

Arbres et Forêts Aléatoires

Master parcours SSD - UE Apprentissage Statistique II

Pierre Mahé - bioMérieux & Université de Grenoble-Alpes

Arbres de classification et de régression

- ▶ méthode ancienne, classique et populaire

A la base de la méthode des forêts aléatoires

- ▶ basée sur le **bagging** : aggrégation par bootstrap
- ▶ très performante sur de nombreux problèmes
- ▶ facile à paramétrer

Ce cours : des arbres aux forêts aléatoires

- ▶ **arbres** : motivation, formalisation, algorithmes et limites
- ▶ aggrégation d'arbres et **bagging**
- ▶ du bagging aux **forêts aléatoires**

Introduction

CART

Bagging

Random Forests

Conclusion

Python

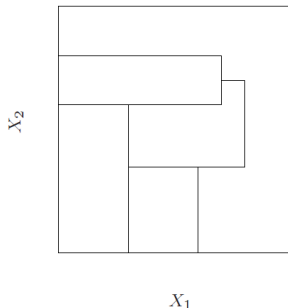
Références

Introduction

Principe

Principe¹ :

1. découper l'espace d'entrée en régions
2. estimer & prédire une valeur par région



⇒ **Régression** : valeur moyenne dans la région

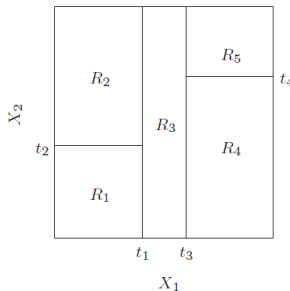
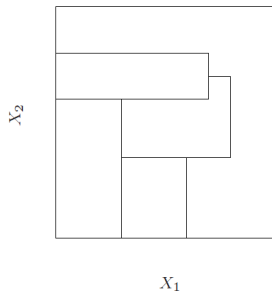
⇒ **Classification** : classe majoritaire dans la région

Principe

1ère question : comment construire ces régions ?

Approche par **partitionnement binaire récursif**

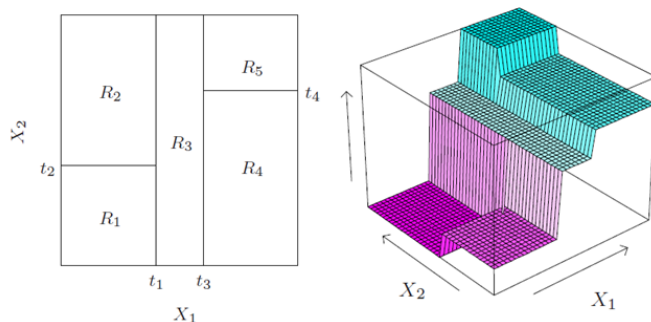
- ▶ binary recursive splitting



⇒ Chaque étape : 1 région divisée selon **1 seuil sur 1 variable**

- ▶ moins général mais plus facile à construire

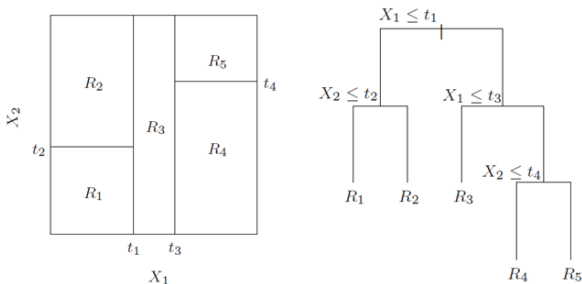
Illustration : régression



⇒ fonction **non linéaire** et **non continue**

Des arbres ?

Le processus récursif peut se traduire comme un **arbre** :



- ▶ **racine** de l'arbre : tout l'espace
- ▶ **feuille** de l'arbre : 1 région
 - ▶ régions imbriquées les unes dans les autres
- ▶ **noeud** (interne) : 1 seuil qui coupe une région en deux
 - ▶ test logique basé sur la valeur d'une **variable**

Avantages :

- ▶ interprétabilité du modèle
- ▶ mécanisme de prédiction proche du processus humain
 - ▶ e.g., domaine médical
- ▶ interprétation en tant qu'arbre valables pour $d \gg 2$
 - ▶ fonction de prédiction très difficile à se représenter

Avantages :

- ▶ interprétabilité du modèle
- ▶ mécanisme de prédiction proche du processus humain
 - ▶ e.g., domaine médical
- ▶ interprétation en tant qu'arbre valables pour $d \gg 2$
 - ▶ fonction de prédiction très difficile à se représenter

Questions :

- ▶ quel(s) critère(s) pour couper une région en deux ?
- ▶ quand arrêter le découpage ?

