

# Une approche probabiliste pour le classement d'objets incomplètement connus dans un arbre de décision

Lamis HAWARAH

Laboratoire TIMC-IMAG  
Grenoble

22 octobre 2008



# Les valeurs manquantes

# Les valeurs manquantes

## Problème

- ▶ Analyse statistique des données
- ▶ Aide à la décision à partir de données

# Les valeurs manquantes

## Problème

- ▶ Analyse statistique des données
- ▶ Aide à la décision à partir de données

## Fouille de données (Data Mining)

### *Construction et exploitation du modèle*

- ▶ Réseaux bayésiens
- ▶ Règles d'association
- ▶ Réseaux de neurones
- ▶ **Arbres de Décision**
- ▶ ....

# Les valeurs manquantes

## Pourquoi ?

- ▶ Oubli de l'enregistrement de la valeur
- ▶ Coût d'acquisition élevé
- ▶ Examen médical non effectué
- ▶ Refus de répondre à certaines questions

# Les valeurs manquantes

## Pourquoi ?

- ▶ Oubli de l'enregistrement de la valeur
- ▶ Coût d'acquisition élevé
- ▶ Examen médical non effectué
- ▶ Refus de répondre à certaines questions

## Conséquence

- ▶ Donnée manquante : donnée complexe
- ▶ Décision non représentative, voire dangereuse (médecine)

## Réponse (1)

Ignorer la valeur manquante : suppression des objets

## Réponse (1)

Ignorer la valeur manquante : suppression des objets

### Conséquences

- ▶ Réduction du volume des données
- ▶ Base non représentative (perte d'information)
- ▶ Information extraite non significative



## Réponse (1)

Ignorer la valeur manquante : suppression des objets

### Conséquences

- ▶ Réduction du volume des données
- ▶ Base non représentative (perte d'information)
- ▶ Information extraite non significative

## Réponse (2)

Imputation : valeur la plus commune, moyenne...

## Réponse (1)

Ignorer la valeur manquante : suppression des objets

### Conséquences

- ▶ Réduction du volume des données
- ▶ Base non représentative (perte d'information)
- ▶ Information extraite non significative

## Réponse (2)

Imputation : valeur la plus commune, moyenne...

### Conséquences

- ▶ Relations entre les attributs modifiées
- ▶ Probabilité élevée de se tromper

## Réponse (1)

Ignorer la valeur manquante : suppression des objets

### Conséquences

- ▶ Réduction du volume des données
- ▶ Base non représentative (perte d'information)
- ▶ Information extraite non significative

## Réponse (2)

Imputation : valeur la plus commune, moyenne...

### Conséquences

- ▶ Relations entre les attributs modifiées
- ▶ Probabilité élevée de se tromper

## Réponse (3)

Traitement de données manquantes internes à l'algorithme

Ex : C4.5, CART...

# Objectif

Aide à la décision en présence de données manquantes  
de *manière probabiliste*

# Objectif

Aide à la décision en présence de données manquantes  
de *manière probabiliste*

## Exemple

En médecine, pour un patient qui pourrait avoir trois maladies  $m_1, m_2, m_3$ , il serait préférable d'estimer la probabilité relatives de  $m_1, m_2, m_3$  plutôt que de lui affecter une seule maladie

# Objectif

Aide à la décision en présence de données manquantes  
de *manière probabiliste*

## Exemple

En médecine, pour un patient qui pourrait avoir trois maladies  $m_1, m_2, m_3$ , il serait préférable d'estimer la probabilité relatives de  $m_1, m_2, m_3$  plutôt que de lui affecter une seule maladie

## Arbre de décision

- ▶ Prédicatif et descriptif
- ▶ Facile à comprendre et simple à construire et à utiliser
- ▶ Proche des structures de connaissances manipulées naturellement par l'esprit humain

# Plan

Problématique et objectif

## Arbres de décision

Algorithme de construction d'un arbre de décision

Valeurs manquantes dans un arbre de décision

## Méthodes de traitement des valeurs manquantes

Arbres d'Attributs Ordonnés (AAO)

C4.5 : Quinlan

## Approches proposées

Objectif

Première proposition : AAOP

Deuxième proposition : AAP

Exemple de Classement

## Expérimentation

Bases de test

Taux d'erreurs

Comparaison des performances de C4.5 et AAP

## Analyse des résultats de classement

## Conclusion et perspectives

# Arbre de Décision

Apprentissage  
inductif, supervisé



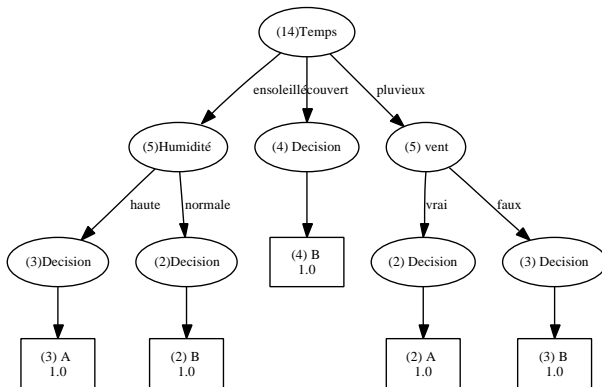
Extrait de la base météo : Ensemble d'apprentissage

id	Temps	Température	Humidité	Vent	Classe
1	Ensoleillé	Elevée	Haute	Faux	A
2	Ensoleillé	Elevée	Haute	Vrai	A
3	Couvert	Elevée	Haute	Faux	B
7	Couvert	Basse	Normale	Vrai	B
12	couvert	Moyenne	Haute	Vrai	B
13	Couvert	Elevée	Normale	Faux	B

L'attribut à prédire (classe) A : Je sors, B : Je reste à la maison



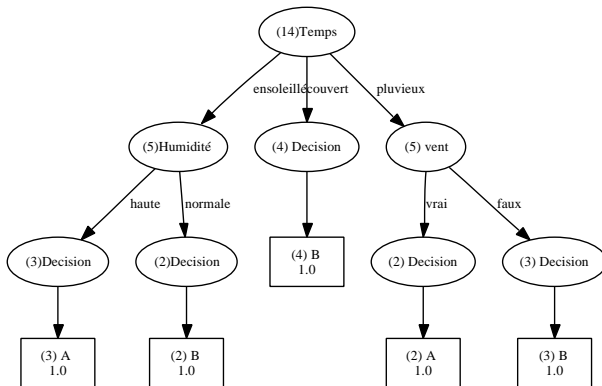
# Arbre de Décision : Phase de construction



# Arbre de Décision : Phase de classement

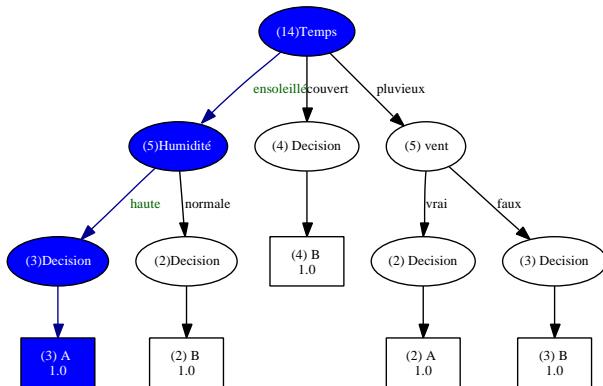
Détermination de la classe d'un nouvel objet **complet**

Temps	Température	Humidité	Vent
Ensoleillé	moyenne	haute	Faux



# Arbre de Décision : Phase de classement

Temps	Température	Humidité	Vent
Ensoleillé	moyenne	haute	Faux



- ▶ Parcours l'arbre depuis sa racine jusqu'à une feuille
- ▶ Classe associée à cette feuille : classe de l'objet : **A**

## Arbre de Décision : Phase de classement

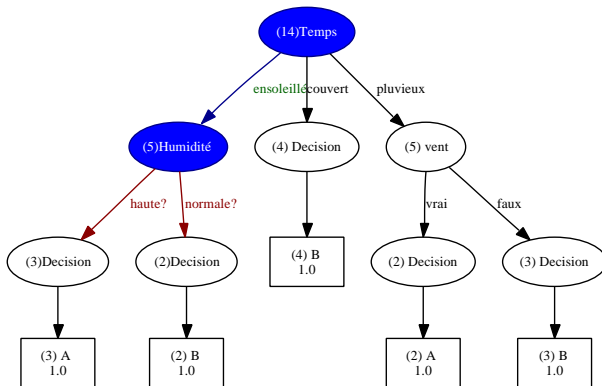
Détermination de la classe d'un nouvel objet **incomplet**

