

# Arbres de Décision

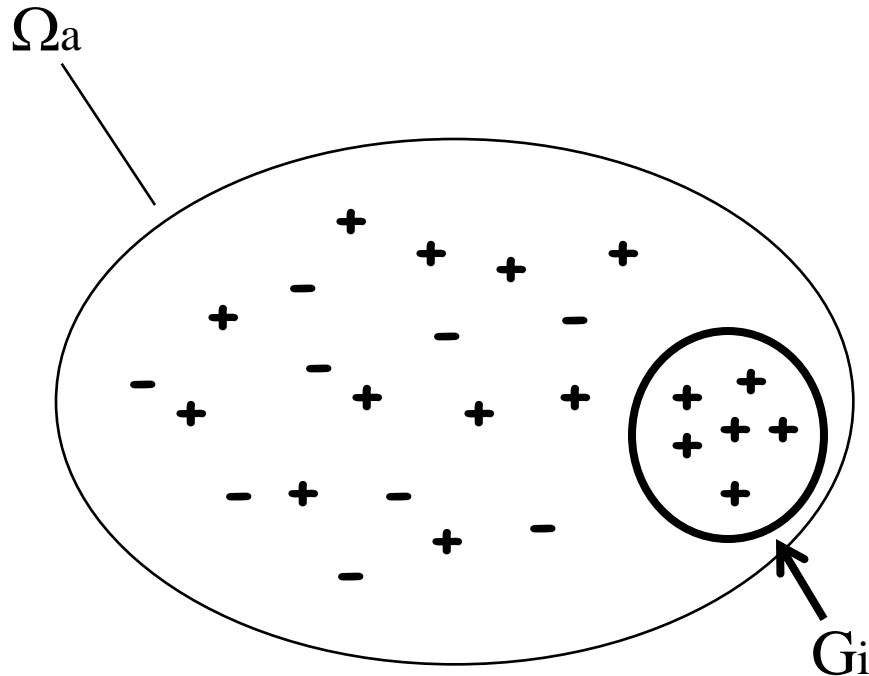
Ricco RAKOTOMALALA

ricco.rakotomalala@univ-lyon2.fr

# Construction d'un arbre de décision

# Arbres de décision – Apprentissage par partitionnement

Objectif : on veut construire des sous-groupes les plus « homogènes » du point de vue de la variable à prédire



La variable qualitative Y prend ses valeurs dans {+,-}

Le sous-groupe  $G_i$  est complètement pur du point de vue de Y, il ne possède que des individus portant la valeur + de Y

si ( $\omega \in G_i$ ) alors ( $Y = +$ )

L'idée est de trouver le plus rapidement Possible (avec le moins de variables) des groupes où  $P(Y=+) \neq 1$

La description des sous-groupes repose sur :

- ☞ la fonction f et ses paramètres éventuels  $\alpha$
- ☞ les variables exogènes  $X_i$

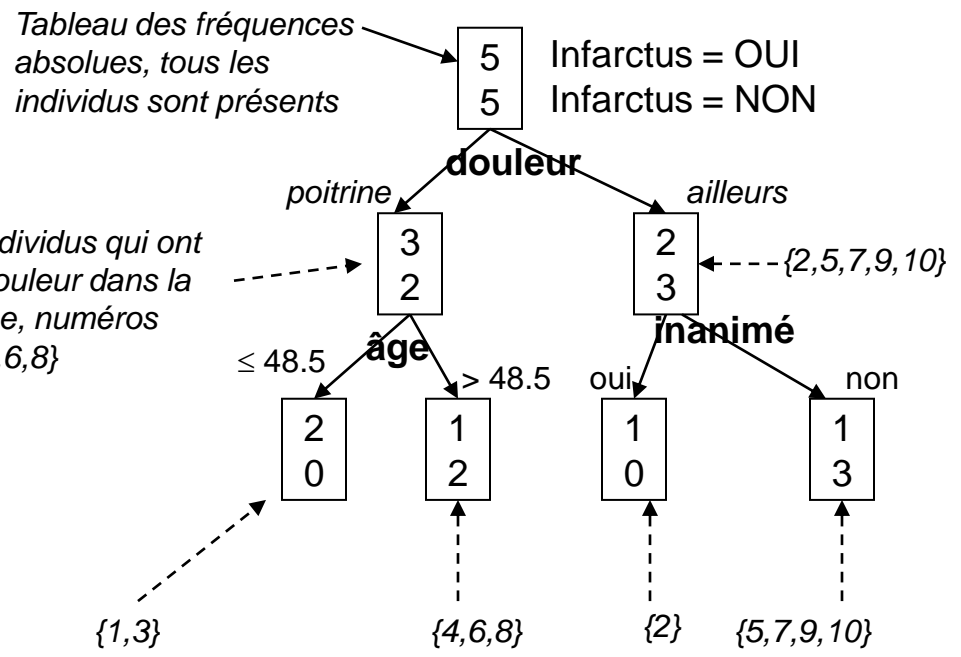
# Arbres de décision – Un exemple

Numéro	Infarctus	Douleur	Age	Inanimé
1	oui	poitrine	45	oui
2	oui	ailleurs	25	oui
3	oui	poitrine	35	non
4	oui	poitrine	70	oui
5	oui	ailleurs	34	non
6	non	poitrine	60	non
7	non	ailleurs	67	non
8	non	poitrine	52	oui
9	non	ailleurs	58	non
10	non	ailleurs	34	non

Y

X

- A résoudre :
- choix de la variable de segmentation
  - traitement des variables continues
  - règle d'arrêt dans la construction
  - décision sur une feuille