⇒ algorithme **CART**

CART : Classification And Regression Trees

Algorithme récursif

Initialisation : noeud racine

- ▶ contient l'ensemble du jeu de données / tout l'espace

On introduit un noeud :

- ▶ en séparant en deux groupes les observations affectées à un noeud de l'arbre
- ▶ selon un critère objectif
- ▶ tant qu'un critère d'arrêt n'est pas atteint

Construction de l'arbre

Comment **séparer en 2** les observations affectées à un noeud ?

Outline

Apprentissage Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

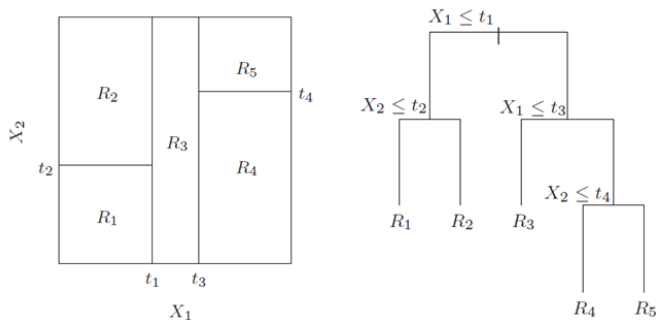
Python

Références

Construction de l'arbre

Comment **séparer en 2** les observations affectées à un noeud ?

⇒ en considérant **1 seuil** défini sur **1 variable**



⇒ principe décliné pour la **classification** et la **régression**

Arbres de classification - formalisation

Comment **séparer au mieux** les instances affectées au noeud ?

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Arbres de classification - formalisation

Comment **séparer au mieux** les instances affectées au noeud ?
⇒ en cherchant **deux ensembles les plus purs possibles**.

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Arbres de classification - formalisation

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Comment **séparer au mieux** les instances affectées au noeud ?

⇒ en cherchant **deux ensembles les plus purs possibles**.

Formalisation (pour la classification) :

- ▶ problème de classification à K catégories
- ▶ le noeud i contient n_i instances
- ▶ on définit le vecteur $P_i \in \mathbb{R}^K$:

$$P_i[k] = \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbf{1}(y_{C_i(j)} = k)$$

où C_i est l'ensemble des instances du noeud i

⇒ **les proportions des différentes classes** au sein du noeud

Arbres de classification - critères d'impureté

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Critères d'impureté du noeud i :

1. indice de Gini :

$$G(i) = 1 - \sum_{k=1}^K P_i[k]^2$$

2. entropie :

$$H(i) = - \sum_{k=1}^K P_i[k] \log_2(P_i[k])$$

(avec par convention $0 \log(0) = 0$)

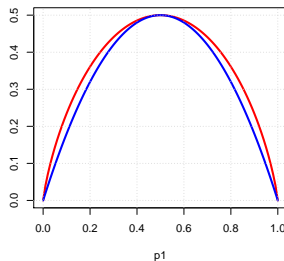
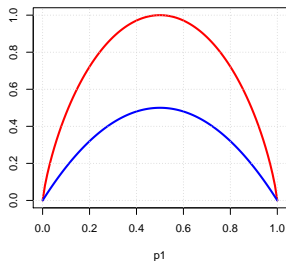
⇒ distribution uniforme : valeurs élevées

⇒ classe(s) majoritaire(s) : valeurs faibles

Arbres de classification - critères d'impureté

Illustration : classification binaire

$$\blacktriangleright P_i = [P_i[1] \ P_i[2]] = [P_i[1] \ 1 - P_i[1]] = [p_1 \ 1 - p_1]$$



► à gauche : **entropie** et **Gini** en fonction de p_1

► à droite : idem avec **entropie** divisée par 2

⇒ critères très similaires

► entropie = mesure du désordre, théorie de l'information

Arbres de classification - critère objectif

Comment **séparer au mieux** les instances affectées au noeud ?
⇒ en cherchant **deux ensembles les plus purs possibles**.