Temps	Température	Humidité	Vent
Ensoleillé	moyenne	?	Faux

# Arbre de Décision : Phase de classement

Détermination de la classe d'un nouvel objet **incomplet**

Temps	Température	Humidité	Vent
Ensoleillé	moyenne	?	Faux



# Plan

Problématique et objectif

## Arbres de décision

Algorithme de construction d'un arbre de décision

Valeurs manquantes dans un arbre de décision

Méthodes de traitement des valeurs manquantes

Arbres d'Attributs Ordonnés (AAO)

C4.5 : Quinlan

Approches proposées

Objectif

Première proposition : AAOP

Deuxième proposition : AAP

Exemple de Classement

Expérimentation

Bases de test

Taux d'erreurs

Comparaison des performances de C4.5 et AAP

Analyse des résultats de classement

Conclusion et perspectives

# Algorithme de construction d'un arbre de décision

- ▶ Partition récursive de l'ensemble d'apprentissage en sous-ensembles plus homogènes
- ▶ Choix des attributs à tester
  - ▶ Attributs permettant le mieux de classer les exemples
  - ▶ Mesure formelle : choix des attributs pertinents et élimination des attributs inutiles
- ▶ **Information Mutuelle [Shannon, 1949] :**
  - ▶ Mesure de la force de la relation entre deux attributs
  - ▶ Réduction de l'incertitude sur C lorsque A est connu

$$\text{gain}(A, C) = \text{IM}(A, C) = H(C) - H(C|A)$$

$$H(C) = \sum_{i=1}^{i=k} P_i \log_2 P_i$$

$$H(C|A) = \sum_{i=1}^{i=n} P(A = a_i) H(C|a_i)$$

# Algorithme de construction d'un arbre de décision

$$IM(\text{Classe}, \text{Temps}) = 0.246$$

$$IM(\text{Classe}, \text{Température}) = 0.030$$

$$IM(\text{Classe}, \text{Humidité}) = 0.154$$

$$IM(\text{Classe}, \text{Vent}) = 0.049$$

- ▶ Sélection de l'attribut pertinent qui maximise le gain d'information

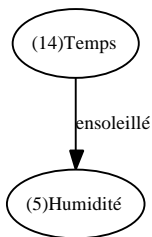
(14)Temps



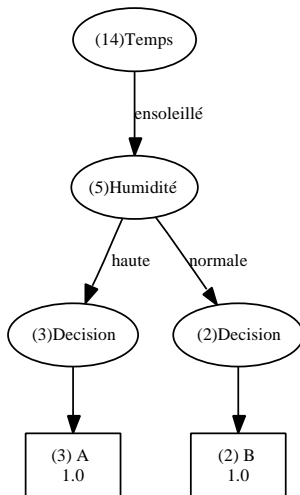
- ▶ Partition de l'ensemble d'apprentissage en 3 sous-ensembles selon les valeurs de l'attribut *Temps*

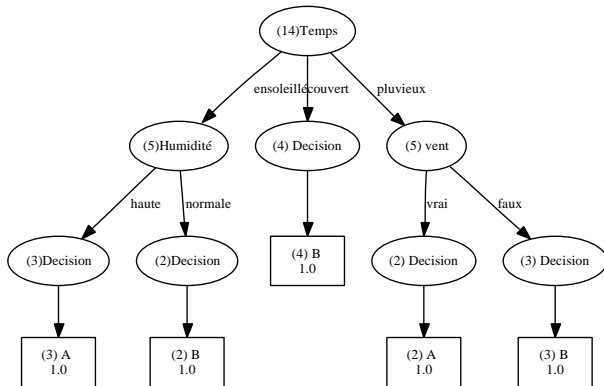
id	Temps	Température	Humidité	Vent	Classe
1	Ensoleillé	Elevée	Haute	Faux	A
2	Ensoleillé	Elevée	Haute	Vrai	A
8	Ensoleillé	Moyenne	Haute	Faux	A
9	Ensoleillé	Basse	Normale	Faux	B
11	Ensoleillé	Moyenne	Normale	Vrai	B

- ▶ Répétition du processus pour chaque sous-ensemble



- Création d'une feuille avec la valeur de la classe la plus probable si un critère d'arrêt est vérifié





# Plan

Problématique et objectif

## Arbres de décision

Algorithme de construction d'un arbre de décision

### Valeurs manquantes dans un arbre de décision

Méthodes de traitement des valeurs manquantes

Arbres d'Attributs Ordonnés (AAO)

C4.5 : Quinlan

Approches proposées

Objectif

Première proposition : AAOP

Deuxième proposition : AAP

Exemple de Classement

Expérimentation

Bases de test

Taux d'erreurs

Comparaison des performances de C4.5 et AAP

Analyse des résultats de classement

Conclusion et perspectives

# Valeurs manquantes

## Pendant la phase de **construction**

- ▶ Calcul du gain d'information pour choisir l'attribut test
- ▶ Partition de l'ensemble d'apprentissage selon l'attribut test choisi

# Valeurs manquantes

## Pendant la phase de **construction**

- ▶ Calcul du gain d'information pour choisir l'attribut test
- ▶ Partition de l'ensemble d'apprentissage selon l'attribut test choisi

## Pendant la phase de **classement**

Si la valeur d'un attribut test est manquante, il est impossible de décider quelle branche on doit choisir pour classer l'objet

# Plan

Problématique et objectif

Arbres de décision

Algorithme de construction d'un arbre de décision

Valeurs manquantes dans un arbre de décision

**Méthodes de traitement des valeurs manquantes**

**Arbres d'Attributs Ordonnés (AAO)**

C4.5 : Quinlan

Approches proposées

Objectif

Première proposition : AAOP

Deuxième proposition : AAP

Exemple de Classement

Expérimentation

Bases de test

Taux d'erreurs

Comparaison des performances de C4.5 et AAP

Analyse des résultats de classement

Conclusion et perspectives

# Arbres d'Attributs Ordonnés (AAO) : Lobo et Numao 2000

## Phase de construction : principe

- ▶ Ordonner les attributs par ordre croissant en fonction de leur Information Mutuelle avec la classe
- ▶ Construire un arbre de décision (appelé arbre d'attribut) pour chaque attribut



# Arbres d'Attributs Ordonnés (AAO) : Lobo et Numao 2000

## Phase de construction : principe

- ▶ Ordonner les attributs par ordre croissant en fonction de leur Information Mutuelle avec la classe
- ▶ Construire un arbre de décision (appelé arbre d'attribut) pour chaque attribut
  - ▶ Sous-ensemble : instances (valeurs connues pour cet attribut)
  - ▶ Attributs utilisés déjà traités

# Arbres des Attributs Ordonnés (AAO)

Extrait de la base météo

id	Temps	Température	Humidité	Vent	Classe
11	Ensoleillé	Elevée	Haute	Faux	A
2	Ensoleillé	Elevée	Haute	Vrai	A
3	Couvert	Elevée	Haute	Faux	B
4	Pluvieux	?	Haute	Faux	B
5	Pluvieux	Basse	Normale	Faux	B
6	Pluvieux	Basse	Normale	Vrai	A
7	Couvert	Basse	Normale	Vrai	B
8	Ensoleillé	Moyenne	Haute	Faux	A
9	Ensoleillé	Basse	Normale	Faux	B
10	Pluvieux	Moyenne	Normale	Faux	B
11	Ensoleillé	Moyenne	Normale	Vrai	B
...	...	...	...	...	...

## Arbres des Attributs Ordonnés (AAO)

$IM(\text{Classe}, \text{Température}) = 0.030$

$IM(\text{Classe}, \text{Vent}) = 0.049$

$IM(\text{Classe}, \text{Humidité}) = 0.154$

$IM(\text{Classe}, \text{Temps}) = 0.246$

*Température, Vent, Humidité, Temps*

# Arbres des Attributs Ordonnés (AAO)

$IM(\text{Classe}, \text{Température}) = 0.030$

$IM(\text{Classe}, \text{Vent}) = 0.049$

$IM(\text{Classe}, \text{Humidité}) = 0.154$

$IM(\text{Classe}, \text{Temps}) = 0.246$

*Température, Vent, Humidité, Temps*

id	Température
1	Elevée
2	Elevée
3	Elevée
5	Basse
6	Basse
8	Moyenne
10	Moyenne
11	Moyenne
12	Moyenne
..	..

