \left(1 - \frac{3}{10}\right) = \frac{21}{100}$$

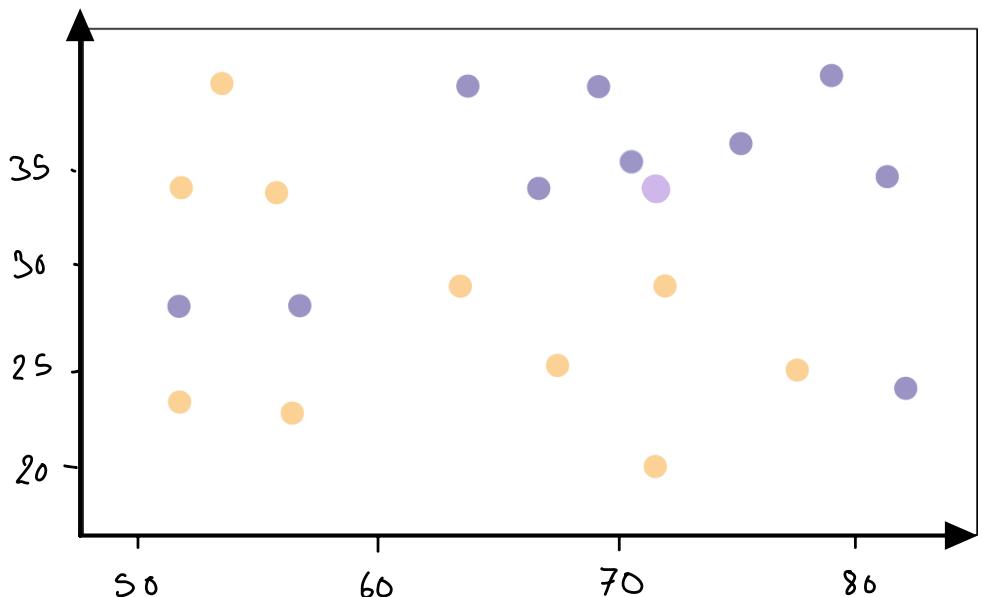
Prop A



$$\frac{2}{7} \times \left(1 - \frac{2}{7}\right)$$

$$\frac{S}{13} \left(1 - \frac{S}{13}\right) = p_i (1 - p_i)$$

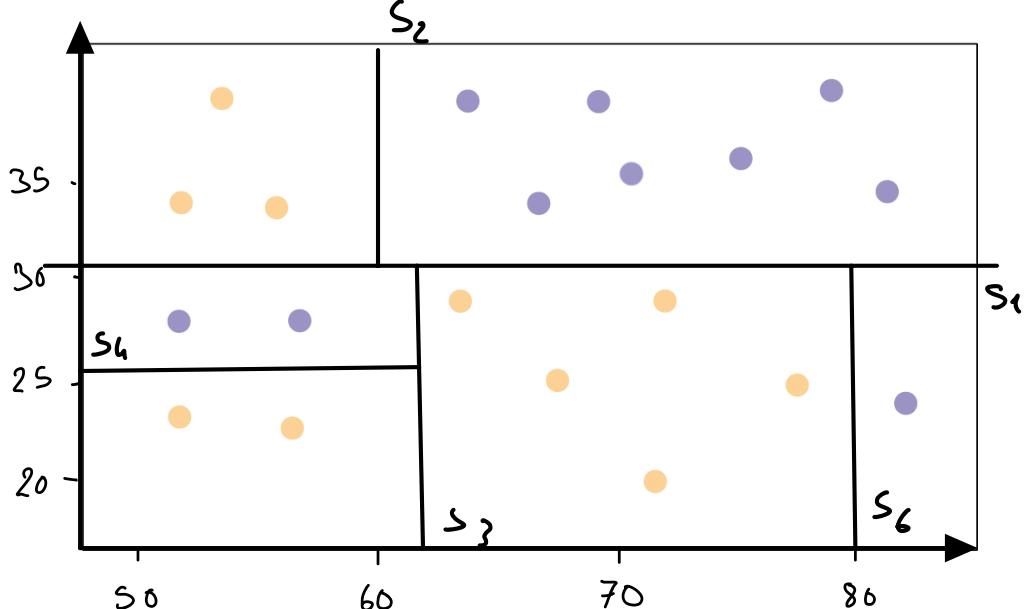
Building trees : CART



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.