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Arbres de classification - critère objectif

Comment **séparer au mieux** les instances affectées au noeud ?

⇒ en cherchant **deux ensembles les plus purs possibles**.

Critère objectif : moyenne pondérée des impuretés

$$\frac{n_G}{n_i} I_G + \frac{n_D}{n_i} I_D$$

- ▶ $n_i = \#$ d'instances affectées au noeud i (à couper)
- ▶ $n_G, n_D = \#$ d'instances affectées aux noeuds fils
 - ▶ Gauche et Droit
- ▶ $I_G, I_D =$ mesure d'impureté des noeud fils
 - ▶ i.e., entropie ou Gini

Arbres de classification - critère objectif

Comment **séparer au mieux** les instances affectées au noeud ?

⇒ en cherchant **deux ensembles les plus purs possibles**.

Critère objectif : moyenne pondérée des impuretés

$$\frac{n_G}{n_i} I_G + \frac{n_D}{n_i} I_D$$

- ▶ $n_i = \#$ d'instances affectées au noeud i (à couper)
- ▶ $n_G, n_D = \#$ d'instances affectées aux noeuds fils
 - ▶ Gauche et Droit
- ▶ $I_G, I_D =$ mesure d'impureté des noeud fils
 - ▶ i.e., entropie ou Gini

⇒ problème : **minimiser ce critère**

Arbres de classification - critère objectif

Comment **séparer au mieux** les instances affectées au noeud ?

⇒ en cherchant **deux ensembles les plus purs possibles**.

Critère objectif : moyenne pondérée des impuretés

$$\frac{n_G}{n_i} I_G + \frac{n_D}{n_i} I_D$$

⇒ minimiser ce critère = **algorithme glouton** (greedy) :

- ▶ on considère **tous les couples** (variable, seuil)
- ▶ on évalue le critère
- ▶ on sélectionne le **meilleur couple** (variable, seuil)

Rappel : (1 variable / 1 seuil) = principe de l'algorithme

Arbres de classification - critères d'arrêt

Quand arrêter le découpage des noeuds ?

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Arbres de classification - critères d'arrêt

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Quand arrêter le découpage des noeuds ?

Critère d'arrêt automatique : plus de gain en pureté

- ▶ le noeud est déjà parfaitement pur
- ▶ on ne peut pas le purifier davantage

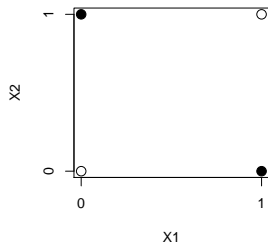
Arbres de classification - critères d'arrêt

Quand arrêter le découpage des noeuds ?

Critère d'arrêt automatique : plus de gain en pureté

- ▶ le noeud est déjà parfaitement pur
- ▶ on ne peut pas le purifier davantage

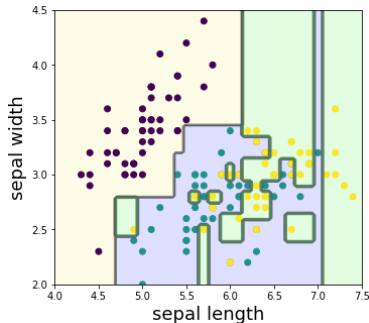
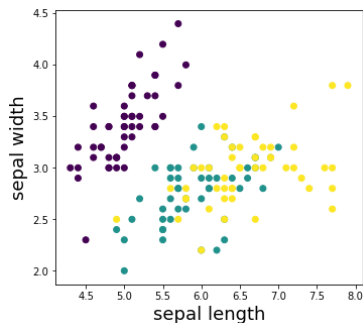
⇒ Exemple : XOR



- ▶ impureté maximum
 - ▶ $p_1 = p_2 = 0.5$
- ▶ constante si on découpe

Arbres de classification - critères d'arrêt

Illustration : Iris dataset



⇒ risque de **sur-apprentissage** si on laisse trop croître l'arbre

Arbres de classification - critères d'arrêt

Eviter le **sur-apprentissage** = limiter la **complexité de l'arbre**
⇒ compromis biais/variance

- ▶ simple : biais élevé ; complexe : variance forte

Arbres de classification - critères d'arrêt

Eviter le **sur-apprentissage** = limiter la **complexité de l'arbre**
⇒ compromis biais/variance

- ▶ simple : biais élevé ; complexe : variance forte

Critères d'arrêt classiques pour l'empêcher de trop croître :

1. s'arrêter quand le **gain en pureté** est marginal
2. critères sur la **structure de l'arbre**
 - ▶ profondeur maximale
 - ▶ nombre de feuilles maximal
 - ▶ nombre minimum d'instances par feuille
 - ▶ nombre minimum d'instances pour couper un noeud
 - ▶ ...

