

Le contexte



Data Mining

- En français : Fouille de données
- Découverte de connaissances dans de gros volumes de données
- Exemples (non réels) :
 - Dans les données d'un site de vente en ligne : Une personne qui achète le livre A achète souvent le livre B (dans 40% des cas)
 - Si l'on s'intéresse à l'achat de voitures : Si la personne est jeune et habite en ville alors la voiture est de type sport

Les domaines d'application...

Potentiellement tous les domaines qui stockent un nombre important de données

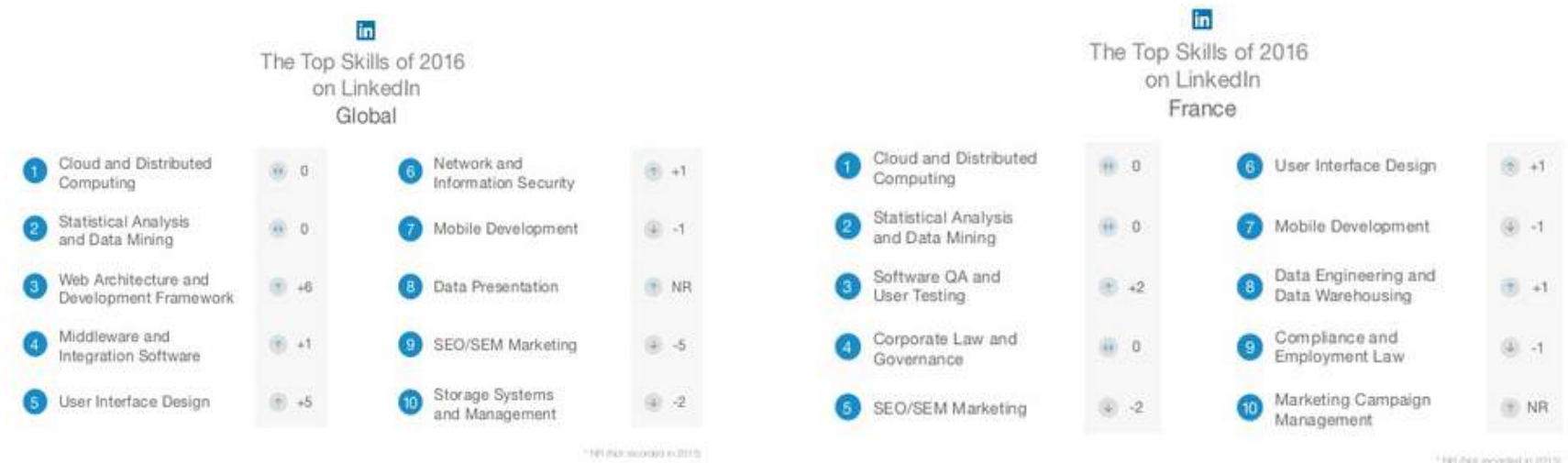
- Banques (prêts), assurances
- Commerce en ligne (profils, préférences des clients)
- Images satellites, médicales (reconnaissance d'objets)
- Environnement

Ce ne sont que quelques exemples !

Une compétence recherchée

Les 10 compétences les plus recherchées en 2016 selon LinkedIn (basé sur l'analyse en termes de compétences et d'historique de recrutement de centaines de millions de profils)

<https://www.linkedin.com/slideshare/linkedin/the-top-skills-that-can-get-you-hired-in-2017>



« Statistical analysis and data mining » est classée dans les 5 premières dans tous les pays analysés

Une compétence recherchée

Les 25 compétences les plus recherchées en 2015 selon LinkedIn
<https://blog.linkedin.com/2016/01/12/the-25-skills-that-can-get-you-hired-in-2016>



« Statistical analysis and data mining » est la seule catégorie classée dans les 4 premières dans tous les pays analysés

Une compétence recherchée

Les 25 compétences les plus recherchées en 2014 selon LinkedIn

<http://blog.linkedin.com/2014/12/17/the-25-hottest-skills-that-got-people-hired-in-2014/>

The 25 Hottest Skills of 2014 on LinkedIn

Global

- 1 Statistical Analysis and Data Mining
- 2 Middleware and Integration Software
- 3 Storage Systems and Management
- 4 Network and Information Security
- 5 SEO/SEM Marketing
- 6 Business Intelligence
- 7 Mobile Development
- 8 Web Architecture and Development Framework
- 9 Algorithm Design
- 10 Perl/Python/Ruby

The 25 Hottest Skills of 2014 on LinkedIn

France

- 1 Statistical Analysis and Data Mining
- 2 Cloud and Distributed Computing
- 3 C/C++
- 4 SEO/SEM Marketing
- 5 IBM Mainframe and Systems
- 6 Dispute Resolution law
- 7 Business Intelligence
- 8 Mechanical and Aerospace Engineering
- 9 Corporate Law and Governance
- 10 Computer Graphics and Animation

Une compétence recherchée

Et en 2013...