Gini's impurity



$$\frac{3}{10} \times \frac{7}{10}$$
$$p_i \quad (1 - p_i)$$

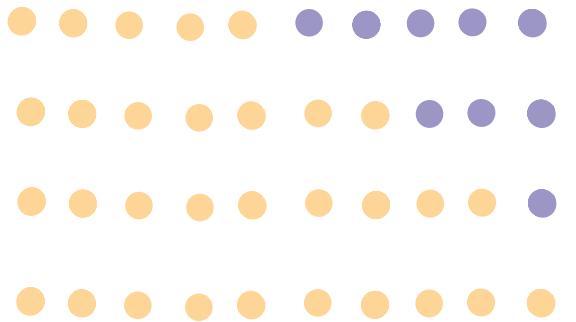
Gini impurity measures **how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset.**

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

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

Group

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



$$50\% \cdot (1 - 50\%) = \frac{1}{4} = 25\%$$

$$\frac{7}{10} \times \frac{3}{10} = 21\%$$

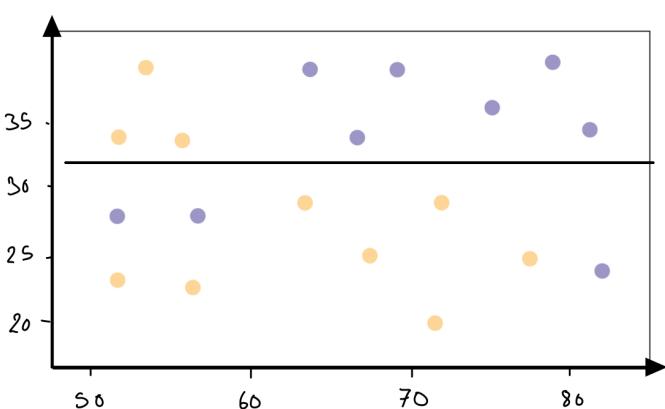
$$\frac{3}{10} \times \frac{1}{10} \approx 3\%$$

$$0\%$$

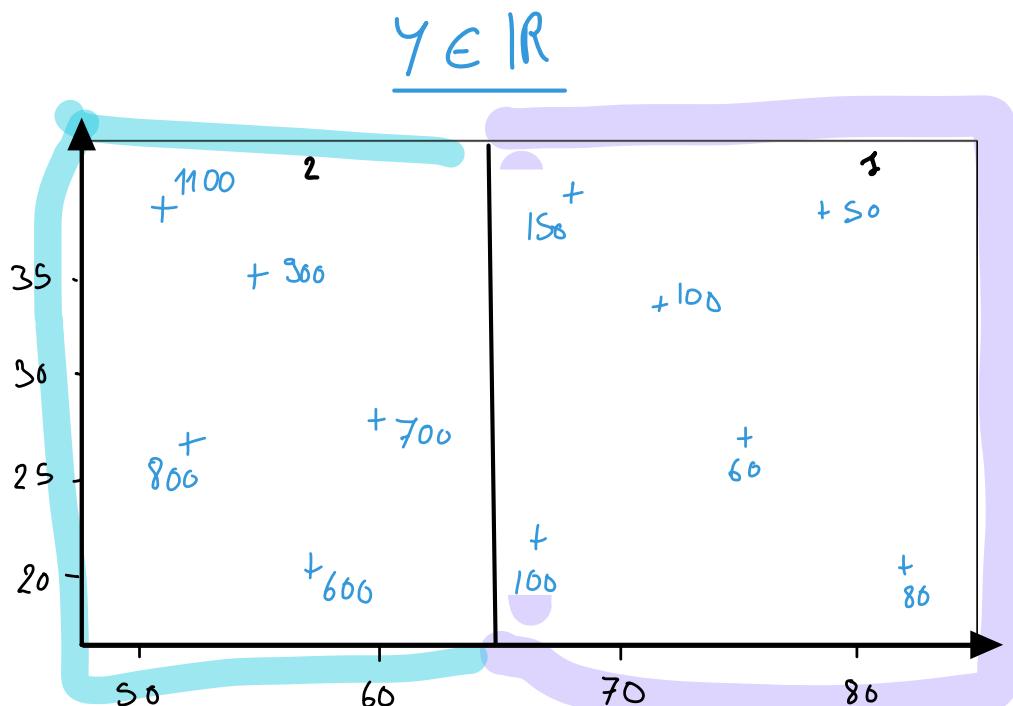
$$\times 2 (\textcolor{yellow}{\bullet} + \textcolor{purple}{\bullet})$$



$$\begin{aligned} & \frac{3}{10} \times \frac{7}{10} \\ & + \frac{3}{10} \times \frac{7}{10} \\ & + \frac{1}{10} \times \frac{6}{10} = \end{aligned}$$



Building trees : CART



$$\text{Variance} = \sum_{i=1}^{11} (\gamma_i - \bar{\gamma})^2$$

400

Variance after first split

$$= \sum_{i=1}^6 (\gamma_i - 100)^2 + \overline{\text{Var}_L}$$

$$+ \sum_{i=1}^5 (\gamma_i - 800)^2$$

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.