# Arbres des Attributs Ordonnés (AAO)

$IM(\text{Classe}, \text{Température}) = 0.030$

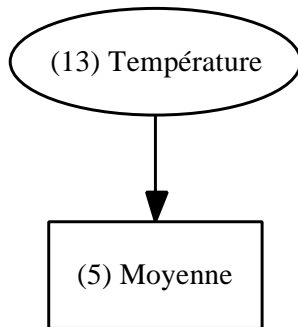
$IM(\text{Classe}, \text{Vent}) = 0.049$

$IM(\text{Classe}, \text{Humidité}) = 0.154$

$IM(\text{Classe}, \text{Temps}) = 0.246$

*Température, Vent, Humidité, Temps*

id	Température
1	Elevée
2	Elevée
3	Elevée
5	Basse
6	Basse
8	Moyenne
10	Moyenne
11	Moyenne
12	Moyenne
..	..



# Arbres d'Attributs Ordonnés (AAO)

*Température, Vent, Humidité, Temps*

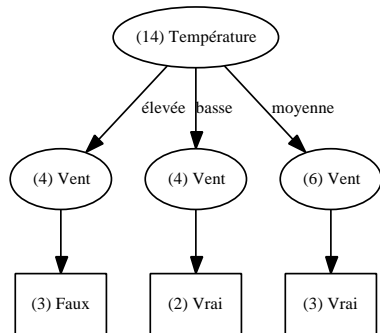
id	Température	Vent
1	Elevée	Faux
2	Elevée	Vrai
3	Elevée	Faux
4	<b>Moyenne</b>	<b>Faux</b>
5	Basse	Faux
6	Basse	Vrai
7	Basse	Vrai
..	..	..

# Arbres d'Attributs Ordonnés (AAO)

Température, *Vent*, Humidité, Temps

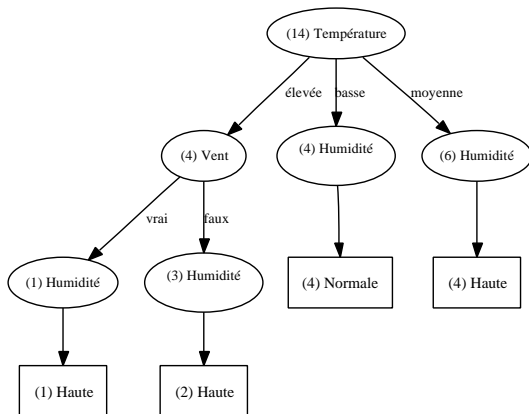
id	Température	Vent
1	Elevée	Faux
2	Elevée	Vrai
3	Elevée	Faux
4	<b>Moyenne</b>	<b>Faux</b>
5	Basse	Faux
6	Basse	Vrai
7	Basse	Vrai
..	..	..

L'AAO de Vent



# Arbres d'Attributs Ordonnés (AAO)

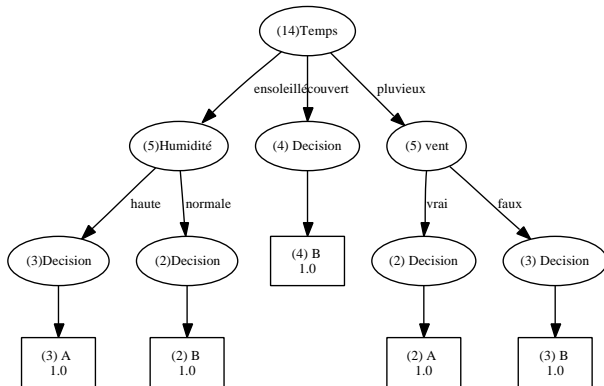
*Température, Vent, Humidité, Temps*  
Construit en utilisant *Vent* et *Température*





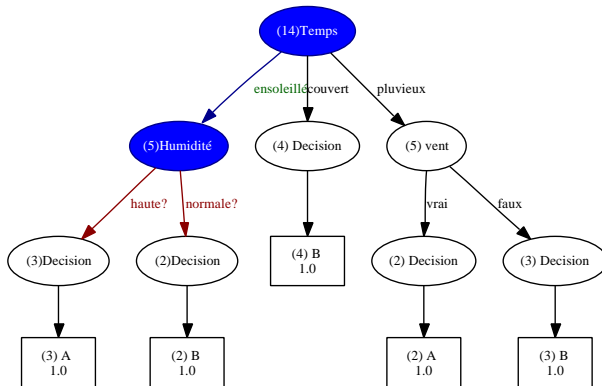
# AAO : Exemple de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	?	Faux	?



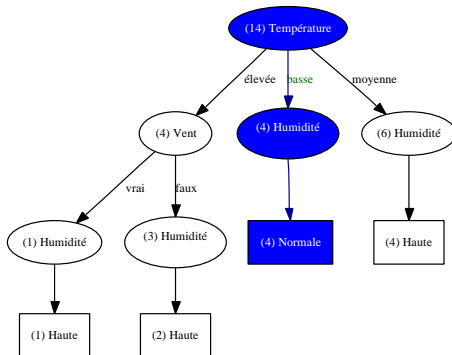
# AAO : Exemple de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	?	Faux	?



# AAO : Exemple de classement

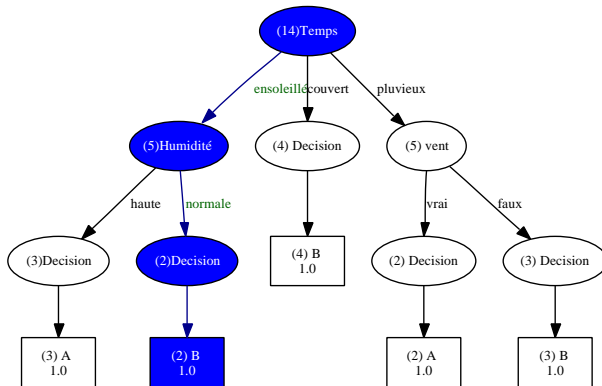
Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	?	Faux	?



Humidité est Normale

# AAO : Exemple de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	normale	Faux	?



L'objet appartient à la classe B

# AAO : Avantages et Inconvénients

## Avantages

# AAO : Avantages et Inconvénients

## Avantages

- ▶ Utilise les arbres de décision pour déterminer la valeur manquante d'un attribut

# AAO : Avantages et Inconvénients

## Avantages

- ▶ Utilise les arbres de décision pour déterminer la valeur manquante d'un attribut
- ▶ Commence par l'attribut le moins dépendant de la classe

# AAO : Avantages et Inconvénients

## Avantages

- ▶ Utilise les arbres de décision pour déterminer la valeur manquante d'un attribut
- ▶ Commence par l'attribut le moins dépendant de la classe

## Inconvénients



# AAO : Avantages et Inconvénients

## Avantages

- ▶ Utilise les arbres de décision pour déterminer la valeur manquante d'un attribut
- ▶ Commence par l'attribut le moins dépendant de la classe

## Inconvénients

- ▶ Attribut manquant remplacé par une seule valeur

# AAO : Avantages et Inconvénients

## Avantages

- ▶ Utilise les arbres de décision pour déterminer la valeur manquante d'un attribut
- ▶ Commence par l'attribut le moins dépendant de la classe

## Inconvénients

- ▶ Attribut manquant remplacé par une seule valeur
- ▶ Objet incomplet affecté à une seule classe

# AAO : Avantages et Inconvénients

## Avantages

- ▶ Utilise les arbres de décision pour déterminer la valeur manquante d'un attribut
- ▶ Commence par l'attribut le moins dépendant de la classe

## Inconvénients

- ▶ Attribut manquant remplacé par une seule valeur
- ▶ Objet incomplet affecté à une seule classe
- ▶ Ne prend pas en compte les dépendances entre les attributs

# Plan

Problématique et objectif

Arbres de décision

Algorithme de construction d'un arbre de décision

Valeurs manquantes dans un arbre de décision

## Méthodes de traitement des valeurs manquantes

Arbres d'Attributs Ordonnés (AAO)

### C4.5 : Quinlan

Approches proposées

Objectif

Première proposition : AAOP

Deuxième proposition : AAP

Exemple de Classement

Expérimentation

Bases de test

Taux d'erreurs

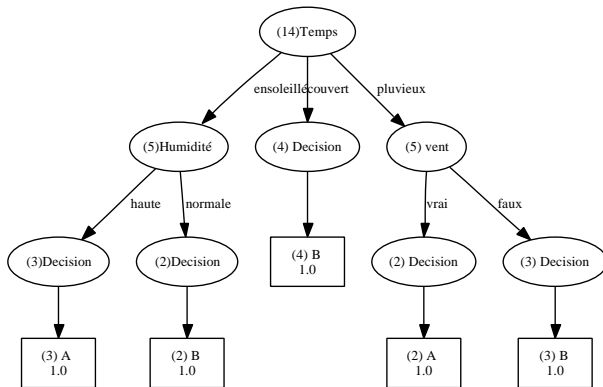
Comparaison des performances de C4.5 et AAP

Analyse des résultats de classement

Conclusion et perspectives

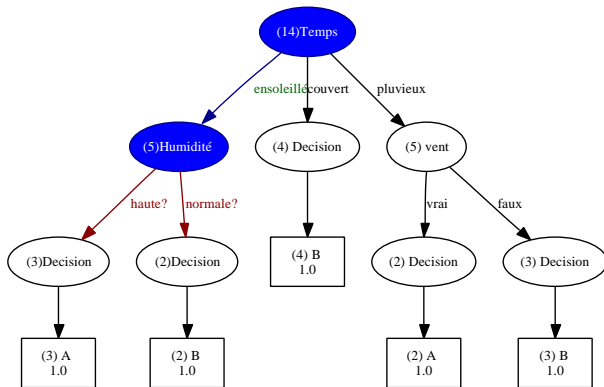
## C4.5 : Phase de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	?	Faux	?



