

Decision Trees and Random Forests

Data science Certificate

Aymeric DIEULEVEUT

May 2020

- goal in ML
- how it is defined
- parameters
- implement it in Python.
- + / -

1 Decision Trees



2 Random Forests

Outline

1 Decision Trees

2 Random Forests

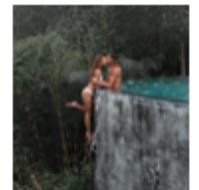
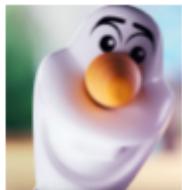
Random Forests

Supervised

ML

framework

- Task: do these pictures contain individuals that look nice (yes or no), suspicious (yes or no)?
- Train algorithms using data (a lot of these):



Nice: yes
Susp: no

Nice: no
Susp: yes

Nice: yes
Susp: yes

Nice: yes
Susp: yes

Nice: yes
Susp: yes

- Pictures: $X_i \in \mathbb{R}^d$, Labels $Y_i \in \{0, 1, 2, 3\}$
- Features: categorical and numerical features (is there a smile? what is the age of the person?)
- Size of the dataset: n for training + t for test.

each image $\in 1$ category.
 $\{0, 1\}$

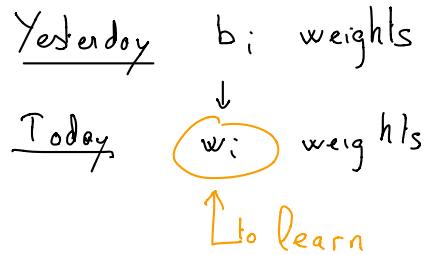
Goals of Decision Trees

Features :

	age	weight	height	salary
25	65	1,8	1M	

Learned : 2 -1 20 0.005

Prix of house : not easily interpretable



- ① Prediction: what is the class of this picture (Classification) ? What is the age of the person (Regression) ?



- ② Provide some understanding: what on this picture allows me to classify it? Ex: presence of a face, human, smiling, etc ...

image: has features $x_1^i \dots x_d^i$

d nbr of features

linear regression: predict output y_i by $\sum_{i=1}^d w_i x_i^i$

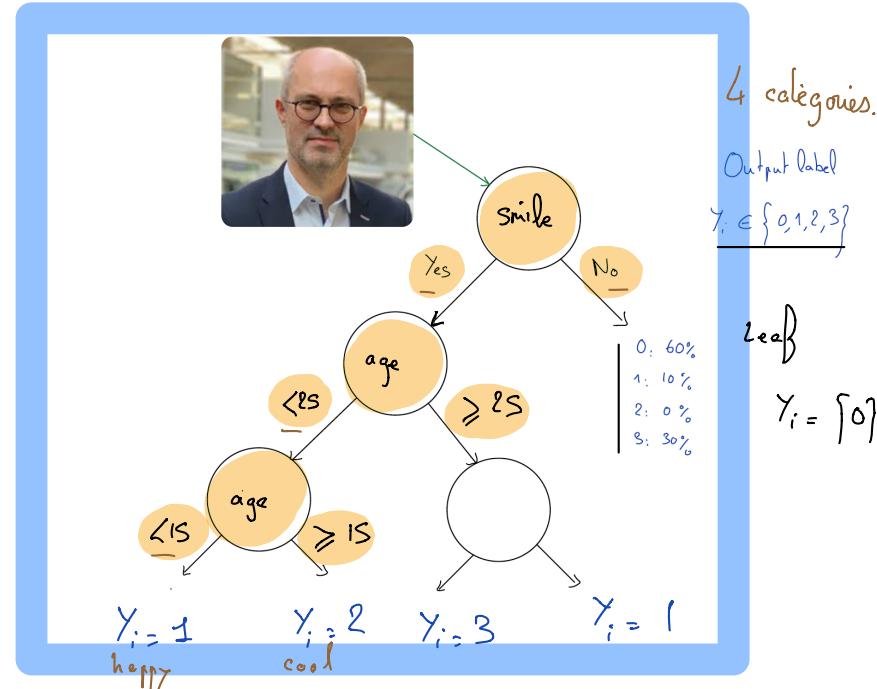
Decision Tree

Training data : images and their features
④ category

- Nodes = tests
- Branches = possible outcomes
- Leaves = final decision.