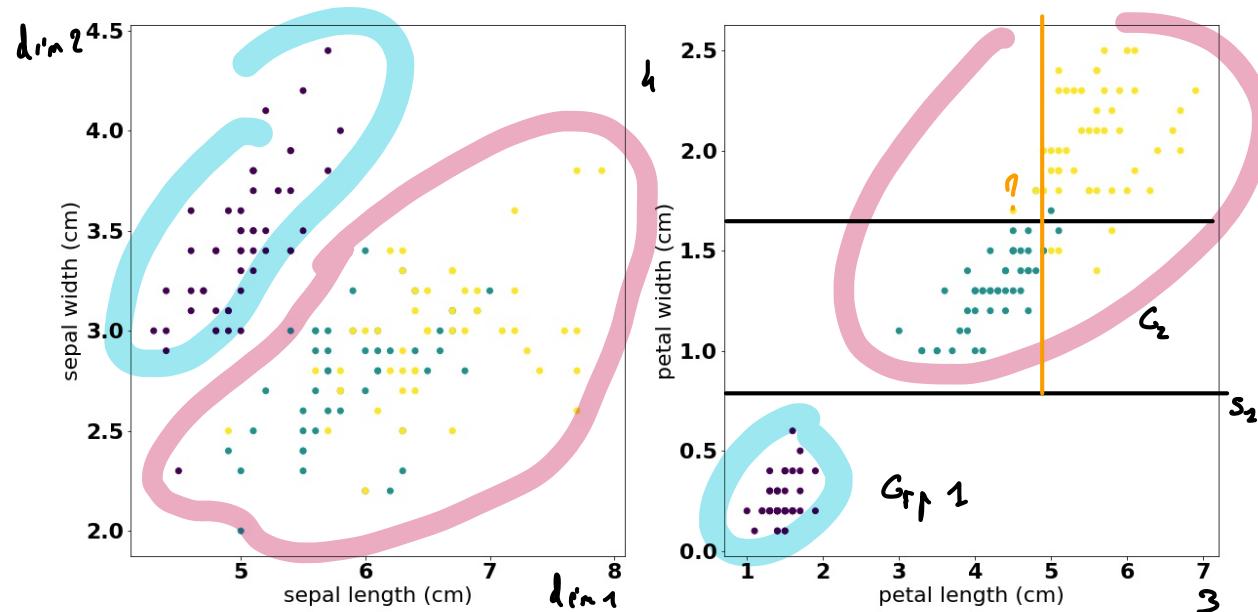
Example: Iris Dataset



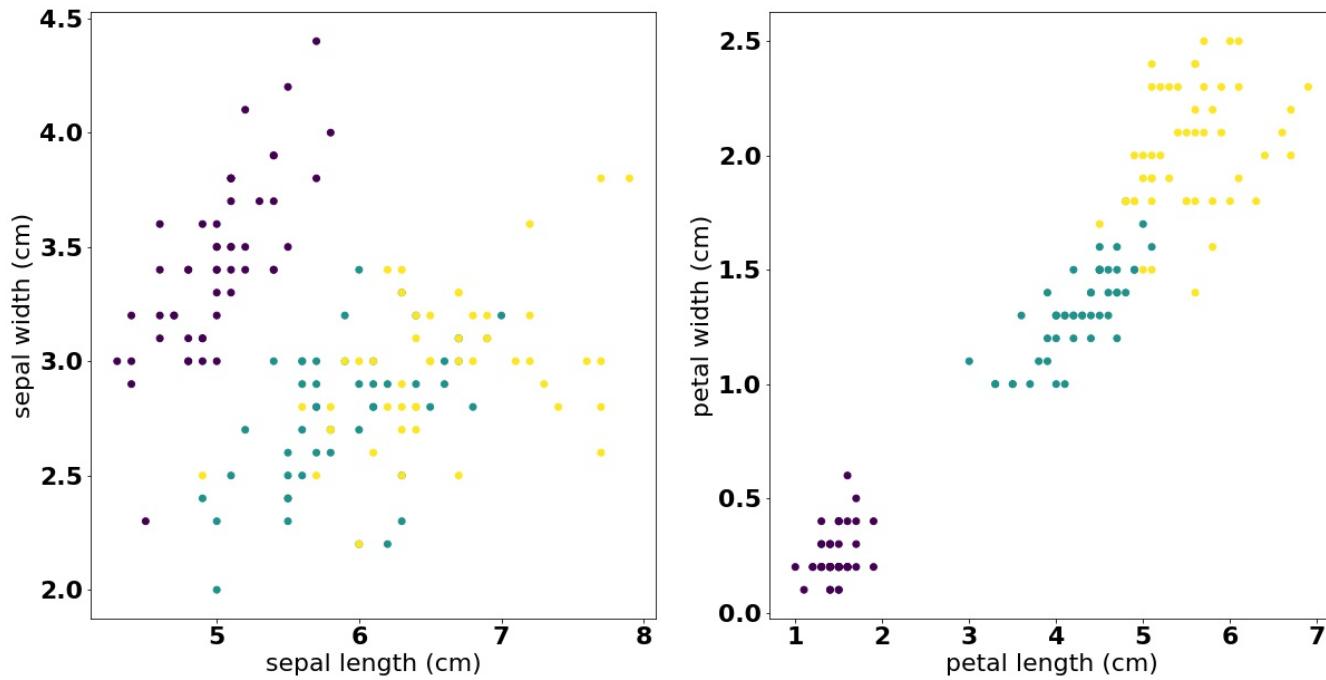
Figure: Iris: setosa, versicolor, virginica

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 ?