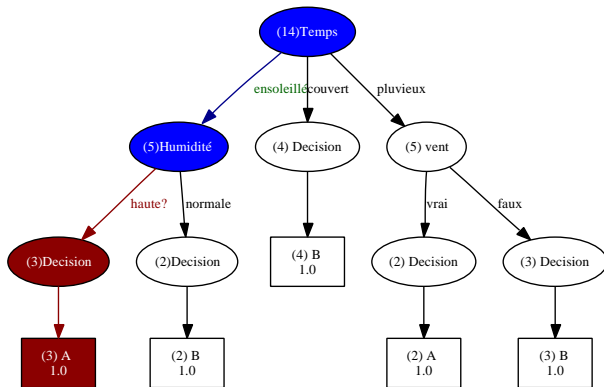
## C4.5 : Phase de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	?	Faux	?



## C4.5 : Phase de classement

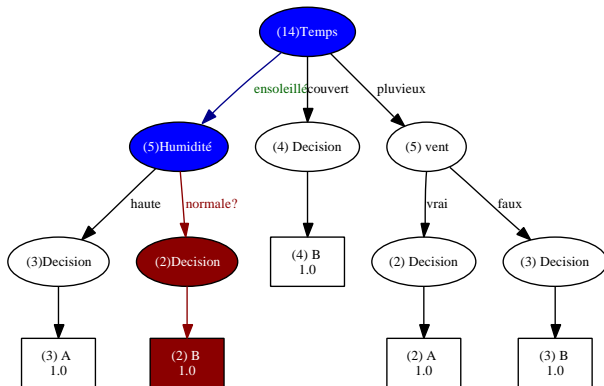
Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	?	Faux	?



$$P(A) = P(A|haute) \times P(haute) = 1 \times \left(\frac{3}{5}\right) = 0.6$$

## C4.5 : Phase de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	basse	?	Faux	?



$$P(B) = P(B|normale) \times P(normale) = 1 \times \left(\frac{2}{5}\right) = 0.4$$



## C4.5 : Avantages et Inconvénients

## C4.5 : Avantages et Inconvénients

### Avantages

- ▶ Simple à utiliser
- ▶ Distribution de probabilité de classe

## C4.5 : Avantages et Inconvénients

### Avantages

- ▶ Simple à utiliser
- ▶ Distribution de probabilité de classe

### Inconvénient

- ▶ Ne prend pas en compte les dépendances éventuelles entre les attributs

# Plan

- Problématique et objectif

- Arbres de décision

  - Algorithme de construction d'un arbre de décision

  - Valeurs manquantes dans un arbre de décision

- Méthodes de traitement des valeurs manquantes

  - Arbres d'Attributs Ordonnés (AAO)

  - C4.5 : Quinlan

- Approches proposées

  - Objectif

  - Première proposition : AAOP

  - Deuxième proposition : AAP

  - Exemple de Classement

- Expérimentation

  - Bases de test

  - Taux d'erreurs

  - Comparaison des performances de C4.5 et AAP

- Analyse des résultats de classement

- Conclusion et perspectives

## Approche proposée : Objectif

Classement probabiliste d'objets incomplètement connus dans un arbre de décision

## Approche proposée : Objectif

Classement probabiliste d'objets incomplètement connus dans un arbre de décision

- ▶ Valeur manquante d'un attribut est prédite sous forme d'une **distribution de probabilité**

# Approche proposée : Objectif

Classement probabiliste d'objets incomplètement connus dans un arbre de décision

- ▶ Valeur manquante d'un attribut est prédite sous forme d'une **distribution de probabilité**
- ▶ Utilisation des **dépendances** entre les attributs

# Approche proposée : Objectif

Classement probabiliste d'objets incomplètement connus dans un arbre de décision

- ▶ Valeur manquante d'un attribut est prédite sous forme d'une **distribution de probabilité**
- ▶ Utilisation des **dépendances** entre les attributs
- ▶ Résultat du classement est **probabiliste**



# Approche proposée : Objectif

Classement probabiliste d'objets incomplètement connus dans un arbre de décision

- ▶ Valeur manquante d'un attribut est prédite sous forme d'une **distribution de probabilité**
- ▶ Utilisation des **dépendances** entre les attributs
- ▶ Résultat du classement est **probabiliste**

## Arbres de Décision Probabiliste

Utilisation d'un arbre de décision probabiliste au lieu d'un arbre de décision classique en gardant sur chaque feuille la distribution de probabilités de classe au lieu de la classe la plus probable

# Plan

Problématique et objectif

Arbres de décision

Algorithme de construction d'un arbre de décision

Valeurs manquantes dans un arbre de décision

Méthodes de traitement des valeurs manquantes

Arbres d'Attributs Ordonnés (AAO)

C4.5 : Quinlan

## Approches proposées

Objectif

**Première proposition : AAOP**

Deuxième proposition : AAP

Exemple de Classement

Expérimentation

Bases de test

Taux d'erreurs

Comparaison des performances de C4.5 et AAP

Analyse des résultats de classement

Conclusion et perspectives

# AAOP : Arbres d'Attributs Ordonnés Probabilistes

## AAOP

- ▶ **Extension** de la méthode des AAOs (Lobo et Numao)
- ▶ *Attributs par ordre croissant selon l'Information Mutuelle relativement à la classe*

# AAOP : Arbres d'Attributs Ordonnés Probabilistes

## AAOP

- ▶ **Extension** de la méthode des AAOs (Lobo et Numao)
- ▶ *Attributs par ordre croissant selon l'Information Mutuelle relativement à la classe*
- ▶ Pour chaque attribut, construire un **AAOP** au lieu d'un AAO

# AAOP : Arbres d'Attributs Ordonnés Probabilistes

## AAOP

- ▶ **Extension** de la méthode des AAOs (Lobo et Numaou)
- ▶ *Attributs par ordre croissant selon l'Information Mutuelle relativement à la classe*
- ▶ Pour chaque attribut, construire un **AAOP** au lieu d'un AAO
- ▶ Utilisation des attributs *déjà traités* et **dépendants** de l'attribut courant

# AAOP : Arbres d'Attributs Ordonnés Probabilistes

## AAOP

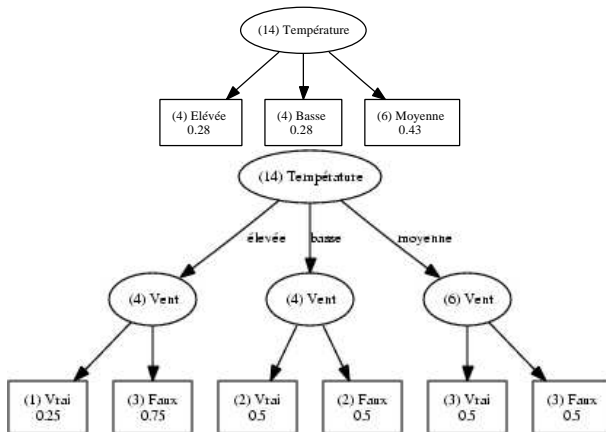
- ▶ **Extension** de la méthode des AAOs (Lobo et Numao)
- ▶ *Attributs par ordre croissant selon l'Information Mutuelle relativement à la classe*
- ▶ Pour chaque attribut, construire un **AAOP** au lieu d'un AAO
- ▶ Utilisation des attributs *déjà traités* et **dépendants** de l'attribut courant

## Différence par rapport à AAO (Lobo et Numao)

- + Feuilles probabilistes
- + Suppression des attributs déjà traités et **indépendants** de l'attribut courant