Les individus qui ont une douleur dans la poitrine, numéros {1,3,4,6,8}

Premier sous-groupe, complètement homogène du point de vue de la variable à prédire : il est constitué exclusivement d'individus qui ont un infarctus



# Arbres de décision – Choix de la variable de segmentation

On choisit la variable  $X^*$  telle qu'elle est la plus liée (corrélée) avec  $Y$   
 👉 on utilise la quantité du  $\chi^2$  calculée sur le tableau de contingence (croisement de  $Y$  avec  $X_i$ ) pour quantifier cette liaison

	$x_{i,1}$	...	$x_{i,L_i}$
$Y_1$			
$\vdots$		$n_{k,l} = \text{card}(\{\omega \in \Omega_a / Y(\omega) = Y_k \text{ et } X_i(\omega) = X_{i,l}\})$	
$Y_K$			

$$X^* = \arg \max_{i=1, \dots, p} \chi^2_{Y, X_i}$$

**Amélioration** : la mesure du  $\chi^2$  augmente avec

- 👉  $n$ , l'effectif sur le nœud à segmenter
- 👉 le nombre de lignes
- 👉 le nombre de colonnes

Ces valeurs sont constantes dans les comparaisons deux à deux du  $\chi^2$

Les variables qui ont beaucoup de modalités (et ainsi induisent beaucoup de colonnes dans le tableau de contingence) sont avantagés

$$t^2_{Y, X_i} = \frac{\chi^2_{Y, X_i}}{n \sqrt{(K-1)(L_i-1)}}$$

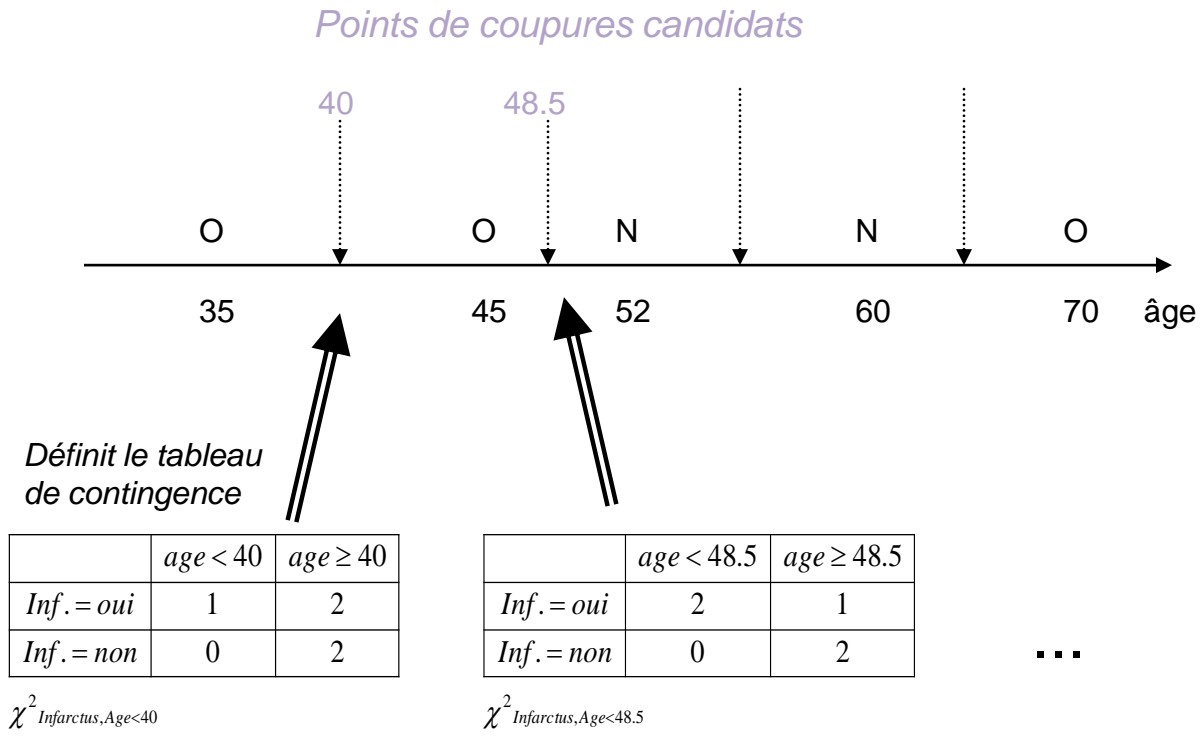
(le  $t$  de Tschuprow varie entre 0 et 1)

# Arbres de décision – Traitement des variables continues

Comment est réalisé le choix du point de coupure  
(ex: d'où vient la valeur 48.5 de découpage de l'âge dans l'arbre exemple)

Point de coupure : borne de discrétisation

- il doit toujours être situé entre deux points consécutifs sur l'axe de la variable quantitative
- il permet de définir un tableau de contingence



# Arbres de décision – Règle d'arrêt

Quand décider qu'un sommet devient une feuille ?

Homogénéité des groupes : critère de précision (confiance)

Pureté d'un sommet

Seuil de spécialisation (ex. si une classe est représentée à 98% -> stop)

Effectif des groupes : critère de support

Taille minimale pour segmenter (ex. en dessous de 10 obs, on ne segmente plus)

Effectif d'admissibilité (ex. si un des sommets produit couvre moins de 2 obs. -> refus)

Test d'indépendance du CHI-2 : démarche statistique

$$\begin{cases} H_0 : Y \text{ et } X^* \text{ indépendants} \\ H_1 : Y \text{ est lié avec } X^* \end{cases}$$

Comment fixer le  
risque du test ?

L'idée est surtout de contrôler la profondeur de l'arbre !