Hyperparameters: \rightarrow depth

→ avoid overfilling

Output:

$66\% \times 150 \text{ pts}$

100 hair pts

*before split gini = 0.662 }
100 train pts*

*nb of Setosa, Ver, Virg
35 28 37*

*spilt 1.2
impurity of {0, 28, 37}*

True

False

pure class

petal width (cm) <= 0.8
gini = 0.662
samples = 100
value = [35, 28, 37]
class = virginica

petal width (cm) <= 1.65
gini = 0.49
samples = 65
value = [0, 28, 37]
class = virginica

petal width (cm) <= 0.8
gini = 0.0
samples = 35
value = [35, 0, 0]
class = setosa

petal width (cm) <= 1.65
gini = 0.18
samples = 30
value = [0, 27, 3]
class = versicolor

petal width (cm) <= 1.65
gini = 0.056
samples = 35
value = [0, 1, 34]
class = virginica

Example:

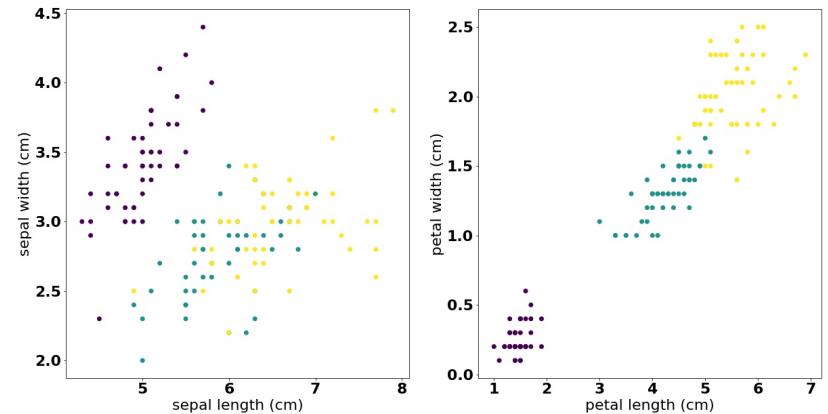
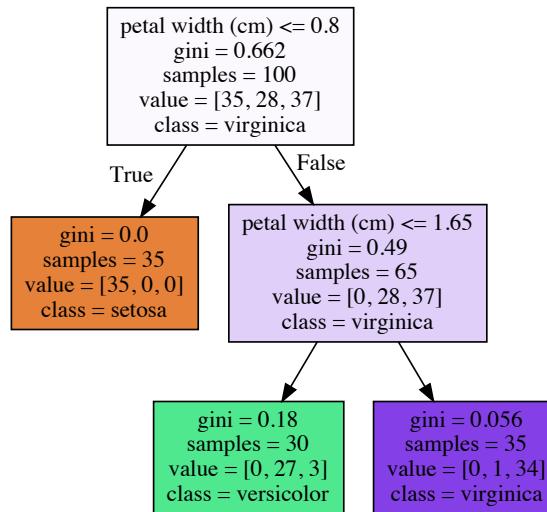
```
# Check  
X_test[0,:], iris.target_names[y_test[0]]
```

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

Code: <https://colab.research.google.com/drive/1dXYIjacNAeASgl54vIWjMQlfgigSVA?usp=sharing#scrollTo=VeJHkJd7RsDw>

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>

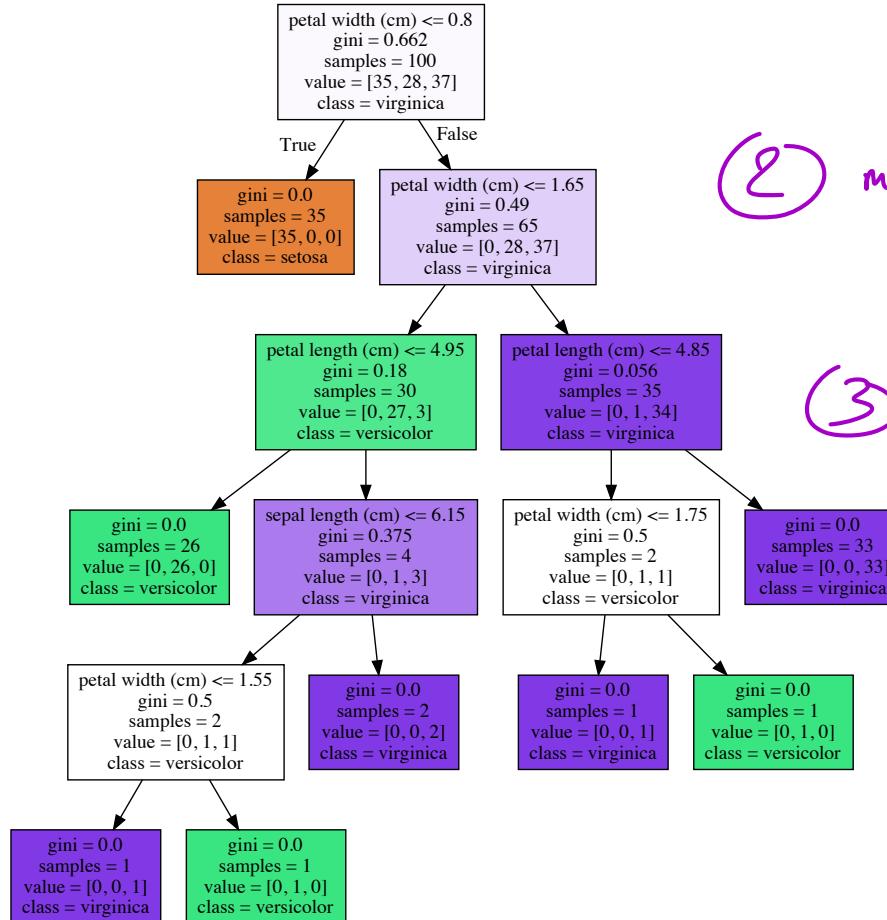
What if I use a deeper tree ? Or if I do not specify the depth in the model definition ?