Arbres de classification - critères d'arrêt

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Eviter le **sur-apprentissage** = limiter la **complexité de l'arbre**
⇒ compromis biais/variance

- ▶ simple : biais élevé ; complexe : variance forte

Critères d'arrêt classiques pour l'empêcher de trop croître :

1. s'arrêter quand le **gain en pureté** est marginal
2. critères sur la **structure de l'arbre**
 - ▶ profondeur maximale
 - ▶ nombre de feuilles maximal
 - ▶ nombre minimum d'instances par feuille
 - ▶ nombre minimum d'instances pour couper un noeud
 - ▶ ...

⇒ des **hyperparamètres** à définir et à optimiser

- ▶ e.g., par validation croisée

Pour la [régression](#) : même mécanisme

Critères d'impureté du noeud i :

1. [erreur quadratique moyenne](#) (MSE) :

$$\text{MSE}(i) = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{C_i(j)} - \bar{y}_i)^2$$

où $n_i = \#$ instances affectées au noeud i , \bar{y}_i la moyenne de leurs réponses et $\{C_i(j)\}_{j=1:n_i}$ leurs indices.

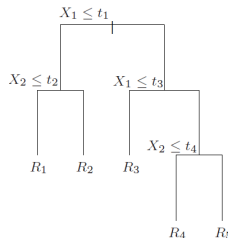
2. [erreur absolue médiane](#) (MAE) :

$$\text{MAE}(i) = \text{median} \left\{ |y_{C_i(j)} - \tilde{y}_i| \right\}$$

où \tilde{y}_i = la médiane des réponses des instances du noeud i .

Mécanisme "top down" pour la **prédiction** :

1. on part de la **racine**
2. on progresse dans l'arbre en évaluant les tests définis aux **noeuds internes**
3. jusqu'à atteindre une **feuille**



⇒ **règles de prédiction** :

- ▶ **classification** : classe majoritaire dans la feuille/région
 - ▶ **prédiction probabiliste** grâce à sa proportion
- ▶ **régression** : valeur moyenne dans la feuille/région

Arbres de décision - remarques

Avantages :

- ▶ interprétabilité
- ▶ simple : pas de préparation des données ; multiclasse

Limites :

- ▶ instabilité
- ▶ fonction non-continue pour la régression

Attention aux prédicteurs catégoriels

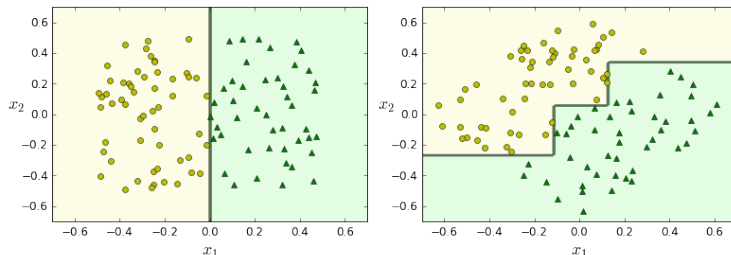
- ▶ relation d'ordre implicite si encodés comme $1, \dots, K$

Approche "non paramétrique"

- ▶ le # de paramètres augmente avec le jeu de données

Algorithme CART le plus répandu (alternatives = C4.5, ID3)

Illustration : **instabilité** à l'"orientation" des données² :



Plus généralement : procédure sensible à de petites variations des données d'apprentissage