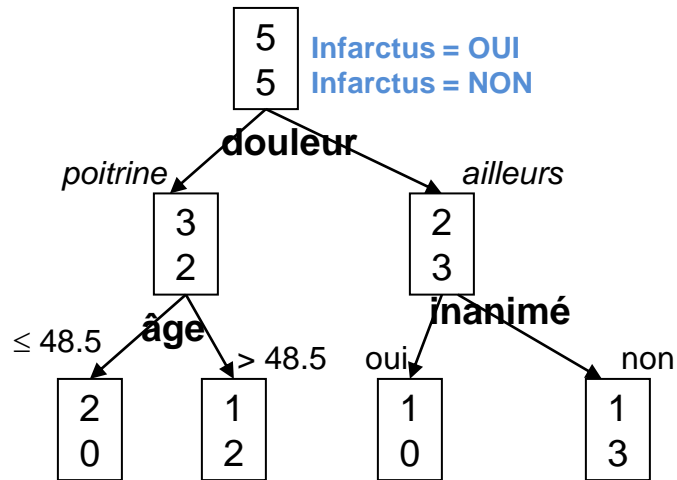
La taille de l'arbre influe fortement sur ses performances

# Les arbres de décision dans la pratique

## Post-traitement, analyse des résultats



# Lecture des règles (1)



Chaque chemin partant de la racine vers les feuilles constitue une règle → 4 règles ici

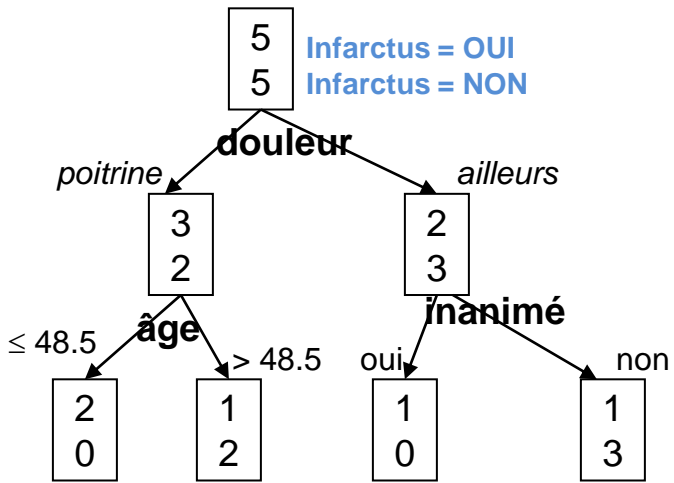
- A. Les règles sont mutuellement exclusives c.-à-d. un individu ne peut déclencher qu'une et une seule règle.
- B. L'ensemble de règles couvre tout l'espace des valeurs possibles c.-à-d. tout individu à classer va déclencher une des règles.

1. **SI** douleur = poitrine **ET** âge ≤ 48.5 **ALORS** infarctus = oui
2. **SI** douleur = poitrine **ET** âge > 48.5 **ALORS** infarctus = non
3. **SI** douleur = ailleurs **ET** inanimé = oui **ALORS** infarctus = oui
4. **SI** douleur = ailleurs **ET** inanimé = non **ALORS** infarctus = non

⇒ **SI** prémisse **ALORS** conclusion

Conjonction de propositions de type « attribut .comparaison. valeur »

# Lecture des règles (2)



L'arbre peut se lire également comme une cascade de règles (propositions) imbriquées.

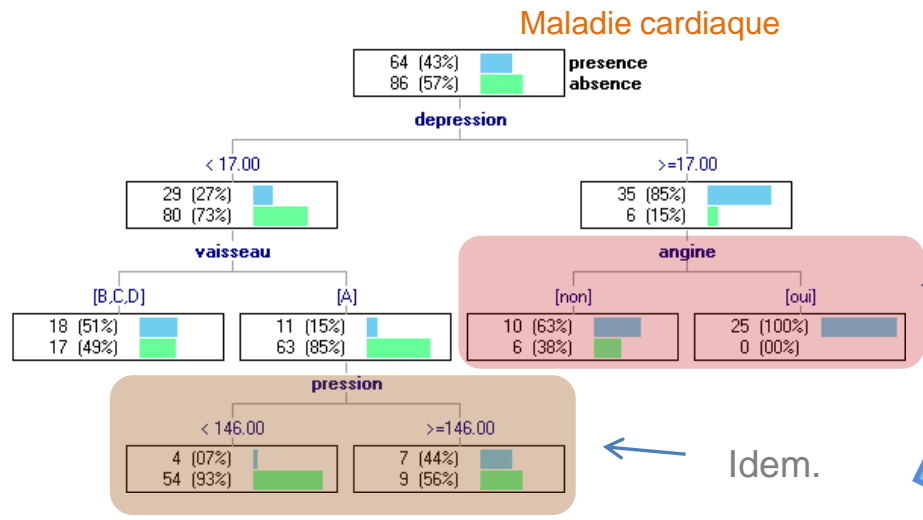
SINON signifie que « douleur ≠ poitrine » c.-à-d. « douleur = ailleurs » dans notre fichier de données.

**SI** poitrine = douleur  
**ALORS**  
**SI** âge ≤ 48.5  
     **ALORS** infarctus = oui  
     **SINON** infarctus = non  
**SINON**  
**SI** inanimé = oui  
     **ALORS** infarctus = oui  
     **SINON** infarctus = non

Mode de lecture intéressant si l'on souhaite traduire les règles en fonction SI(...) sous Excel.