Deeper trees

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

① max_depth = 2

For any node, at most 2 questions



② min - samples - leaf

③ min - impurity - decrease

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

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.

} "taille" prune the complete tree.
→ other soln.

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

Pros:

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

I could do by hand given enough time

\longrightarrow naive implement. $O(nd)$ complexity per split \rightarrow easy to reduce

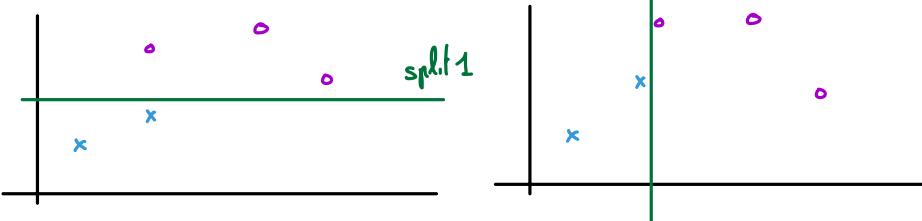
Cons:

- 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

Solution : Using Random Forests !

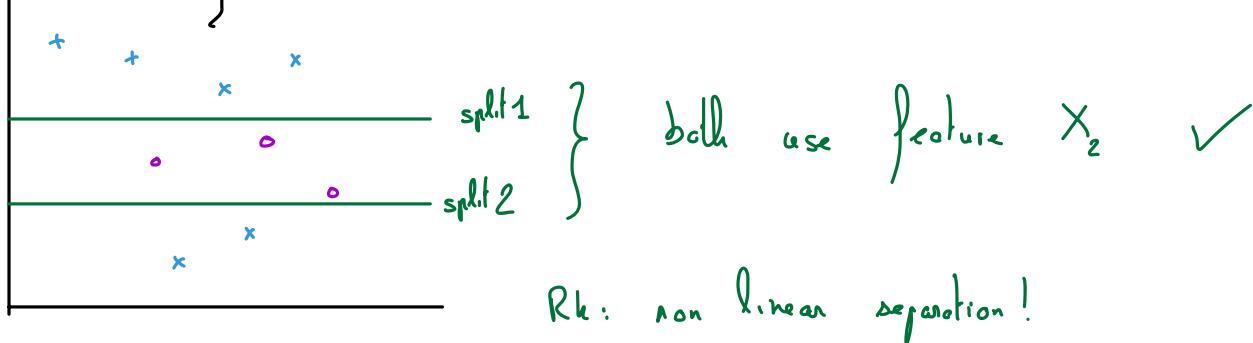
Woodchop.