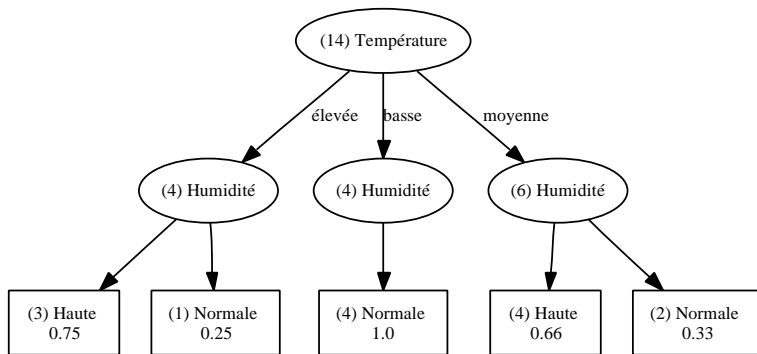
# Exemple : AAOP

## L'AAOP de Température et l'AAOP de Vent



# Exemple : AAOP

## L'AAOP d'Humidité



N'utilise pas l'attribut *Vent*



# Plan

Problématique et objectif

Arbres de décision

Algorithme de construction d'un arbre de décision

Valeurs manquantes dans un arbre de décision

Méthodes de traitement des valeurs manquantes

Arbres d'Attributs Ordonnés (AAO)

C4.5 : Quinlan

## Approches proposées

Objectif

Première proposition : AAOP

**Deuxième proposition : AAP**

Exemple de Classement

Expérimentation

Bases de test

Taux d'erreurs

Comparaison des performances de C4.5 et *AAP*

Analyse des résultats de classement

Conclusion et perspectives

## AAP : Arbres d'Attributs Probabilistes

La première proposition (AAOP) ne prend pas en compte toutes les dépendances

# AAP : Arbres d'Attributs Probabilistes

La première proposition (AAOP) ne prend pas en compte toutes les dépendances

## AAP

- ▶ Calcul pour chaque attribut de ses attributs dépendants en utilisant l'Information Mutuelle :

$$Dep(A_i) = \{A_j \mid IM_N(A_i, A_j) > Seuil \}$$

Seuil  $\geq$  l'Information Mutuelle Normalisée moyenne [Lobo et Numao 2000]

# AAP : Arbres d'Attributs Probabilistes

La première proposition (AAOP) ne prend pas en compte toutes les dépendances

## AAP

- ▶ Calcul pour chaque attribut de ses attributs dépendants en utilisant l'Information Mutuelle :

$$Dep(A_i) = \{A_j \mid IM_N(A_i, A_j) > Seuil \}$$

Seuil  $\geq$  l'Information Mutuelle Normalisée moyenne [Lobo et Numao 2000]

- ▶ Construction d'un AAP pour chaque attribut en utilisant les attributs dont il dépend

# AAP : Arbres d'Attributs Probabilistes

La première proposition (AAOP) ne prend pas en compte toutes les dépendances

## AAP

- ▶ Calcul pour chaque attribut de ses attributs dépendants en utilisant l'Information Mutuelle :

$$Dep(A_i) = \{A_j \mid IM_N(A_i, A_j) > Seuil \}$$

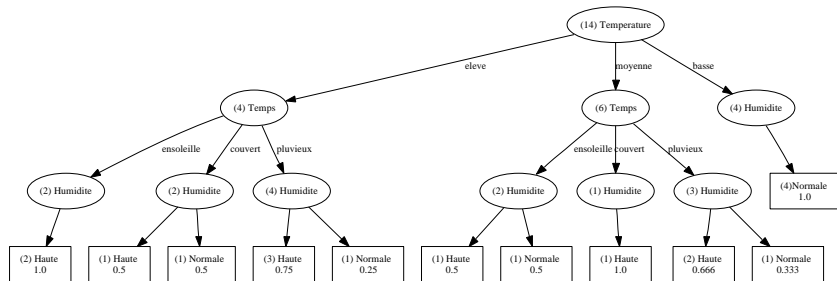
Seuil  $\geq$  l'Information Mutuelle Normalisée moyenne [Lobo et Numao 2000]

- ▶ Construction d'un AAP pour chaque attribut en utilisant les attributs dont il dépend

## Différence par rapport à AAOP

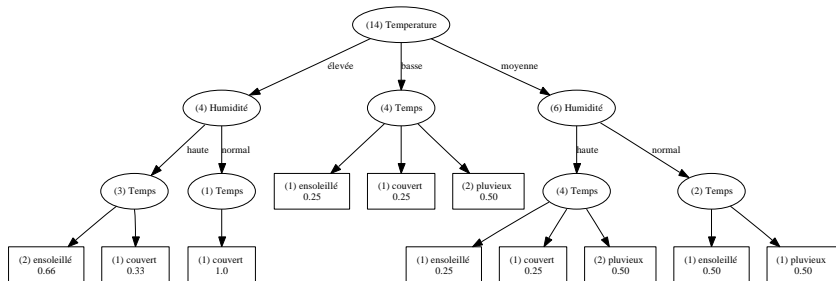
- + Pas d'ordre de construction
- + Arbre d'attribut construit en utilisant ses attributs dépendants

# Exemple : AAP



L'arbre d'Humidité construit avec Temps et Température

# Exemple : AAP



L'arbre de Temps construit avec **Température** et **Humidité**

# Problème avec AAP

## Cycle

Deux attributs **dépendants** et **manquants** en même temps



# Problème avec AAP

## Cycle

Deux attributs **dépendants** et **manquants** en même temps

## Solution

- ▶ Attribut le moins dépendant de la classe : appel de son arbre d'attribut construit selon *AAOP*
- ▶ Autre attribut : appel de son arbre d'attribut construit selon *AAP*

# Plan

Problématique et objectif

Arbres de décision

Algorithme de construction d'un arbre de décision

Valeurs manquantes dans un arbre de décision

Méthodes de traitement des valeurs manquantes

Arbres d'Attributs Ordonnés (AAO)

C4.5 : Quinlan

## Approches proposées

Objectif

Première proposition : AAOP

Deuxième proposition : AAP

### Exemple de Classement

Expérimentation

Bases de test

Taux d'erreurs

Comparaison des performances de C4.5 et AAP

Analyse des résultats de classement

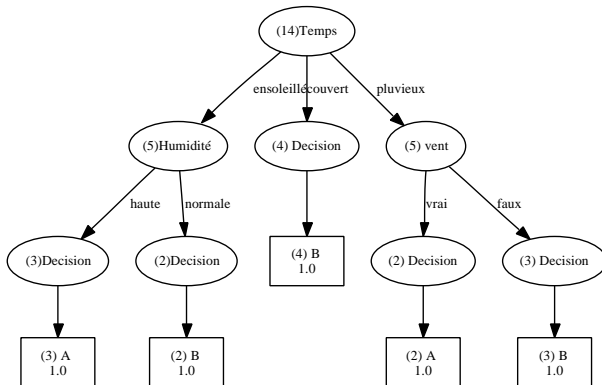
Conclusion et perspectives

## Exemple de classement

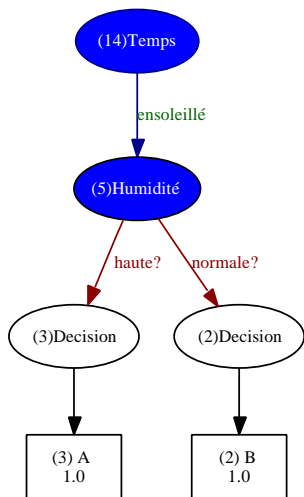
Temps	Température	Humidité	Vent	Classe
Ensoleillé	?	?	Faux	?