# Post-élagage manuel – Suppression des feuilles non pertinentes

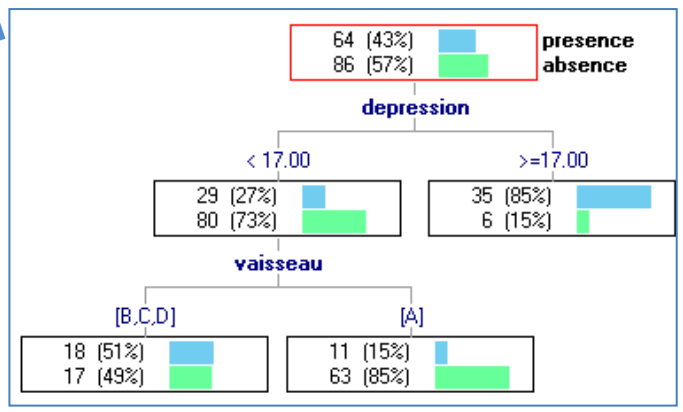


➔ A priori, on aurait 5 feuilles = 5 règles.

Est-ce que ces 2 feuilles sont pertinentes ?  
 Est-ce qu'il est vraiment nécessaire d'induire deux règles à partir de ces feuilles ?



Idem.



Arbre à 3 règles. Totalement équivalent du point de vue du classement.

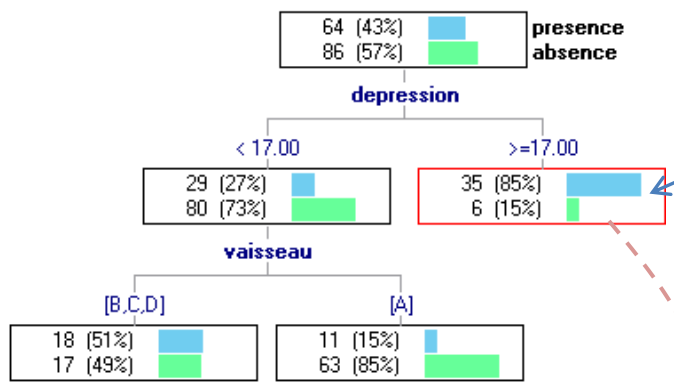
**Post-élagage.** Suppression des feuilles sœurs (issues du même père) portant des conclusions identiques. Suppression de proche en proche c.-à-d. on part des feuilles et on remonte jusqu'à la racine. On arrête lorsque la suppression n'est plus possible.



Certaines méthodes de construction d'arbres intègrent une procédure automatisée de post-élagage (ex. C4.5, CART)

# Expertise des règles – Rôle des autres variables quantitatives

Une règle délimite une sous-population, décrite par la succession de propositions partant de la racine à la feuille.



La sous-population est caractérisée par les variables apparaissant sur le chemin. Est-ce que les autres variables (non-sélectionnées sur le chemin) nous permettent de mieux comprendre la formation du groupe ?

Ces personnes ont une « depression » plus élevée que les autres. Prévisible oui/non au vu de l'arbre ?

Ces personnes ont un « pic » plus élevé. On pouvait le deviner ?

Elles ont un « taux\_max » significativement plus bas. On pouvait le deviner ?

Informations on : Level 2, Node 2

IF depression >=17.00

Characterization

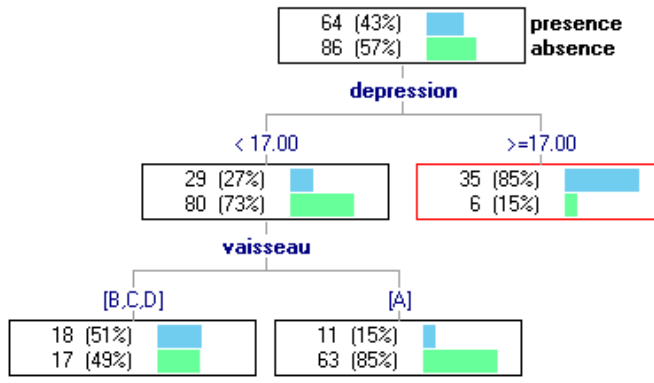
Continuous attributes

Attribute	Strength	Local Avg	Global Avg
depression	9.96	25.1707	10.3600
pic	6.16	2.0488	1.5600
age	1.79	57.0488	54.7467
pression	1.23	135.3415	132.2267
cholester	0.76	254.1951	248.6867
taux_max	-4.60	135.0488	149.5733

41 examples (27.33% of the learning set)

**VALEUR TEST** ( $\approx t$  de Student du test de comparaison de moyennes)

# Expertise des règles – Rôle des autres variables qualitatives



**J-Measure.** Pour distinguer les variables les plus caractérisantes.

**Valeur-Test** (≈ comparaison de proportions [80% vs. 63%]) pour distinguer les modalités les plus caractérisantes.

170 individus en tout dans le fichier. 95 hommes et 55 femmes.