Learn : For each node :

- which feature to use
- decision criterion.



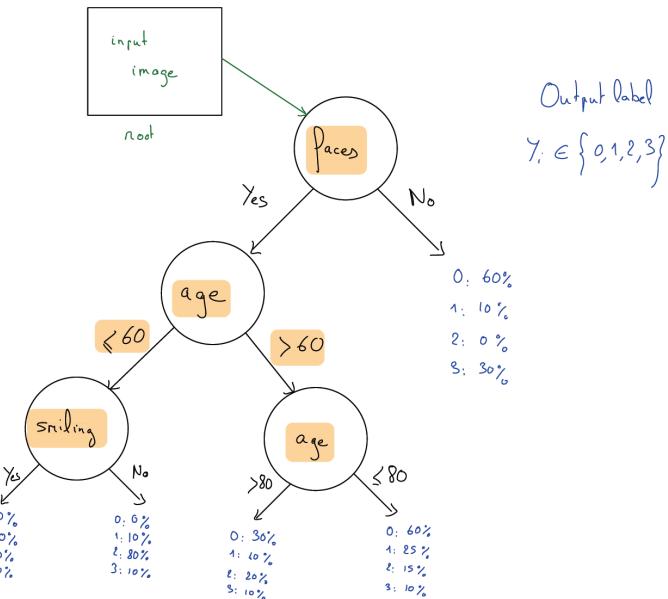
Output:
* classification : category
* regression : continuous
what you want to predict.

Decision Tree

- Nodes = tests

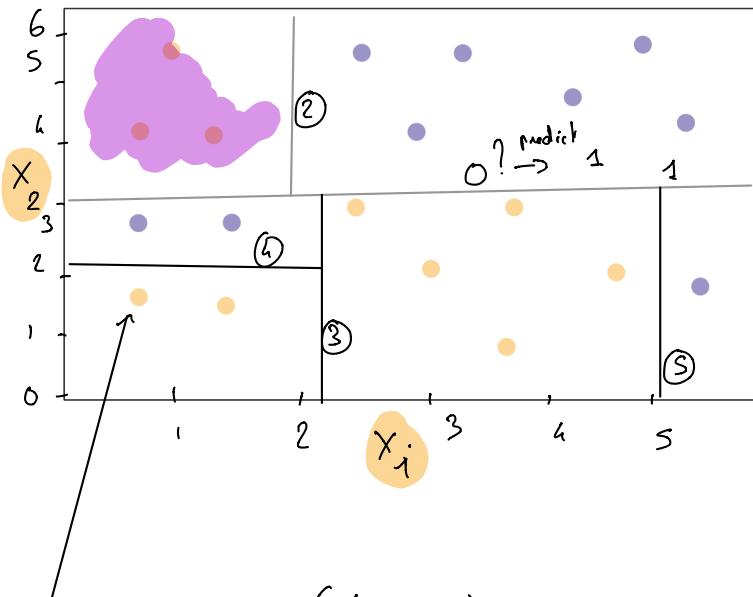
- Branches = possible outcomes

- Leaves = final decision.