Not stable

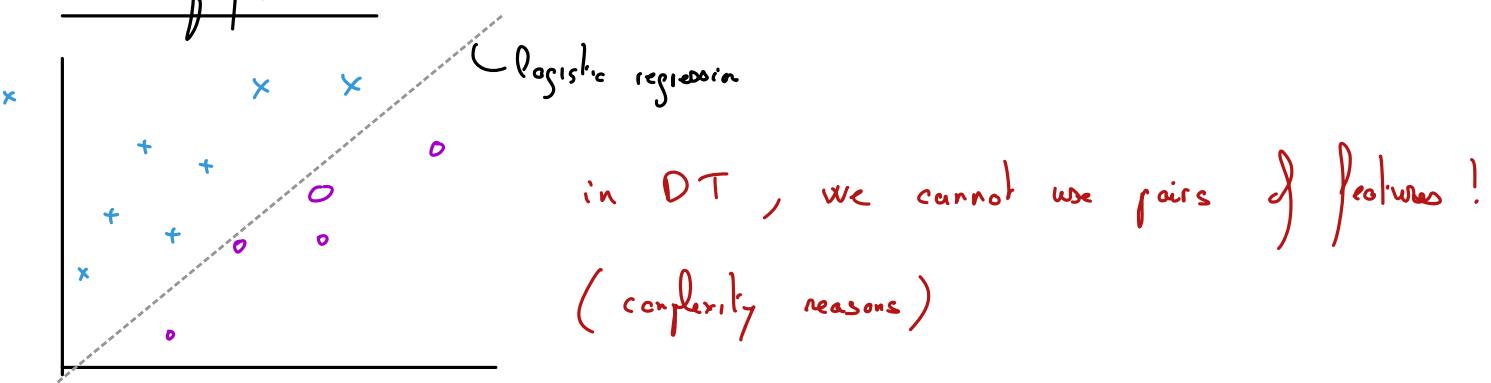


changing one data point
may drastically change
the decision tree.