Informations on : Level 2, Node 2

IF depression >=17.00

Characterization

Continuous attributes

coeur ( 0.1547 )				
Values	Strength	Local Dist.	Global Dist.	Recall
presence	6.46	35 (85%)	64 (43%)	55%
absence	6.46	6 (15%)	86 (57%)	7%

vaisseau ( 0.0673 )				
Values	Strength	Local Dist.	Global Dist.	Recall
B	-0.09	8 (20%)	30 (20%)	27%
A	-3.73	14 (34%)	88 (59%)	16%
C	3.82	13 (32%)	21 (14%)	62%
D	2.10	6 (15%)	11 (7%)	55%

angine ( 0.0506 )				
Values	Strength	Local Dist.	Global Dist.	Recall
non	-3.90	16 (39%)	96 (64%)	17%
oui	3.90	25 (61%)	54 (36%)	46%

sexe ( 0.0275 )				
Values	Strength	Local Dist.	Global Dist.	Recall
masculin	2.67	33 (80%)	95 (63%)	35%
feminin	-2.67	8 (20%)	55 (37%)	15%

41 exemples (27.33% of the learning set)

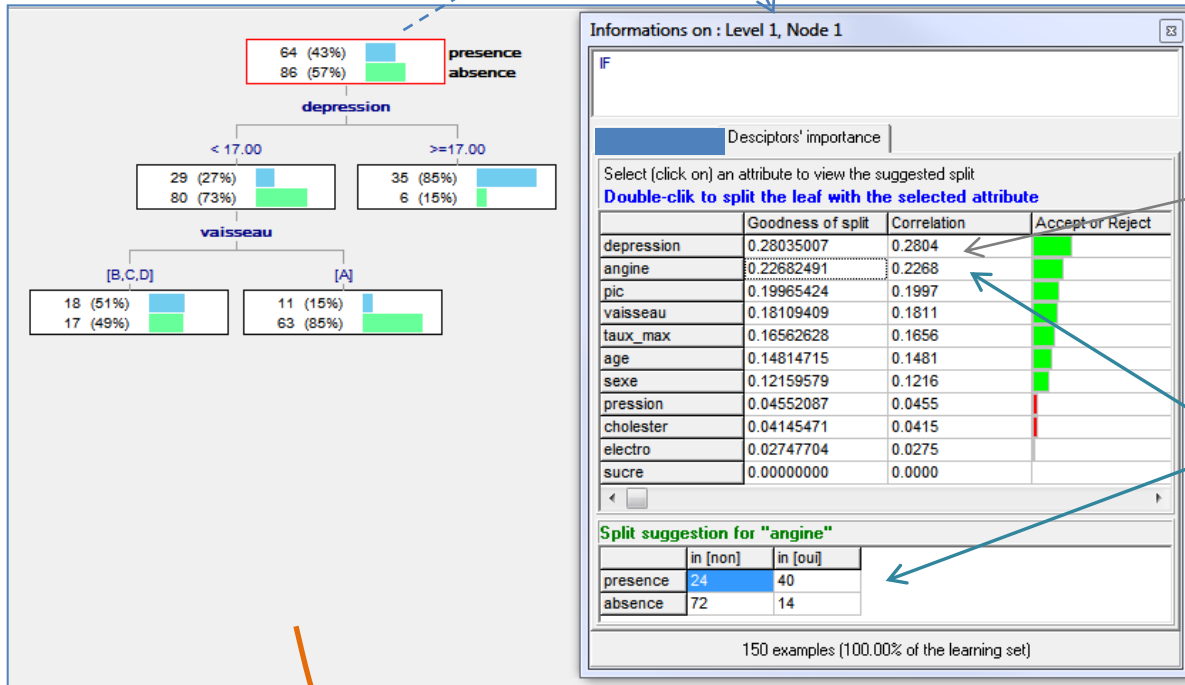
$P(\text{Homme} / \text{Sommet}) = 80\%$

$P(\text{Sommet} / \text{Homme}) = 35\%$

Sur le sommet (c.-à-d. individus avec « depression ≥ 17 »), il y a 41 personnes, dont 33 hommes et 8 femmes.

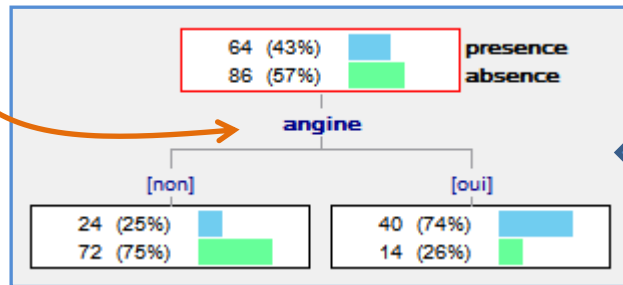


# Construction interactive des arbres



« depression » est la meilleure variable de segmentation avec un impact = 0.2804

« engine » est la 2<sup>ème</sup> meilleure variable avec impact = 0.2268. On obtiendrait le partitionnement ci-dessous.



On peut modifier interactivement l'arbre compte tenu de la qualité numérique (impact) et de la pertinence par rapport aux connaissances du domaine.



On associe connaissances expertes et critères numériques pour construire un modèle efficace et interprétable.



Très peu d'outils libres proposent la construction interactive des arbres.

# Arbres de décision - Bilan

# Arbres de décision

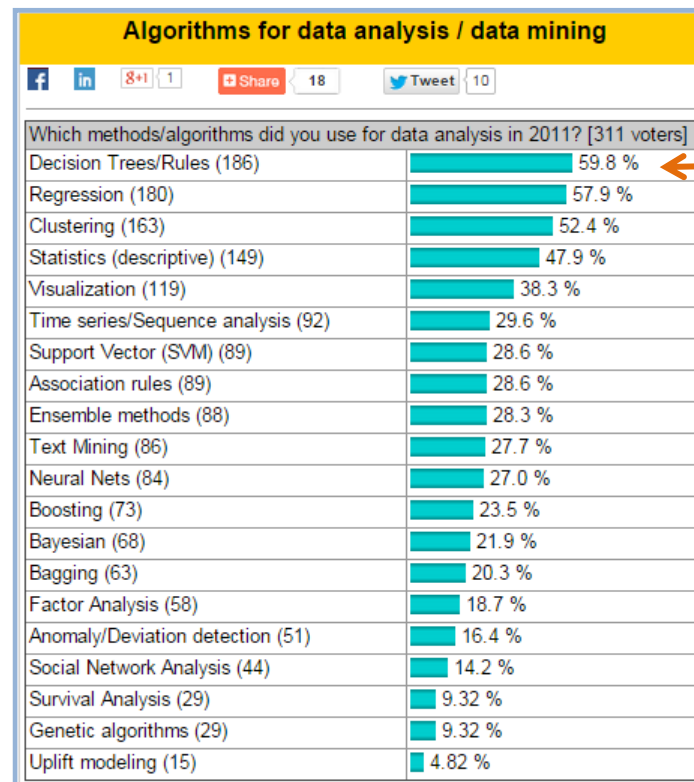
## Une méthode très populaire

### Kdnuggets

« [Methods/algorithms used for data analysis/data mining in 2011](#) », Oct-Nov, 2011.