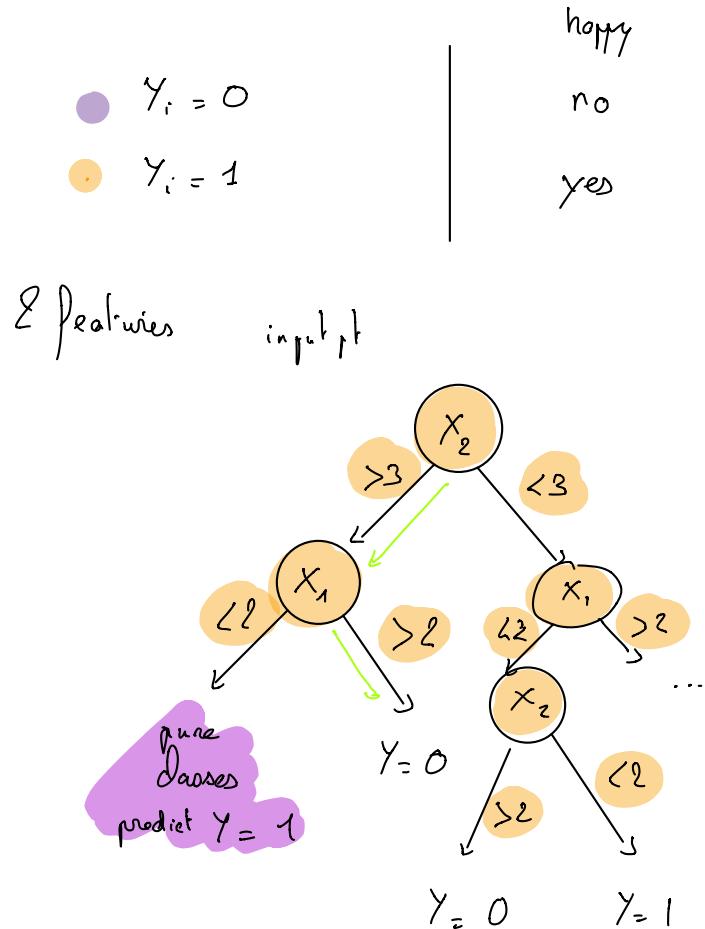


Building trees : CART

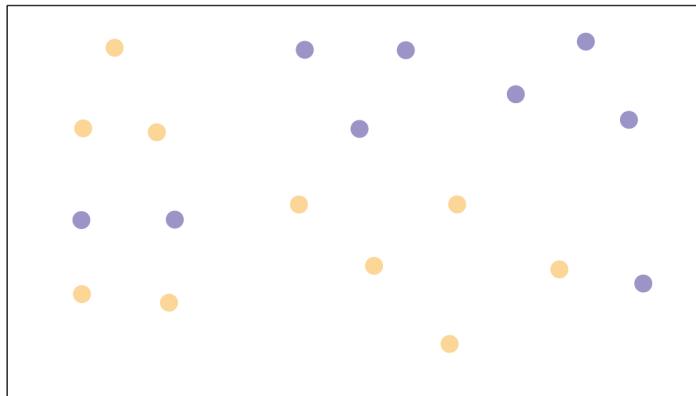
1 algo to learn tree



obs n° 1: (X_1^1, X_2^1) , label "no" = 0



Building trees : CART



$$\begin{array}{c} \text{Y} = 0.5 \\ + 0 \\ - 1 \\ \hline \end{array} \quad \left. \begin{array}{c} \text{Y} = 15.5 \\ + 10 \\ - 21 \\ \hline \end{array} \right\} \text{average} \quad \frac{32}{4} = 8$$

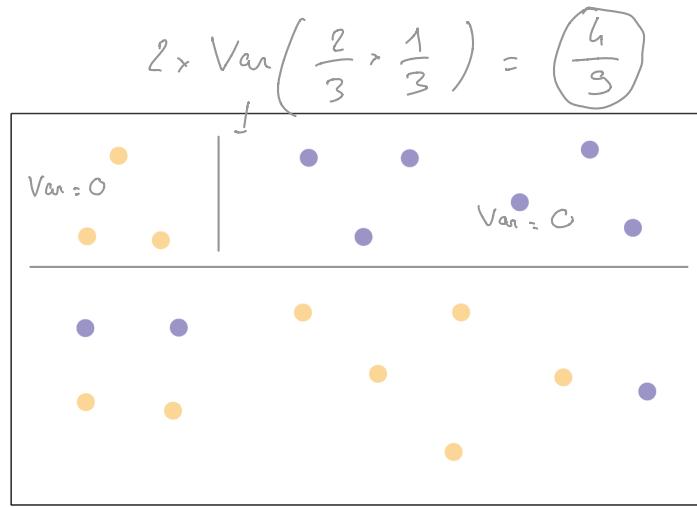
no split : $MSE = (0.8)^2 + (1-8)^2 + (10-8)^2 + (21-8)^2$