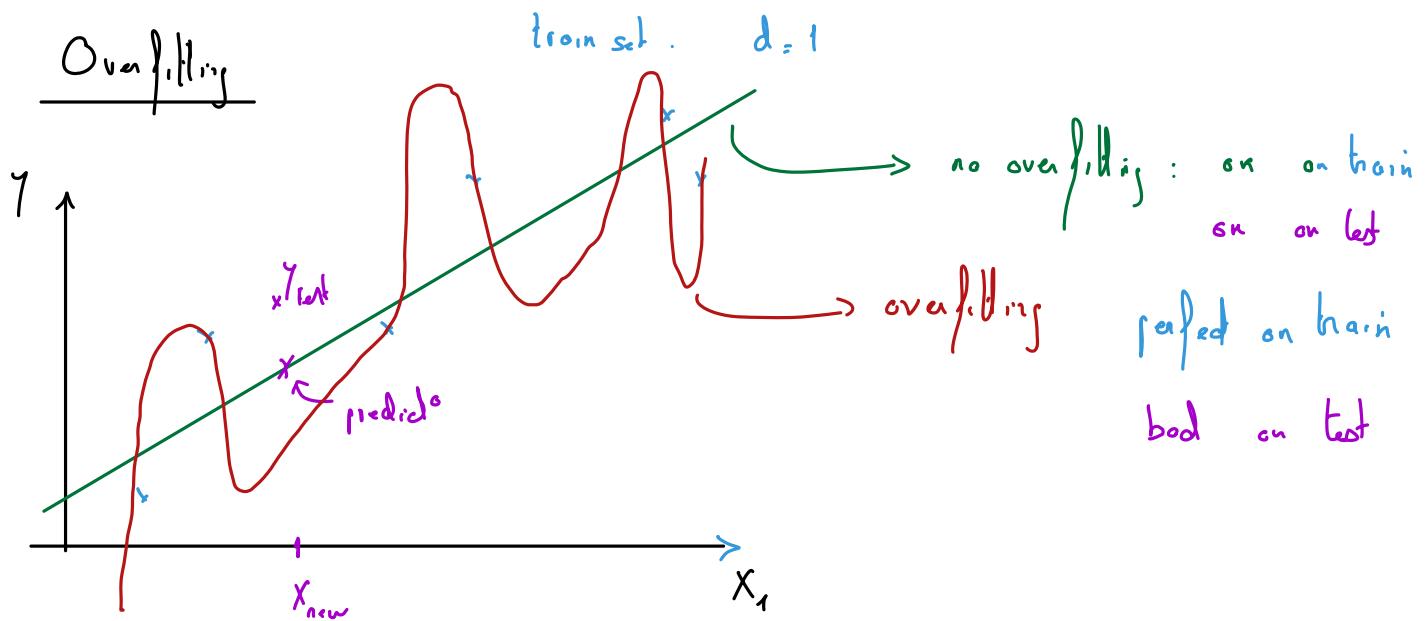
Some Feature Twice

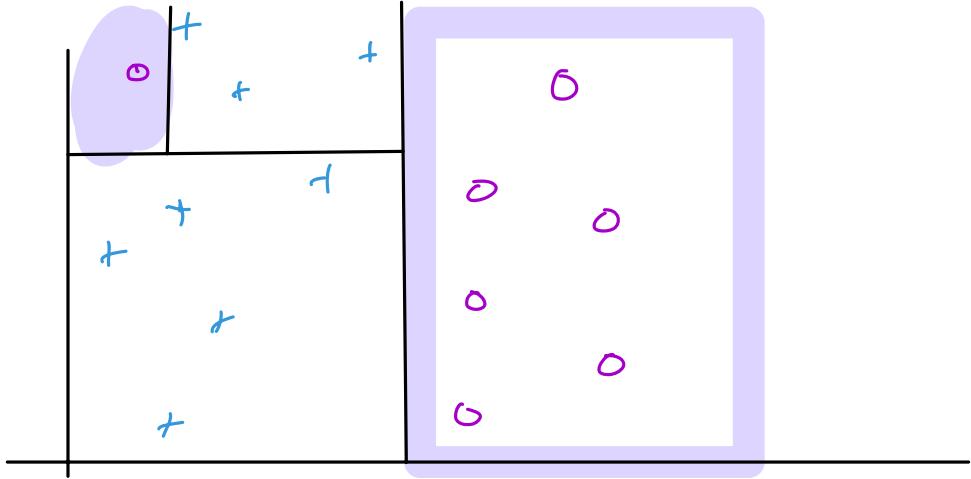


Pairs of Features



Overfitting





Train / Validation / Test

Avoid overfitting

Outline

1 A Reminder on Supervised Learning

2 Decision Trees

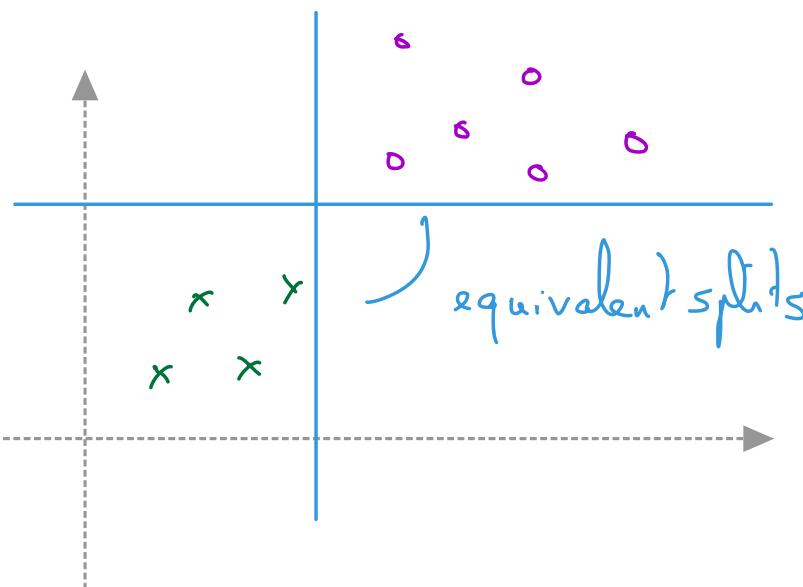
3 Random Forests

Random Forests

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

How to create variability ?

① At each step test only some features



② Bootstrapping :

Tree 1: split 1,
only uses X_1

Tree n°2 split n°1
only uses X_2 .

each tree use a
part of data to learn

Random Forests

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

How to create variability ?

- ① • Instead of the most important features, search the best feature among a random subset (~~ntry~~). ~~max_feature~~
- ② • Work on subsets of data (bootstrap=True).
- More <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>

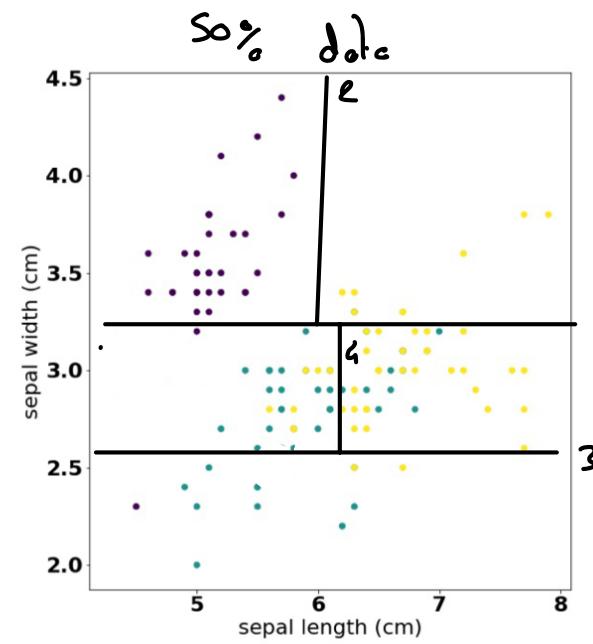
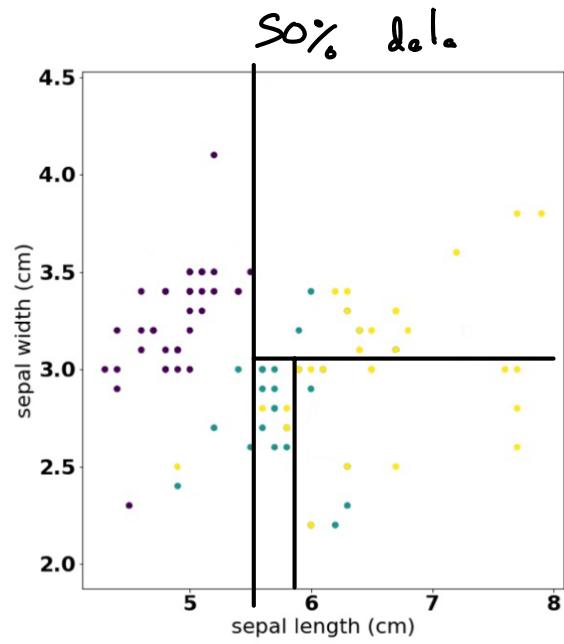
```
# Creating and fitting a Random Forest
from sklearn.ensemble import RandomForestClassifier

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

clf.fit(X_train,y_train)
```