### Site data mining et scoring

« [Quelle méthode de data mining utilisez-vous le plus ?](#) »





# Pourquoi cette popularité ? Avantages et inconvénients

## Avantages :

- connaissances « intelligibles » -- validation d'expert (si arbre pas trop grand)
- traduction directe de l'arbre vers une base de règles
- sélection automatique des variables pertinentes
  
- non paramétrique
- traitement indifférencié selon le type des variables prédictives
- robuste face aux données aberrantes, solutions pour les données manquantes
- robuste face aux variables redondantes
- rapidité et capacité à traiter des très grandes bases
  
- enrichir l'interprétation des règles à l'aide des variables non sélectionnées
- possibilité pour le praticien d'intervenir dans la construction de l'arbre

## Inconvénients :

- problème de stabilité sur les petites bases de données (feuilles à très petits effectifs)
- recherche « pas-à-pas » : difficulté à trouver certaines interactions (ex. xor)
- peu adapté au « scoring »
- performances moins bonnes en général par rapport aux autres méthodes (en réalité, performances fortement dépendantes de la taille de la base d'apprentissage)

# Arbres de décision - Logiciels

The screenshot shows the SIPINA software interface. The main window displays a decision tree for the 'engine' attribute. The tree structure is as follows:

- Root Node: 64 (43%) presence, 86 (57%) absence
- Level 1: 'engine' attribute
  - [non]: 24 (25%) presence, 72 (75%) absence
  - [oui]: 40 (74%) presence, 14 (26%) absence
- Level 2, Node 1: Selected node in the [non] branch.

The 'Informations on : Level 2, Node 1' panel shows the following table:

Attribute	Goodness of split	Correlation	Accept
age	0.18881119	0.1888	
vaisseau	0.18279570	0.1828	
taux_max	0.16483516	0.1648	
depression	0.15000000	0.1500	
pic	0.09371720	0.0937	
sexe	0.07113463	0.0711	
cholester	0.05427230	0.0543	
pression	0.04512100	0.0451	
sucre	0.00000000	0.0000	
engine	0.00000000	0.0000	
electro	0.00000000	0.0000	

The 'Split suggestion for "age"' panel shows the following table:

	< 53.50	>= 53.50
presence	2	22
absence	42	30

96 examples (64.00% of the learning set)

- + Association forte avec les tableurs (Excel, Libre / Open Office)
- + Le seul outil gratuit à proposer des fonctionnalités utilisatrices interactives de qualité
- Outil universitaire, peu d'outils pour l'exploitation des résultats et le reporting
- Stand alone, peu d'interaction avec les autres méthodes de data mining

SPAD 8.0 - Interactive Tree

Projet Edition Affichage Diagramme Dessin Outils Aide 35%

Interactif Tree

- Diagrammes de traitements
  - Diagramme
- Modèles
  - tree\_result
  - acp\_result

Méthodes

- Analyses factorielles
- Classifications - Typologies
- Amado - Graphiques de Bertin
- Scoring et Modélisation
- Arbres de décision - Segmentations
  - Arbres de décision
  - ICT - Typologie par Arbre de décision
- Tableaux multiples
- Text Mining

Exécutions

Niveau	Elément	Stop
1	ACP	✓
0	iris.iris	✓
2	Archivage axes et ...	✓
3	Arbres de décision	✓
4	Archivage du mod...	✓
2	Archivage du mod...	✓

Visualisation graphique des données

Graphiques Edition Affichage Outils Aide 21%

Galerie Arbre de décision

Rapport Graphique

setosa  
versicolor  
virginica  
Echantillon : Apprentissage

1 36 / 34%  
ACP1 < -1.159  
100%

0 106 / 100%  
Racine

2 70 / 66%  
ACP1 >= -1.159  
49% 51%

3 44 / 42%  
ACP1 < 1.367  
77%

4 26 / 25%  
ACP1 >= 1.367  
100%

Distribution

Effectif : 44 Tableau standard

Modal...	Effectif	Pourc...	Rappel...	Rappel...
setosa	0	0	0	0
versic...	34	77.273	100	100
virgini...	10	22.727	27.778	27.778

Variable de découpage potentiel

Variable	Qu...	Rôle	Comment...
ACP1	0.5721		
ACP3	0.4824		
ACP4	0.3512		Non signi...
ACP2	0.1942		Non signi...