after split $MSE = (0-0.5)^2 + (1-0.5)^2 + (10-15.5)^2 + (21-15.5)^2$

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.

Gini's impurity



9 blue pts
10 orange pts

Gini's impurity ≈ 0.5

for a class in ratio, probability of misclassif = $(1 - p_i)$

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.**

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

probabil i probabil i
misclassif misclassif.

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

$$\text{Var}(\mathcal{B}(p)) = p(1-p)$$

$$I_G(p) = 2 \text{Var}(\mathcal{B}(p))$$

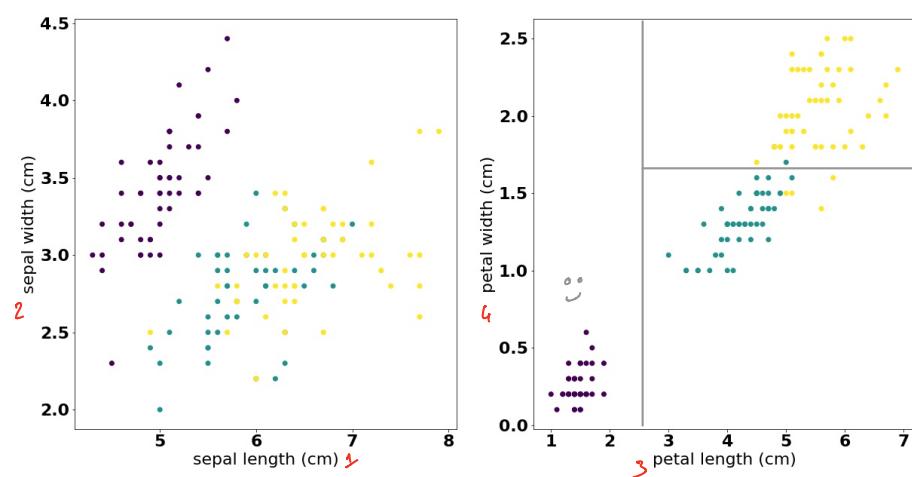
Example: Iris Dataset



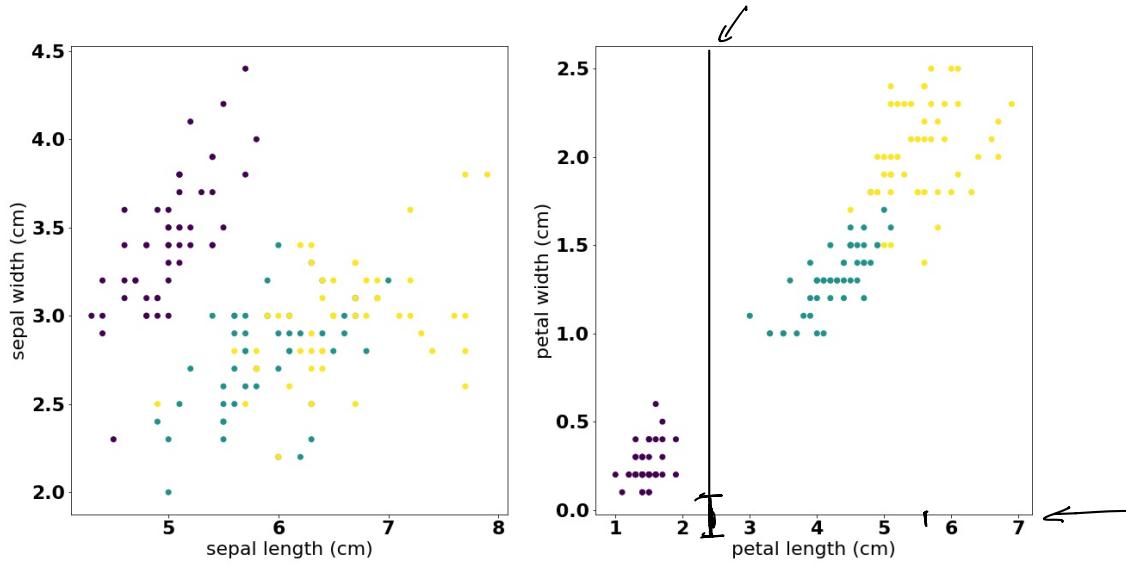
Figure: Iris: setosa, versicolor, virginica