# Exemple de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	?	?	Faux	?



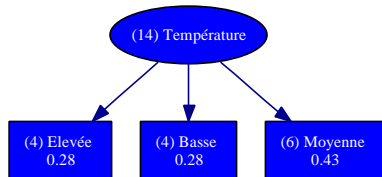
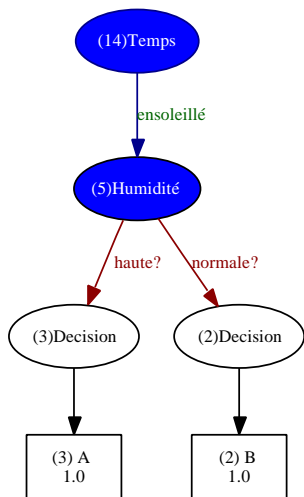
# Exemple de classement



$$P(A) = P(A|haute) P(haute)$$
$$P(B) = P(B|normale) P(normale)$$
$$P(A|haute) = 1 \implies$$
$$P(A) = P(haute)$$
$$P(B|normale)=1 \implies$$
$$P(B) = P(normale)$$

Humidité et Température :  
dépendants  
*AAP* d'Humidité  
*AAOP* de Température

# Exemple de classement



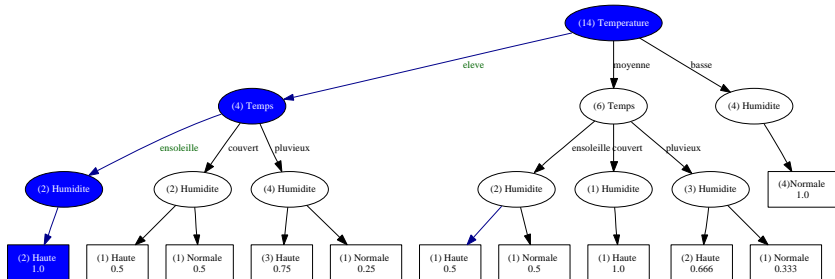
$$P(\text{élevée}) = 0.28$$

$$P(\text{moyenne}) = 0.43$$

$$P(\text{basse}) = 0.28$$

# Exemple de classement

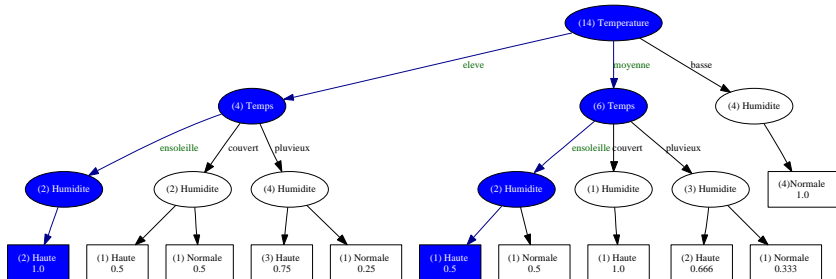
Temps	Température	Humidité	Vent	Classe
Ensoleillé	?	?	Faux	?



$$P(\text{haute}) = P(\text{haute} | \text{ensoleillé, élevée}) \times P(\text{ensoleillé, élevée})$$

# Exemple de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	?	?	Faux	?

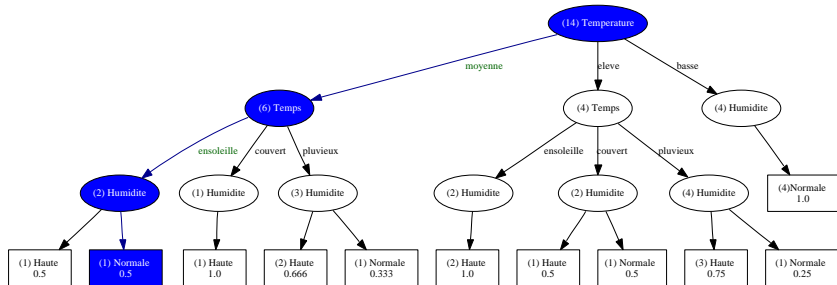


$$\begin{aligned}
 P(\text{haute}) &= P(\text{haute} \mid \text{ensoleillé, élevée}) \times P(\text{ensoleillé, élevée}) \\
 &\quad + P(\text{haute} \mid \text{ensoleillé, moyenne}) \times P(\text{ensoleillé, moyenne}) \\
 &= 1 \times 0.28 + 0.5 \times 0.43 = 0.495
 \end{aligned}$$



# Exemple de classement

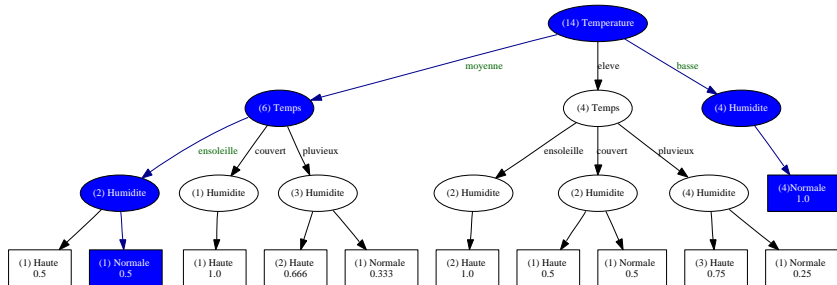
Temps	Température	Humidité	Vent	Classe
Ensoleillé	?	?	Faux	?



$$P(\text{normale}) = P(\text{normale} | \text{ensoleillé, moyenne}) \times P(\text{ensoleillé, moyenne})$$

# Exemple de classement

Temps	Température	Humidité	Vent	Classe
Ensoleillé	?	?	Faux	?



$$\begin{aligned} P(\text{normale}) &= P(\text{normale}|\text{ensoleillé, moyenne}) \times P(\text{ensoleillé, moyenne}) \\ &\quad + P(\text{normale}|\text{basse}) \times P(\text{basse}) \\ &= 0.5 \times 0.43 + 1 \times 0.28 = 0.495 \end{aligned}$$

# Avantages et Inconvénients

## Avantages

- ▶ Par rapport à C4.5 : prendre en compte d'autres attributs qui ne sont pas forcément dans l'arbre  
*Température* lors du calcul de la probabilité de *Humidité*

# Avantages et Inconvénients

## Avantages

- ▶ Par rapport à C4.5 : prendre en compte d'autres attributs qui ne sont pas forcément dans l'arbre  
*Température* lors du calcul de la probabilité de *Humidité*
- ▶ Par rapport à AAO : utilise un arbre d'attribut construit en fonction de ses attributs dépendants  
*Temps* lors du calcul de la probabilité de *Humidité*

# Avantages et Inconvénients

## Avantages

- ▶ Par rapport à C4.5 : prendre en compte d'autres attributs qui ne sont pas forcément dans l'arbre  
*Température* lors du calcul de la probabilité de *Humidité*
- ▶ Par rapport à AAO : utilise un arbre d'attribut construit en fonction de ses attributs dépendants  
*Temps* lors du calcul de la probabilité de *Humidité*

## Inconvénient

Complexité de classement **exponentielle** en fonction du nombre d'attributs manquants dans l'objet à classer

$$\text{Complexity}(\text{Trees}) = O(m \times \bar{v}^m \log(\bar{L}_T))$$

# Plan

- Problématique et objectif

- Arbres de décision

  - Algorithme de construction d'un arbre de décision

  - Valeurs manquantes dans un arbre de décision

- Méthodes de traitement des valeurs manquantes

  - Arbres d'Attributs Ordonnés (AAO)

  - C4.5 : Quinlan

- Approches proposées

  - Objectif

  - Première proposition : AAOP

  - Deuxième proposition : AAP

  - Exemple de Classement

- Expérimentation**

  - Bases de test**

  - Taux d'erreurs

  - Comparaison des performances de C4.5 et AAP

- Analyse des résultats de classement

- Conclusion et perspectives

# Bases de test

Test de notre approche, C4.5 et AAO sur 11 bases

- ▶ *Attributs dépendants* : Vote, Breast-cancer, Zoo, Lymphography, Mushroom, Dermatology, Splice, Iris, Breast-w
- ▶ *Attributs indépendants* : Nursery, Car

# Bases de test

Test de notre approche, C4.5 et AAO sur 11 bases

- ▶ *Attributs dépendants* : Vote, Breast-cancer, Zoo, Lymphography, Mushroom, Dermatology, Splice, Iris, Breast-w
- ▶ *Attributs indépendants* : Nursery, Car
  
- ▶ Base de test construite à partir de la base d'apprentissage
- ▶ Attributs pertinents inconnus
- ▶ Tests sur plusieurs seuils de dépendance
- ▶ Comparaison avec les résultats de classement donnés par C4.5 et AAO pour la même base
- ▶ Calcul de la matrice de confusion de chaque méthode



# Bases de test

## La base *mushroom*

- ▶ Champignons sont comestibles ou toxiques
- ▶ Taille base d'apprentissage : 5644 instances, 22 attributs discrets
- ▶ Taille base de test : 73 instances
- ▶ La racine : *odor*, 75.34% de valeurs manquantes



Mushroom	Seuil	bien classés	mal classés	50%
<b>AAP</b>	<b>0.1</b> 0.2	<b>80.82%</b> 75.34%	<b>19.17 %</b> 24.65 %	
<b>C4.5</b>		<b>58.90%</b>	<b>41.09%</b>	
<b>AAO</b>		67.12%	32.87%	