⇒ le **bagging** permet de limiter cette instabilité

Bagging

Instabilité des arbres de décisions

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

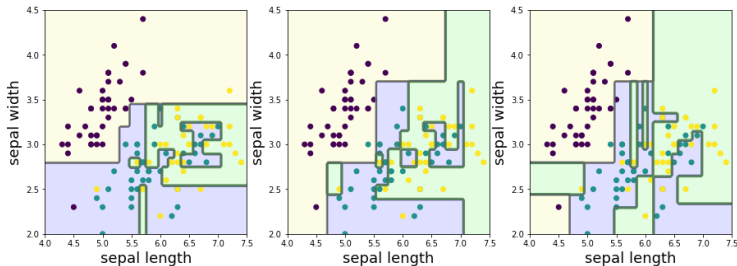
Conclusion

Python

Références

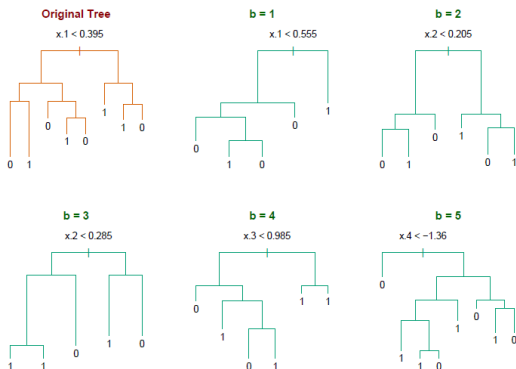
Instabilité des arbres de décision : illustration (1/2)

- jeu de données Iris, 3 échantillons bootstrap



Instabilité des arbres de décision : illustration (2/2)

- ▶ données simulées, 5 descripteurs (Hastie et al., 2001)
- ▶ 5 échantillons bootstrap



Instabilité = **limitation sévère** des arbres de décision

- ▶ compromet la généralisation
- ▶ remet en partie en cause leur interprétabilité

Instabilité = **limitation sévère** des arbres de décision

- ▶ compromet la généralisation
- ▶ remet en partie en cause leur interprétabilité

Stratégie de **bagging** :

1. générer plusieurs **échantillons bootstrap** du jeu d'apprentissage
2. construire les arbres correspondants
3. agréger les prédictions de tous les arbres

⇒ **bagging** = bootstrap **agg**regating

Bagging : principe générique

- ▶ pas spécifique aux arbres...
- ▶ ...mais souvent utilisé avec les arbres

Règles de prédiction (typiques) :

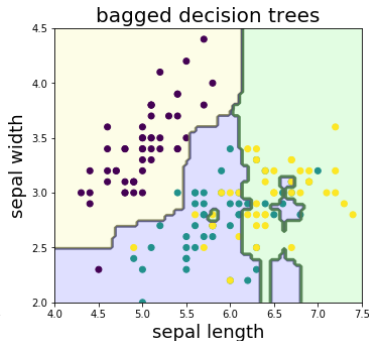
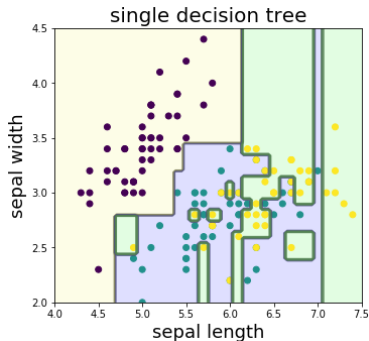
- ▶ **régression** : moyenne des prédictions
- ▶ **classification** : classe majoritaire

Prédiction probabiliste pour la classification ?

- ▶ la proportion de votes peut être optimiste
 - ▶ e.g., si ils prédisent tous la bonne classe avec $p = 0.6$
- ▶ considérer la **moyenne des probabilités** de chaque arbre

Illustration : Iris dataset

- **gauche** : 1 arbre ; **droite** : bagging de 2000 arbres



⇒ capture des caractéristiques plus locales

Bagging et erreur "Out Of Bag"

Outline

Apprentissage Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Echantillon bootstrap : ne contient pas tout le jeu de donnée

Instances hors de l'échantillon = **échantillon indépendant**

- ▶ on peut les classer par l'arbre correspondant

⇒ on peut **estimer la généralisation** lors de l'apprentissage

Bagging et erreur "Out Of Bag"

Echantillon bootstrap : ne contient pas tout le jeu de donnée

Instances hors de l'échantillon = échantillon indépendant

- ▶ on peut les classer par l'arbre correspondant

⇒ on peut estimer la généralisation lors de l'apprentissage

Erreur Out Of Bag (OOB) - pour chaque instance :

1. extraire les arbres qui ne l'ont pas utilisée
2. obtenir les prédictions correspondantes
3. les agréger et comparer à la vraie réponse