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



cut is always a straight line

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 ?

Step 1

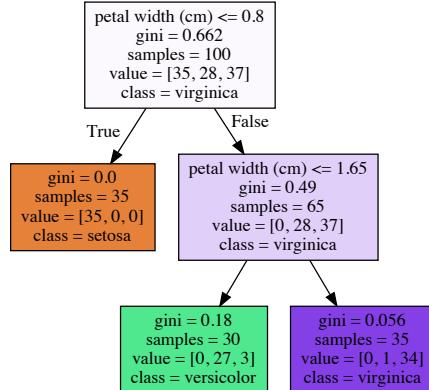
Step 1 : specify

- * impurity criterion : Gini is default
- * stopping criterion: tree_depth = 2

Fit model: X_{train} , Y_{train}

Step 2: $\text{clf} = \text{DecisionTreeClassifier}(\dots)$ \rightarrow learns the tree

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]]
```

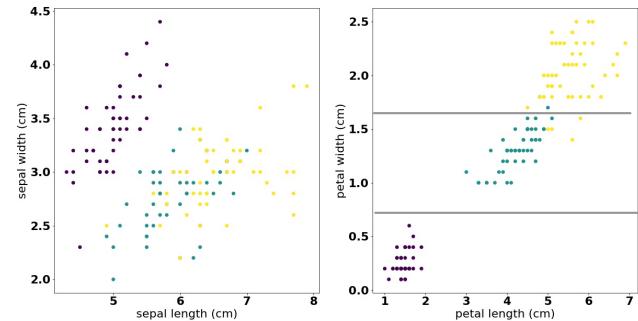
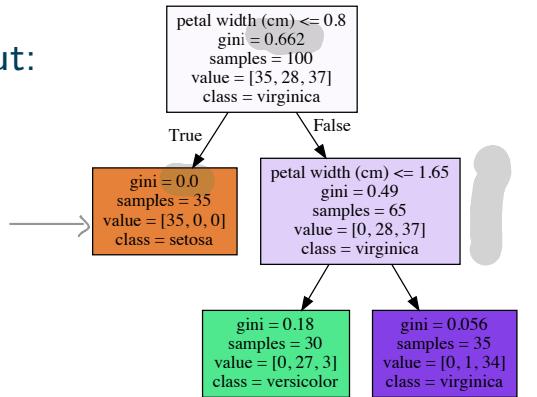
Returns:

```
(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

bad point

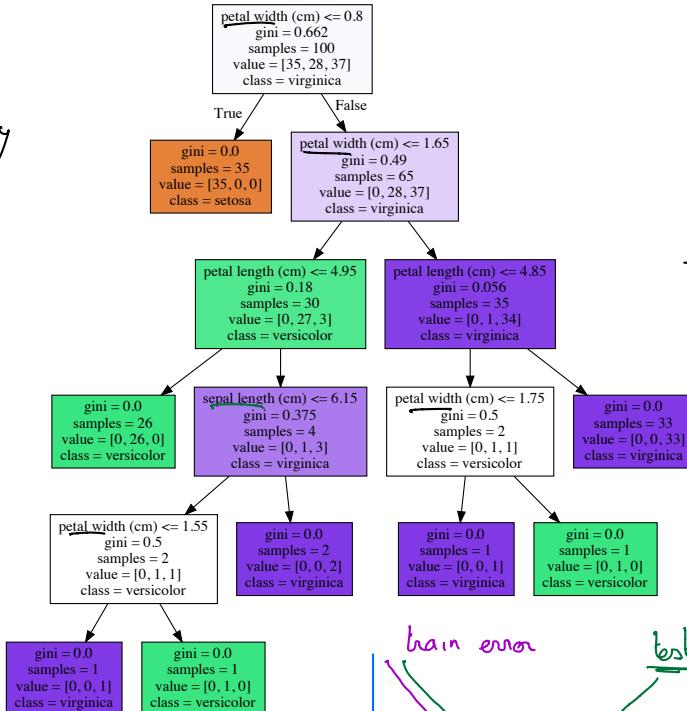
```
# Choose the Learning algorithm  
clf = DecisionTreeClassifier(random_state=0)
```

{ defining model : no max-depth

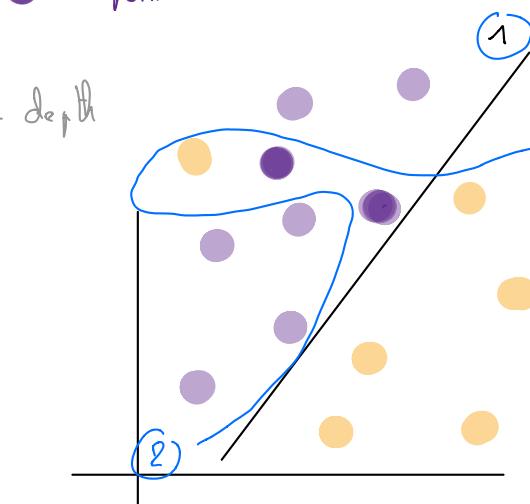
Overfitting.

* very good accuracy
on training

* degraded acc on
test



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



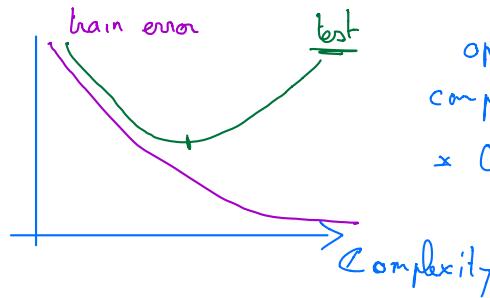
option 1 is simple

* but 1 error train
set

* less errors on test
set

option 2 is more
complex

* 0 error train set



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.



Learn in practice

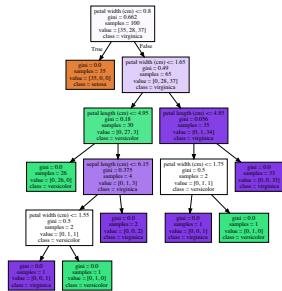
② Pruning

- ▶ 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

11h 15.

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.

Cut is always using
only one feature

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 iterative construction
- ④ Biased if some classes dominate

Solution : Using Random Forests !



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 ?