# Plan

- Problématique et objectif

- Arbres de décision

  - Algorithme de construction d'un arbre de décision

  - Valeurs manquantes dans un arbre de décision

- Méthodes de traitement des valeurs manquantes

  - Arbres d'Attributs Ordonnés (AAO)

  - C4.5 : Quinlan

- Approches proposées

  - Objectif

  - Première proposition : AAOP

  - Deuxième proposition : AAP

  - Exemple de Classement

- Expérimentation**

  - Bases de test

  - Taux d'erreurs**

  - Comparaison des performances de C4.5 et AAP

- Analyse des résultats de classement

- Conclusion et perspectives

# Taux d'erreurs

## Root Mean Squared Error

Pour une instance  $x$  le *RMES* est donné par l'équation suivante :

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^{j=n} (t(j|x) - P(j|x))^2}$$

$x$  est l'instance,  $j$  est la valeur de classe

$t(j|x)$  est la vraie probabilité de la classe  $j$  pour  $x$

$P(j|x)$  est la probabilité estimée par la méthode pour l'instance  $x$  et la classe  $j$

# Taux d'erreurs

DataBase	AAP	C4.5	AAO
Zoo	<b>0.133817</b>	0.44812	0.245
Mushroom	<b>0.412543</b>	0.643147	0.535865
Dermatology	<b>0.25183339</b>	0.4108025	0.3068860
Vote	<b>0.310443</b>	0.52039	0.315079
Breast-cancer	<b>0.525338</b>	0.632477	0.6238318
Lymphography	<b>0.260477</b>	0.477835	0.420603
Splice	0.35292	0.38507	<b>0.22929</b>
Iris	0.290535	0.534436	<b>0.226616</b>
Breast-w	<b>0.29610</b>	0.55673	0.32818
Nursery	0.4456728	0.44999	0.436149
Car	0.3292286	0.32743	0.35666

# Plan

- Problématique et objectif

- Arbres de décision

  - Algorithme de construction d'un arbre de décision

  - Valeurs manquantes dans un arbre de décision

- Méthodes de traitement des valeurs manquantes

  - Arbres d'Attributs Ordonnés (AAO)

  - C4.5 : Quinlan

- Approches proposées

  - Objectif

  - Première proposition : AAOP

  - Deuxième proposition : AAP

  - Exemple de Classement

- Expérimentation**

  - Bases de test

  - Taux d'erreurs

  - Comparaison des performances de C4.5 et AAP**

- Analyse des résultats de classement

- Conclusion et perspectives

# Comparaison de performances entre C4.5 et AAP

## Test McNemar

Soient deux hypothèses de classement  $h_1$  (AAP) et  $h_2$  (C4.5) :

$$M = \frac{(|n_{01} - n_{10}| - 1)^2}{n_{10} + n_{01}}$$

Ce test suit une loi  $\chi^2$  à 1 degré de liberté. L'hypothèse que  $h_1$  et  $h_2$  aient le même taux d'erreur peut être rejetée avec une probabilité supérieure à 95% si  $|M| > 3.84$

où

- ▶  $n_{11}$  : le nombre d'instances bien classées par  $h_1$  et  $h_2$  ;
- ▶  $n_{01}$  : le nombre d'instances mal classées par  $h_1$  et non par  $h_2$  ;
- ▶  $n_{10}$  : le nombre d'instances mal classées par  $h_2$  et non par  $h_1$  ;
- ▶  $n_{00}$  : le nombre d'instances mal classées par  $h_2$  et  $h_1$ .

## Les résultats du test McNemar sur les bases de données

$n_{01}$  : nombre d'instances mal classées par **AAP** et non par **C4.5**

$n_{10}$  : nombre d'instances mal classées par **C4.5** et non par **AAP**

Base	$n_{01}$	$n_{10}$	McNemar	Différence
Vote	12	32	8.205	significative
Breast-cancer	5	14	3,368	non significative
Lymphography	2	11	4.923	significative
Zoo	1	23	18.375	très significative
Mushroom	0	12	10.083	très significative
Dermatology	1	8	4.000	significative
Splice	2	5	0.571	non significative
Iris	1	19	14.450	très significative
Breast-w	3	19	10.227	très significative
Nursery	2	0	3.448	non significative
Car	1	5	1.500	non significative

# Plan

- Problématique et objectif

- Arbres de décision

  - Algorithme de construction d'un arbre de décision

  - Valeurs manquantes dans un arbre de décision

- Méthodes de traitement des valeurs manquantes

  - Arbres d'Attributs Ordonnés (AAO)

  - C4.5 : Quinlan

- Approches proposées

  - Objectif

  - Première proposition : AAOP

  - Deuxième proposition : AAP

  - Exemple de Classement

- Expérimentation

  - Bases de test

  - Taux d'erreurs

  - Comparaison des performances de C4.5 et AAP

- Analyse des résultats de classement**

- Conclusion et perspectives



# Analyse des résultats de classement

## Objectif

- ▶ Analyse du résultat de classement de chaque objet test
- ▶ Comparaison des distributions de probabilités données par *AAP* et *C4.5*
- ▶ Choix de la distribution la plus représentative par rapport à la base d'apprentissage

## Algorithme Analyser-Instance

Pour un objet test  $o$

- ▶ Calcul de la distance entre l'objet  $o$  et tous les objets de la base d'apprentissage
- ▶ Si ( $\text{Distance} < \text{Near}$ ), deux objets Plus Proches Voisins (PPV)
- ▶ Fréquence de ses Plus Proches Voisins de chaque classe

# RELIEF

## Relief : Kira et Rendell 1992

- ▶ Mesure d'impureté pour évaluer la qualité d'un attribut
- ▶ Basé entièrement sur l'analyse statistique
- ▶ Dépendances entre les attributs
- ▶ Valeurs manquantes
- ▶ Calcul de la distance entre les objets

# RELIEF

## Fonction de Distance

$$Distance(l_1, l_2) = \sum_{j=1}^{j=n} diff(A_j, l_1, l_2)$$

$$diff(A, l_1, l_2) = \begin{cases} 0 & \text{if } V^{(A, l_1)} = V^{(A, l_2)} \\ 1 & \text{if } V^{(A, l_1)} \neq V^{(A, l_2)} \\ 1 - P(V^{(A, l_2)} / Class_{l_1}) & \text{if } A \text{ is unknown in } l_1 \end{cases}$$

- ▶  $Class_{l_1}$  est la valeur de la classe dans l'instance  $l_1$
- ▶  $1 - P(V^{(A, l_2)} / Class_{l_1})$  est la probabilité que deux instances  $l_1$  et  $l_2$  prennent des valeurs différentes pour l'attribut  $A$  dans le cas où une de ces instances ( $l_1$  ici) possède une valeur inconnue pour  $A$

# Exemple de calcul de distance

## Instance test

id	Temps	Température	Humidité	Vent	Classe
5	Ensoleillé	Moyenne	?	Faux	A

## Base d'apprentissage

id	Temps	Température	Humidité	Vent	Classe
1	Ensoleillé	Moyenne	Haute	Faux	A
2	Pluvieux	Moyenne	Normale	Faux	B
3	Pluvieux	Moyenne	Haute	Vrai	A
4	Ensoleillé	Moyenne	Normale	faux	A

# Exemple de calcul de distance

## Instance test

id	Temps	Température	Humidité	Vent	Classe
5	Ensoleillé	Moyenne	?	Faux	A

## Base d'apprentissage

id	Temps	Température	Humidité	Vent	Classe
1	Ensoleillé	Moyenne	Haute	Faux	A
2	Pluvieux	Moyenne	Normale	Faux	B
3	Pluvieux	Moyenne	Haute	Vrai	A
4	Ensoleillé	Moyenne	Normale	faux	A

$$\begin{aligned} \text{Distance}(5, 1) &= 0 + 0 + (1 - P(\text{Haute}|A)) + 0 \\ &= (1 - \frac{2}{3}) = 0.333 \end{aligned}$$

# Exemple de calcul de distance

## Instance test

id	Temps	Température	Humidité	Vent	Classe
5	Ensoleillé	Moyenne	?	Faux	A

## Base d'apprentissage

id	Temps	Température	Humidité	Vent	Classe
1	Ensoleillé	Moyenne	Haute	Faux	A
2	Pluvieux	Moyenne	Normale	Faux	B
3	Pluvieux	Moyenne	Haute	Vrai	A
4	Ensoleillé	Moyenne	Normale	faux	A