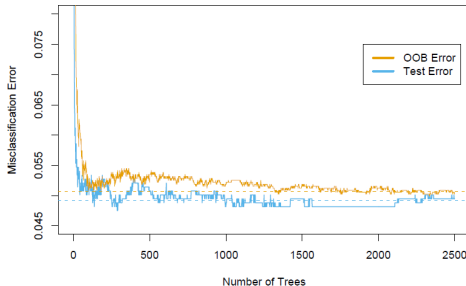
⇒ une caractéristique unique du bagging

Bagging et erreur "Out Of Bag"

Erreur Out Of Bag (OOB) :

- ▶ permet d'estimer la généralisation lors de l'apprentissage
- ▶ en général proche de l'estimation par validation croisée
- ▶ peut aider à savoir quand arrêter la procédure

Illustration (tirée de Hastie et al. (2001)) :



Bagging = bootstrap aggregating

Avantages :

- ▶ stabilité → meilleure performance de prédiction
- ▶ simple à mettre en oeuvre
- ▶ accès à l'erreur Out Of Bag

Inconvénients :

- ▶ coût calculatoire
- ▶ avec les arbres : perte de l'interprétabilité

Random Forests / Forêts Aléatoires

Random Forests

Random Forest = algorithme de bagging basé sur les arbres

- ▶ forêt aléatoire

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Random Forests

Outline

Apprentissage Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Random Forest = algorithme de bagging basé sur les arbres

- ▶ forêt aléatoire

Différence par rapport au bagging "classique" :

quand on ajoute un noeud à un arbre on considère
 $m \leq p$ variables candidates choisies aléatoirement

Random Forest = algorithme de bagging basé sur les arbres

- ▶ forêt aléatoire

Différence par rapport au **bagging "classique"** :

quand on ajoute un noeud à un arbre on considère
 $m \leq p$ variables candidates choisies aléatoirement

⇒ chaque arbre a individuellement un pouvoir prédictif limité

- ▶ i.e., encore plus qu'un arbre classique

⇒ l'étape aléatoire a pour effet de les dé-corréler davantage

- ▶ i.e., encore plus qu'avec uniquement le bootstrap

⇒ les combiner par bagging = un classifieur très performant

Random Forests

Outline

Apprentissage Statistique II

Introduction

CART

Bagging

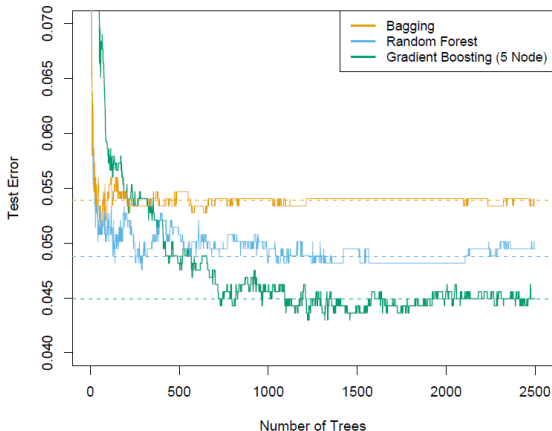
Random Forests

Conclusion

Python

Références

Illustration : SPAMS dataset (Hastie et al., 2001) :



(Gradient boosting : autre algorithme non couvert ici)

Même procédure pour **classification** et **régression**

- ▶ prédiction par bagging - moyenne ou vote majoritaire

Hyper-paramètres :

1. contrôle de la complexité des arbres
 - ▶ idem CART
2. nombre d'arbres à combiner
 - ▶ idem bagging
3. $m \leq p$: nombre de variables candidates à considérer

Hyper-paramètres : à optimiser

- ▶ par validation croisée...ou grâce à l'erreur OOB

Mais :

1. Heuristiques pour le choix de m :
 - ▶ classification : $m = \sqrt{p}$
 - ▶ régression : $m = p/3$
2. nombre d'arbres : en ajouter ne peut pas faire de mal
 - ▶ vrai pour le bagging en général
3. pas nécessaire de contrôler la complexité des arbres (?)
 - ▶ e.g., noeuds de taille min. 1 / 5 pour classif / régression

⇒ simple à paramétrer

⇒ point de départ pour une optimisation plus fine

Critère d'importance des variables

Gain en pureté obtenu au noeud i d'un arbre :

$$\begin{aligned}\Delta(i) &= n_i l_i - (n_G l_G + n_D l_D) \\ &= n_i \left[l_i - \left(\frac{n_G}{n_i} l_G + \frac{n_D}{n_i} l_D \right) \right]\end{aligned}$$

où :