- Instead of the most important features, search the best feature among a random subset (ntry). X_{train}, Y_{train} : to build tree # i , you use a random subset $X_{train}^i; Y_{train}^i$
- 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)
```

← 100 trees in ↑
← all trees have depth 3
← use bootstrap.

1st step: defining the model.
all hyperparam. / options

2nd step: fitting model.

How to aggregate ? How to interpret a Random Forest ?

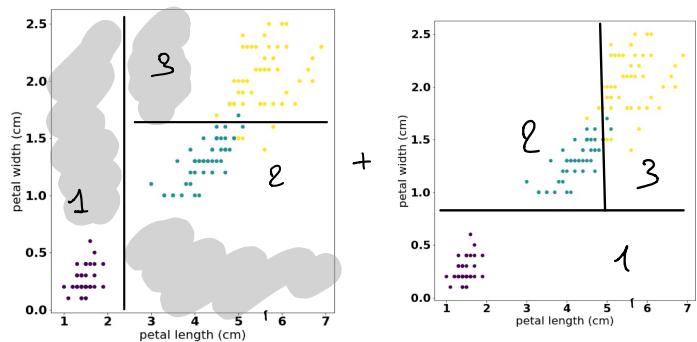
Aggregation

first, t 1 2 3 -
sg 1 0 ← h.g.h conf. dence

2nd 1 SO 40 10 ← less certain

$$\text{Classified} \quad \left\{ 3, 2, 8, 2, 5 \right\} = \frac{\sum_{j \in \text{tree output}}}{n_{\text{tree}}} \quad \begin{matrix} \uparrow & \uparrow \\ \text{Ensembling} & \end{matrix}$$

We using a voting or averaging process !



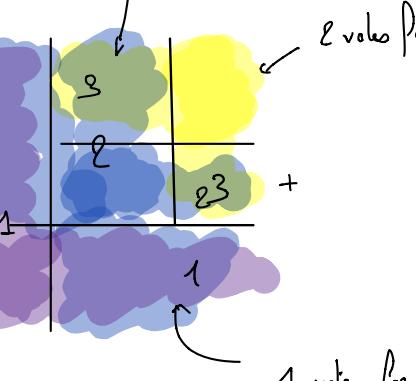
for any μ^t : check how many times over 100 has
it is placed in each category?

60 times cat 1

20 times 2

20 times cat 3

\Rightarrow final classif 1



Feature importance

petal width
 γ

- 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.

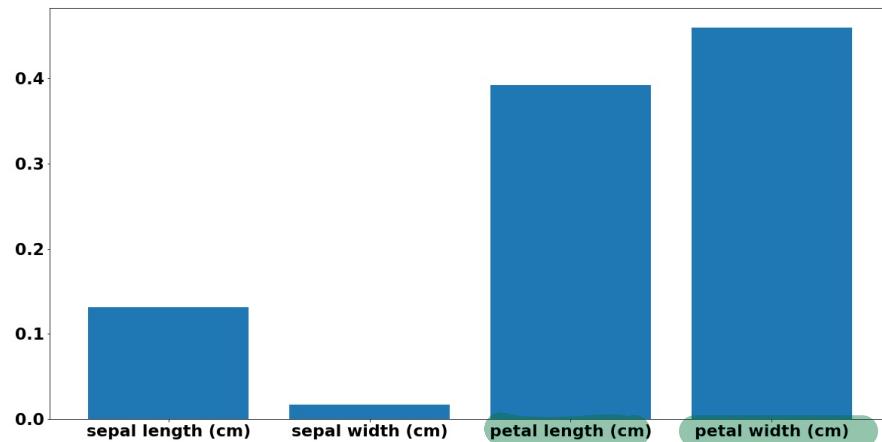


Figure: Feature importance output for a Random Forest Classifier in Python

Comparison with decision trees

RF:

already existed
for DT

Cons:

- Less interpretable ✓
- Slower to run ✓

arguments / hyperparameters

* specific to the forest:
* number of trees
* randomness split

* related to each tree:
* max_depth
* min_sample_leaf

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.

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 ?

→ use MSE instead of Gini for each tree
→ averaging instead of voting to aggregate the trees.

Questions?

Boosting

- ⊕ gain in accuracy / precis°
- ⊖ no distributo possibl of the workload.

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>

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.

create a tree visualize
create a forest

* create a forest
* check the most important feature
* use only a few features.

Melba Lab :

- * PyCharm.
- * Colab : * need to connect to gggl account
* save the nb at beg. & end.

*