Variables continues

Variable	Valeu...	Moyenn...	Moyenn...
ACP2	5.093	-0.03	0.544
ACP1	3.11	-0.014	0.611
ACP4	2.431	0.004	0.044
ACP3	-0.277	0.002	-0.011

Variables nominales (J-Mesure)

Variable	J-Mesure
type	0.352

Modalités (Valeur-test)

Filtre  Aucun  [Valeur-test] > 2.0

Variable	Modal...	Val...	% racine	% noe...
type	versic...	8.359	32.075	77.273
type	virgini...	-2.048	33.962	22.727
type	setosa	-6.191	33.962	0

Moteur de calcul prêt

+ Fonctionnalités interactives  
 + La méthode (les arbres de décision) s'inscrit dans un tout cohérent. Elle peut coopérer avec les autres techniques de data mining (ex. ici arbre à partir des axes factoriels)

➡ Idem pour les autres outils phares du data mining : SAS EM, IBM SPSS Modeler, etc. (Voir [youtube](#))

# R – Plusieurs packages spécialisés (tree, rpart, party, ...)

```
> print(a1)
n= 150

node), split, n, loss, yval, (yprob)
* denotes terminal node

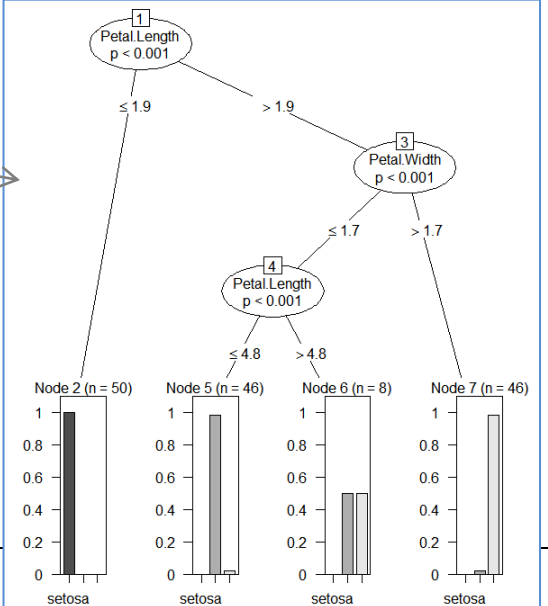
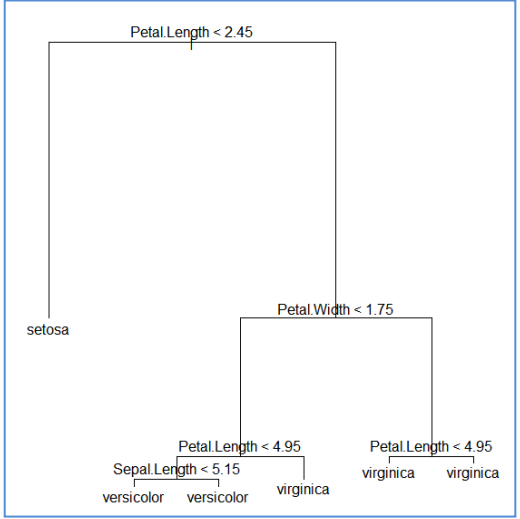
1) root 150 100 setosa (0.33333333 0.33333333 0.33333333)
2) Petal.Length < 2.45 50 0 setosa (1.00000000 0.00000000 0.00000000) *
3) Petal.Length >= 2.45 100 50 versicolor (0.00000000 0.50000000 0.50000000)
6) Petal.Width < 1.75 54 5 versicolor (0.00000000 0.90740741 0.09259259) *
7) Petal.Width >= 1.75 46 1 virginica (0.00000000 0.02173913 0.97826087) *
```

```
data(iris)

library(rpart)
a1 <- rpart(Species ~ ., data = iris)
print(a1)

library(tree)
a2 <- tree(Species ~ ., data = iris)
plot(a2)
text(a2)

library(party)
a3 <- ctree(Species ~ ., data = iris)
plot(a3)
```



Aucune fonctionnalité interactive !

**TANAGRA 1.4.50 - [Supervised Learning 1 (C-RT)]**

File Diagram Component Window Help

Dataset (VisaPremier.xls)

- Define status 1
  - Supervised Learning 1 (C-RT)
  - Supervised Learning 2 (C4.5)
  - Supervised Learning 3 (ID3)
  - Supervised Learning 4 (CS-CRT)

**Tree description**

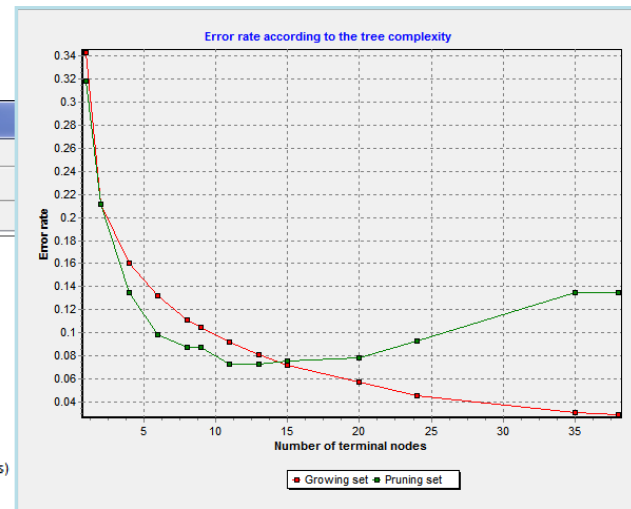
Number of nodes: 21  
Number of leaves: 11