- ▶ le noeud i est coupé en (G, D) (Gauche, Droit)
- ▶ n_i = nombre d'instances affectées au noeud i
- ▶ l_i = critère d'impureté (e.g., Gini, MSE)

$\Rightarrow \sim$ (impureté initiale - moyenne pondérée des impuretés)

\Rightarrow le facteur n_i reflète l'importance du noeud

- ▶ n_i élevé $\rightarrow i$ haut dans l'arbre \rightarrow variable importante

Critère d'importance des variables

Critère d'importance d'une variable : somme du gain en pureté qu'elle apporte dans l'ensemble des arbres

Implémentation :

1. On le calcule pour chaque variable :
 - ▶ on parcourt tous les arbres de la forêt
 - ▶ on parcourt tous les noeuds de l'arbre
 - ▶ si le noeud i est basé sur cette variable on ajoute $\Delta(i)$
2. On le normalise par la somme des valeurs obtenues

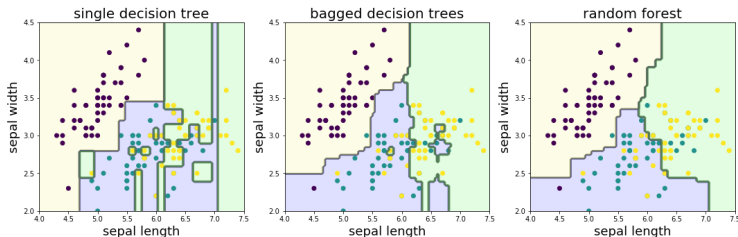
⇒ valeur élevée du critère = variable importante

⇒ somme à 1 sur l'ensemble des variables

(d'autres critères existent, celui-ci est implémenté dans Sk-Learn)

Illustration : Iris dataset

- ▶ **gauche** : 1 arbre ; **milieu** : bagging ; **droite** : random forest



⇒ **Out Of Bag** accuracy :

- ▶ bagging = 74.7%
- ▶ random forest = 76%

Random Forests - remarques

Outline

Apprentissage Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Random Forests : bagging basé sur les arbres de décision

- ▶ sélection aléatoire des variables pour **décorréliser les arbres**

Avantages :

- ▶ performant sur de nombreux problèmes
- ▶ simple à paramétrer (heuristiques)
- ▶ accès à l'erreur Out Of Bag (approche bagging)
- ▶ critère d'importance des variables

Inconvénients

- ▶ coût calculatoire
- ▶ modèle "boîte noire"
 - ▶ +/- car critère d'importance des variables...

Conclusion

Conclusion / bilan (1/2)

Outline

Apprentissage Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Ce cours : des arbres aux forêts aléatoires

Arbres de régression et classification

- ▶ algorithme CART
- ▶ critères d'impureté (Gini, entropie ; MSE, MAE)
- ▶ interprétabilité du modèle
- ▶ sur-apprentissage / instabilité
- ▶ critères de contrôle de complexité

Bagging

- ▶ bootstrap aggregating
- ▶ principe général, décliné ici aux arbres
- ▶ stratégie ensembliste : vote majoritaire / moyenne
- ▶ erreur "Out Of Bag" accessible lors de l'apprentissage

Random Forests

- ▶ bagging d'arbres
- ▶ dimension stochastique pour décorréler les arbres
 - ▶ arguments théoriques
- ▶ critère d'importance des variables
- ▶ performant et simple à paramétrer (heuristiques)

Mise en oeuvre dans Scikit-Learn

Arbres de classification et régression (1/3)

Deux classes du module `tree` :

- ▶ `DecisionTreeClassifier` → classification
- ▶ `DecisionTreeRegressor` → régression

Paramètres principaux (constructeur) :

- ▶ `criterion` : critère d'impureté
 - ▶ `gini` (défaut) ou `entropy`
 - ▶ `mse` (défaut) ou `mae`
- ▶ critères de contrôles de complexité :
 - ▶ `max_depth`, `min_samples_split`,
`min_samples_leaf`, `max_leaf_nodes`, ...
 - ▶ `min_impurity_decrease`
- ▶ (+ `max_features` : pour réduire le nombre de variables candidates - pour les random forests)