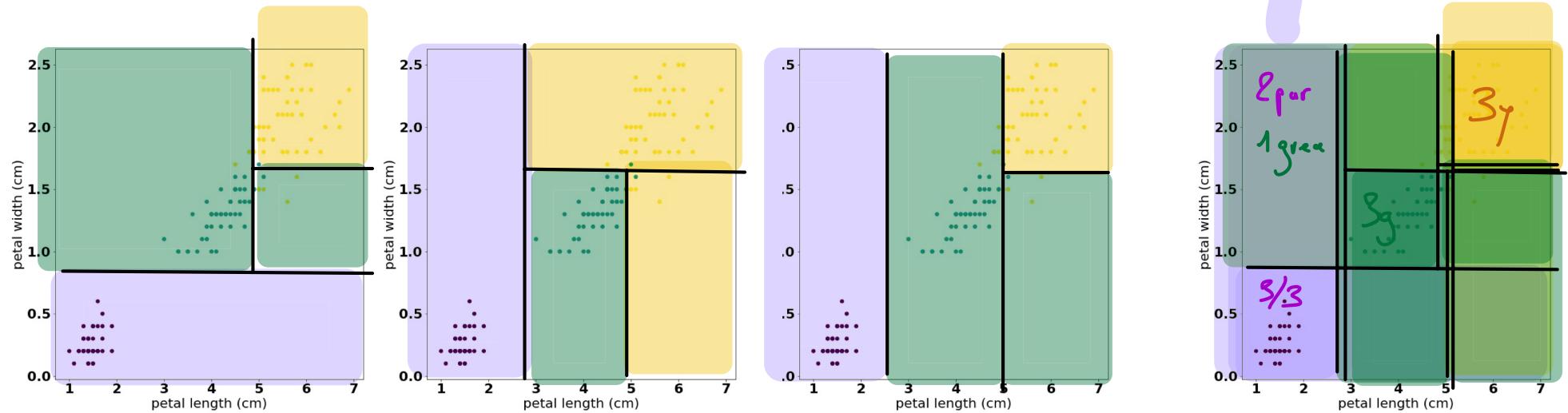
How to aggregate ? How to interpret a Random Forest ?

Bootstrap



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.

good enough

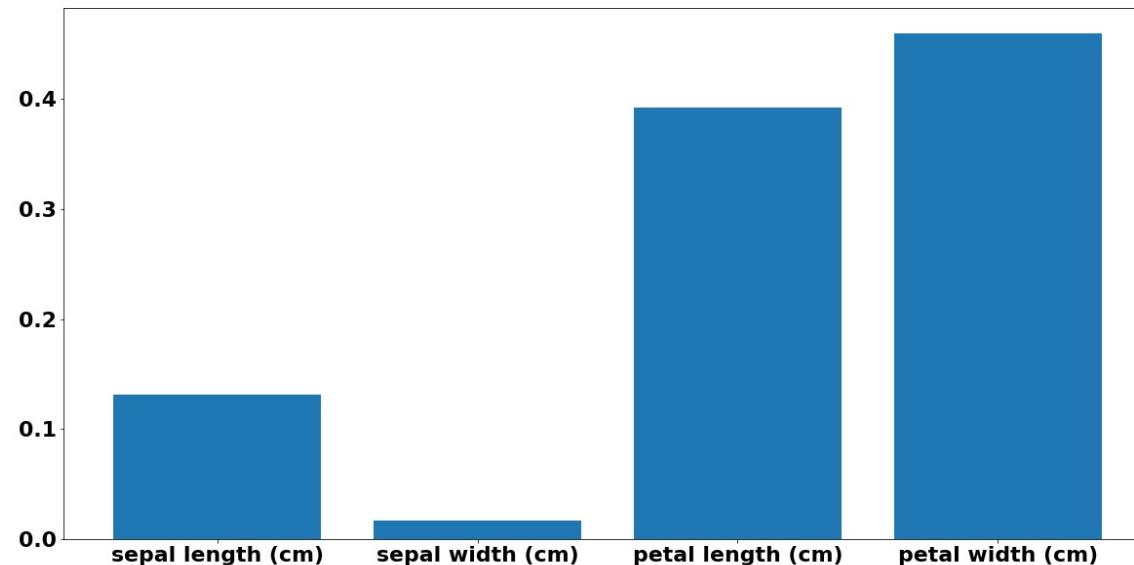


Figure: Feature importance output for a Random Forest Classifier in Python

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.

red a
risk overfitting
will prevent
in place

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 ? *All the same!*

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

Ⓐ we can't distribute computation as much (dependency / order)

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

(between trees)

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

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

Lab :

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