Distance(5, 1)=0.33

Si Near= 3 : Distance(5, 1) < 3, les deux objets sont PPV de même classe

# Exemple de calcul de distance

## Instance test

id	Temps	Température	Humidité	Vent	Classe
5	Ensoleillé	Moyenne	?	Faux	A

## Base d'apprentissage

id	Temps	Température	Humidité	Vent	Classe
1	Ensoleillé	Moyenne	Haute	Faux	A
2	Pluvieux	Moyenne	Normale	Faux	B
3	Pluvieux	Moyenne	Haute	Vrai	A
4	Ensoleillé	Moyenne	Normale	faux	A

$$\begin{aligned} \text{Distance}(5, 2) &= 1 + 0 + (1 - P(\text{normale}|A)) + 0 \\ &= 1 + (1 - \frac{1}{3}) = 1.66666 \end{aligned}$$

# Exemple de calcul de distance

## Instance test

id	Temps	Température	Humidité	Vent	Classe
5	Ensoleillé	Moyenne	?	Faux	A

## Base d'apprentissage

id	Temps	Température	Humidité	Vent	Classe
1	Ensoleillé	Moyenne	Haute	Faux	A
2	Pluvieux	Moyenne	Normale	Faux	B
3	Pluvieux	Moyenne	Haute	Vrai	A
4	Ensoleillé	Moyenne	Normale	faux	A

Distance(5, 2)=1.666

Si Near= 3 : Distance(5, 2) < 3 et les deux objets sont PPV  
des classes différentes



# Exemple : Analyser-Instance

Zoo



- ▶ Base d'apprentissage : 101 objets et 17 attributs
- ▶ Classe : 7 valeurs
- ▶ Base de test : 71 objets

Attributs	Instance1
hair	true
feathers	false
eggs	true
milk	?
airborne	false
aquatic	true
predator	?
toothed	false
backbone	true
breathes	true
venomous	false
fins	false
legs	4
tail	?
domestic	false
catsize	true
classe	mammal

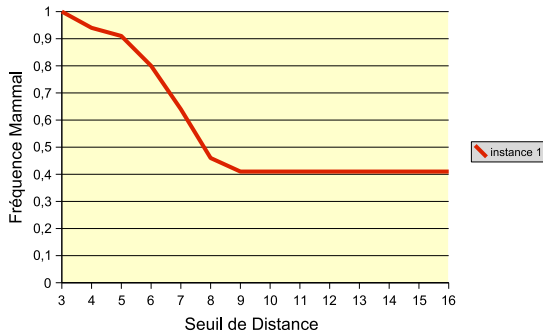
## Exemple : Analyser-Instance

Valeurs de classe :

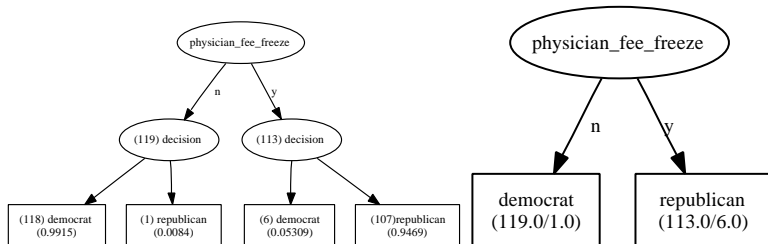
mammal, bird, reptile, fish, amphibian, insect, invertebrate

Near	Résultat	Instances proches
3	( <b>100%</b> , 0%, 0%, 0%, 0%, 0%, 0%)	2
4	( <b>94%</b> , 0%, 5%, 0%, 0%, 0%, 0%)	19
5	( <b>91%</b> , <b>3%</b> , <b>3%</b> , 0%, <b>3%</b> , 0%, 0%)	34
7	( <b>64%</b> , 17%, 5%, 5%, 6%, 0%, 3%)	64
9	(41%, 20%, 5%, 13%, 4%, 7%, 10%)	100
AAP	( <b>100%</b> , 0%, 0%, 0%, 0%, 0%, 0%)	
AAO	(0%, 0%, <b>100%</b> , 0%, 0%, 0%, 0%)	
C4.5	( <b>51%</b> , 0%, <b>27%</b> , 0%, <b>22%</b> , 0%, 0%)	

## Fréquence de la classe *mammifère*



# La base Vote



L'AAP de la base Vote pour un seuil 0.5 et l'arbre de C4.5

# La base Vote

## Vote

- ▶ Base d'apprentissage : 232 objets, 16 attributs
- ▶ Base de test : 240 objets
- ▶ Classe : Democrat et Republican



attributes	instance 1	instance 2	instance 3
physician-fee	?	?	?
el-salvador	y	?	y
education	y	n	n
crime	y	n	n
near=8	<b>(16%, 83%)</b>	<b>(91%, 08%)</b>	<b>(92%, 07%)</b>
near=10	(29%, 70%)	(84%, 15%)	(75%, 24%)
near=12	(38%, 61%)	(70%, 29%)	(57%, 42%)
AAP	(11%,89%)	(99%,01%)	(85%,15%)
AAO	(0%,100%)	(100%,0%)	(100%,0%)
C4.5	(53%,47%)	(53%,47%)	(53%,47%)

# Plan

- Problématique et objectif

- Arbres de décision

  - Algorithme de construction d'un arbre de décision

  - Valeurs manquantes dans un arbre de décision

- Méthodes de traitement des valeurs manquantes

  - Arbres d'Attributs Ordonnés (AAO)

  - C4.5 : Quinlan

- Approches proposées

  - Objectif

  - Première proposition : AAOP

  - Deuxième proposition : AAP

  - Exemple de Classement

- Expérimentation

  - Bases de test

  - Taux d'erreurs

  - Comparaison des performances de C4.5 et AAP

- Analyse des résultats de classement

- Conclusion et perspectives**

# Conclusion 1

- ▶ Traitement du problème des valeurs manquantes
- ▶ Deux propositions :
  - ▶ Arbres d'Attributs Ordonnés Probabilistes (AAOP)
  - ▶ Arbres d'Attributs Probabilistes (AAP) : *Dépendances*
- ▶ Tests sur plusieurs bases de données
- ▶ Comparaison avec les résultats de classement donnés par C4.5 et AAO

# Conclusion 1

- ▶ Traitement du problème des valeurs manquantes
- ▶ Deux propositions :
  - ▶ Arbres d'Attributs Ordonnés Probabilistes (AAOP)
  - ▶ Arbres d'Attributs Probabilistes (AAP) : *Dépendances*
- ▶ Tests sur plusieurs bases de données
- ▶ Comparaison avec les résultats de classement donnés par C4.5 et AAO
  
- ▶ Résultat de classement : une distribution de probabilité
- ▶ Performance meilleure quand les attributs sont dépendants
- ▶ Complexité *exponentielle* en fonction du nombre d'attributs manquants
  - ▶ Délimite l'usage : peu de valeurs manquantes  $\Rightarrow$  amélioration par rapport à C4.5 et AAO est importante



## Conclusion 2

- ▶ Proposition d'un algorithme, *Analyser-Instance*, issu de la méthode des k plus proches voisins
- ▶ Calcul pour chaque objet à classer de la fréquence de ses objets les plus proches de chaque classe
- ▶ Utilisation d'une fonction de distance qui prend en compte les valeurs manquantes dans l'objet à classer

## Conclusion 2

- ▶ Proposition d'un algorithme, *Analyser-Instance*, issu de la méthode des k plus proches voisins
  - ▶ Calcul pour chaque objet à classer de la fréquence de ses objets les plus proches de chaque classe
  - ▶ Utilisation d'une fonction de distance qui prend en compte les valeurs manquantes dans l'objet à classer
- 
- ▶ Résultats de l'algorithme *Analyser-Instance* proches des résultats donnés par notre approche
  - ▶ Mauvaise performance de C4.5 lorsque l'attribut manquant est la racine de l'arbre

# Perspectives

- ▶ Comparaison des distribution de probabilités avec celles obtenues avec les Réseaux Bayésiens et les Règles d'Association
- ▶ Arbres de décision probabilistes

# Perspectives

- ▶ Comparaison des distribution de probabilités avec celles obtenues avec les Réseaux Bayésiens et les Règles d'Association
- ▶ Arbres de décision probabilistes

Simplification de notre approche pour diminuer la complexité de classement

- ▶ Construction d'une seule famille d'arbres *AAP* qui prend en compte les dépendances
- ▶ Pour résoudre le problème de Cycle : une solution hybride qui combine :
  - ▶ Arbres d'Attributs Probabilistes (AAP)
  - ▶ Analyser-Instance (utiliser pour l'attribut le moins dépendant de la classe)

*Merci pour votre attention*