Arbres de classification et régression (2/3)

Méthodes principales :

- ▶ `fit` et `predict`
- ▶ `DecisionTreeClassifier` : `predict_proba` et `predict_log_proba`
 - ▶ `proba` = proportion de la classe majoritaire
- ▶ `decision_path` : chemin suivi par une/des instance(s)

Exemple :

```
# import model
from sklearn.tree import DecisionTreeClassifier
# instancie model
tree = DecisionTreeClassifier(max_depth=3)
# fit model
tree.fit(X_train, y_train)
# make predictions
preds = tree.predict(X_test)
probs = tree.predict_proba(X_test)
```

Arbres de classification et régression (3/3)

Outline

Apprentissage
Statistique II

Introduction

CART

Bagging

Random Forests

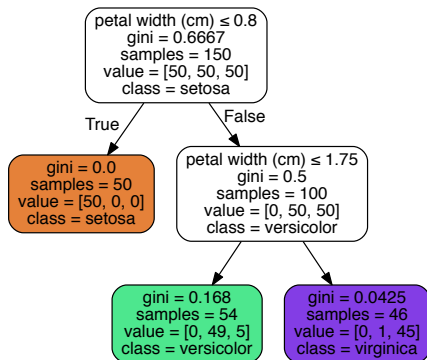
Conclusion

Python

Références

Outils de **visualisation** :

```
from sklearn.tree import export_graphviz
import graphviz
dot_data = tree.export_graphviz(tree, out_file=None)
graph = graphviz.Source(dot_data)
graph.render("figure.pdf")
```



(NB : à customiser avec le nom des classes et des variables)

Bagging (1/2)

Deux classes du module ensemble :

- ▶ `BaggingClassifier` → classification
- ▶ `BaggingRegressor` → régression

Paramètres principaux (constructeur) :

- ▶ `base_estimator` : prédicteur de base
 - ▶ par défaut, un arbre de décision
- ▶ `n_estimators` : nombre de prédicteurs à aggréger
- ▶ `bootstrap` : booléen (vrai par défaut)
 - ▶ si faux, on tire sans remise : on parle de pasting
- ▶ `oob_score` : booléen indiquant si on veut le calculer
 - ▶ on le trouve alors dans le champ `oob_score_`

Bagging (2/2)

Méthodes principales :

- ▶ `fit` et `predict`
- ▶ `BaggingClassifier` : `predict_proba` et `predict_log_proba`

Exemple :

```
from sklearn.ensemble import BaggingClassifier
# define model
bag = BaggingClassifier(
    DecisionTreeClassifier(), n_estimators = 200,
    oob_score = True, random_state = 22)
# fit
bag.fit(X, y)
```


Random Forest (1/2)

Outline

Apprentissage Statistique II

Introduction

CART

Bagging

Random Forests

Conclusion

Python

Références

Deux classes du module ensemble :

- ▶ `RandomForestClassifier` → classification
- ▶ `RandomForestRegressor` → régression

Paramètres principaux (constructeur) :

- ▶ `n_estimators` : nombre d'arbres dans la forêt
- ▶ `criterion` : critère d'impureté
- ▶ `max_features` : paramètre m
 - ▶ défaut : \sqrt{p} pour classification, p pour régression
- ▶ critères de contrôle de la complexité des arbres
 - ▶ `max_depth`, `min_samples_split`,
`min_samples_leaf`, `max_leaf_nodes`, ...
 - ▶ `min_impurity_decrease`
- ▶ `oob_score` : booléen indiquant si on veut le calculer

Random Forest (2/2)

Méthodes principales :

- ▶ `fit` et `predict`
- ▶ `RandomForestClassifier` : `predict_proba` et `predict_log_proba`

Attribut particulier : `feature_importances_`

- ▶ le critère d'importance des variables (après appel à `fit`)

Exemple :

```
# fit random model
from sklearn.ensemble import RandomForestClassifier
# instantiate model
rf = RandomForestClassifier(n_estimators = 500, oob_score = True, random_state= 20)
# learn model
rf.fit(X, y)
```

- A. Géron. *Hands-On Machine Learning with Scikit-Learn & TensorFlow*. O'Reilly, 2017.
- T. Hastie, R. Tibshirani, and J.. Friedman. *The Elements of Statistical Learning*. Springer, 2001.