**Decision tree**

- moycred3 < 20.5000
  - csp in [Pcad,Part,Pret]
    - anciente < 7.5000 then CARTEVP = oui (88.46 % of 26 examples)
    - anciente >= 7.5000
      - agemvt < 12.5000
        - nbc b < 0.5000 then CARTEVP = non (88.24 % of 17 examples)
        - nbc b >= 0.5000 then CARTEVP = oui (100.00 % of 15 examples)
      - agemvt >= 12.5000
        - agemvt < 14.5000 then CARTEVP = non (97.14 % of 70 examples)
        - agemvt >= 14.5000
          - nbc b < 0.5000 then CARTEVP = non (83.33 % of 24 examples)
          - nbc b >= 0.5000 then CARTEVP = oui (80.00 % of 15 examples)
    - csp in [Psan,Pemp,Pouv,Pinc,Pagri] then CARTEVP = non (92.97 % of 123 examples)
  - moycred3 >= 20.5000
    - agemvt < 12.5000 then CARTEVP = oui (98.57 % of 70 examples)
    - agemvt >= 12.5000
      - agemvt < 13.5000
        - mtfactor < 82940.0000 then CARTEVP = non (85.00 % of 60 examples)
        - mtfactor >= 82940.0000 then CARTEVP = oui (77.78 % of 9 examples)
      - agemvt >= 13.5000 then CARTEVP = oui (81.82 % of 99 examples)

**Components**

Data visualization	Statistics	Nonparametric statistics	Instance selection	Feature construction	Feature selection
Regression	Factorial analysis	PLS	Clustering	Spv learning	Meta-spv learning
Spv learning assessment	Scoring	Association			
Binary logistic regression	C-PLS	CS-MC4	Decision List		Linear discriminant analy
BVM	C-RT	C-SVC	ID3		Log-Reg TRIRLS
C4.5	CS-CRT	CVM	K-NN		Multilayer perceptron



Plusieurs approches, dont la méthode CART

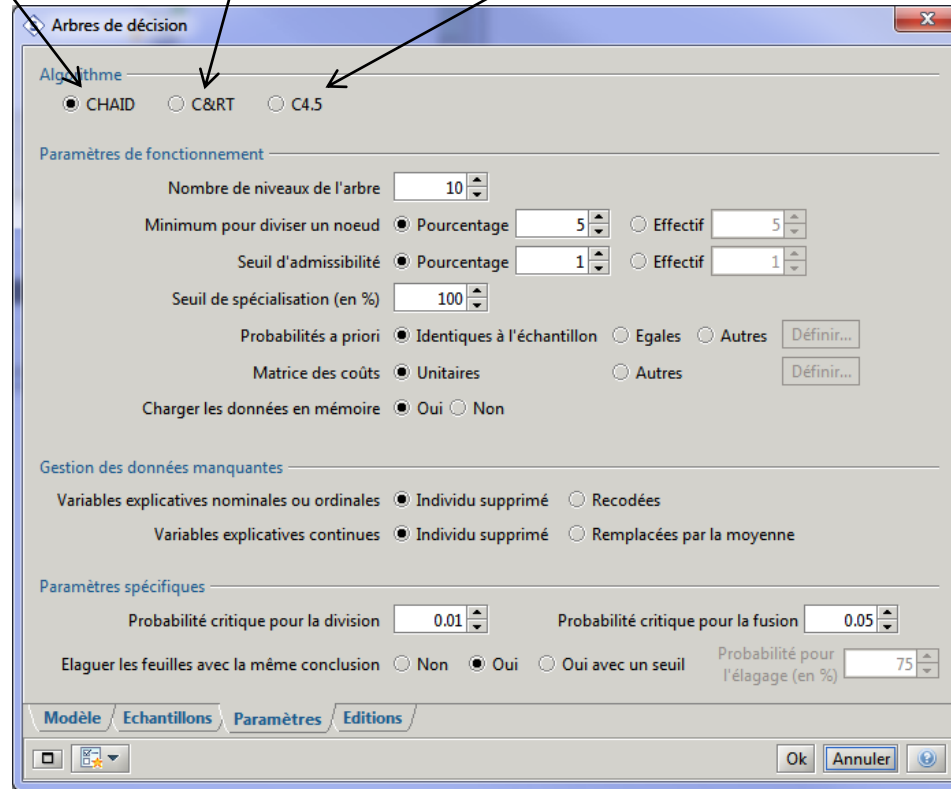
Certaines peuvent prendre en compte les coûts de mauvais classement

Pas de construction interactive des arbres

CART et apparentés (Breiman et al., 1984)

CHAID et apparentés (Kass, 1980)

C4.5 et apparentés (Quinlan, 1993)



Les méthodes se distinguent par : (1) les mesures d'évaluation des segmentations ; (2) les stratégies de regroupement durant la segmentation ; (3) la détermination de la taille « optimale » de l'arbre.

# Bibliographie - Tutoriels



- « [L'add-in Sipina pour Excel 2007 et 2010](#) », août 2014.
- « [Sipina add-on pour OoCalc](#) », mars 2012.
- « [Apprentissage et test avec Sipina](#) », mars 2008.
- « [Analyse interactive avec Sipina](#) », mars 2008.
- « [Sipina – Traitement des très grands fichiers](#) », octobre 2009.
- « [Le format PMML pour le déploiement de modèles](#) », septembre 2010.
- « [Arbres de décision interactifs avec SPAD](#) », janvier 2010.
- « [Nouveaux arbres interactifs dans SPAD 8](#) », août 2014.
- « [Introduction à R – Arbre de décision](#) », mars 2012.

Etc., etc., et n'oublions pas les tutoriels sur [youtube](#)...

# Bibliographie : arbres de décision

« [Arbres de Décision](#) », R. Rakotomalala, Revue MODULAD, 33:163-187, 2005

« Graphes d'induction », D. Zighed et R. Rakotomalala, Hermès, 2000.

« C4.5 – Programs for Machine Learning », Quinlan, 1993.

« Classification and Regression Tree », L. Breiman, J. Friedman, R. Olshen et C. Stone, 1984.

Autres références : <http://sipina-arbres-de-decision.blogspot.fr/p/references.html>