<http://blog.linkedin.com/2013/12/18/the-25-hottest-skills-that-got-people-hired-in-2013/>

The 25 Hottest Skills of 2013 on LinkedIn

- ① Social Media Marketing
- ② Mobile Development
- ③ Cloud and Distributed Computing
- ④ Perl/Python/Ruby
- ⑤ Statistical Analysis and Data Mining
- ⑥ User Interface Design
- ⑦ Digital and Online Marketing
- ⑧ Recruiting
- ⑨ Business Development/Relationship Management
- ⑩ Retail Payment and Information Systems

Machine Learning

- En français : Apprentissage automatique, apprentissage artificiel
- Une partie de la fouille de données
- Concerne « Toute méthode permettant de construire un modèle de la réalité à partir de données » (Apprentissage artificiel – Concepts et algorithmes. Cornuéjols et Miclet. 2010)

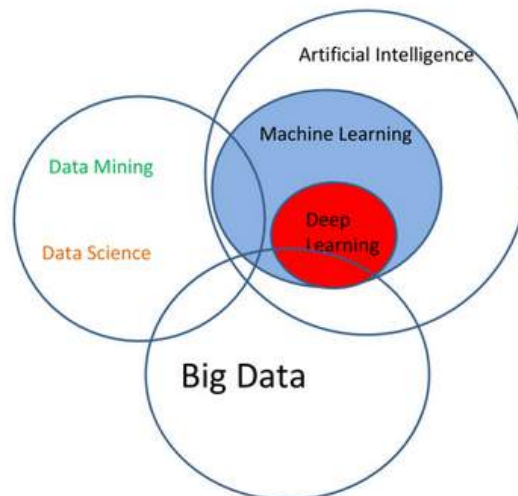
Data Mining/Machine Learning/Data Science

- Des domaines liés, des frontières aux contours mal définis

<http://www.kdnuggets.com/2016/10/battle-data-science-venn-diagrams.html>

Une proposition de représentation

<http://www.kdnuggets.com/2016/03/data-science-puzzle-explained.html>



Des compétences recherchées

- Juin 2015 : Twitter annonce le rachat de Whetlab, startup fondée par des chercheurs et développant des technologies d'apprentissage (analyses prédictives sur le comportement des utilisateurs)

<http://www.zdnet.fr/actualites/machine-learning-twitter-acquiert-whetlab-39821008.htm>

<http://www.journaldunet.com/solutions/cloud-computing/twitter-acquiert-whetlab-0615.shtml>

- Juin 2016 : Twitter acquiert Magic Pony Technology qui a développé une solution de traitement d'image par apprentissage automatique

<http://www.zdnet.fr/actualites/intelligence-artificielle-twitter-met-la-main-sur-magic-pony-technology-39838622.htm>

Des compétences recherchées

- Juin 2015 : Yahoo! Inc. et l'Université Carnegie Mellon annoncent un partenariat de 10 millions de \$ pour que les chercheurs puissent tester de nouvelles façons d'améliorer l'expérience utilisateur avec l'apprentissage automatique notamment

<https://www.ml.cmu.edu/news/yahoo-announcement.html>

Des compétences recherchées

- Plus de 600 publications des chercheurs de Google sur l'intelligence artificielle et l'apprentissage automatique

<http://research.google.com/pubs/MachineIntelligence.html>

- Juin 2016 : Google crée un groupe de recherche en Europe, dédié au machine learning

<http://www.zdnet.fr/actualites/intelligence-artificielle-google-cree-un-groupe-de-recherche-en-europe-39838488.htm>

Des compétences recherchées

- Septembre 2015 : Apple recrute des experts en apprentissage automatique

<http://www.reuters.com/article/2015/09/07/us-apple-machinelearning-idUSKCN0R71H020150907>

- Août 2016 : Apple rachète une start-up (Turi) en pointe dans le « machine learning »

<http://www.lesechos.fr/tech-medias/hightech/0211191589613-apple-rachete-une-start-up-en-pointe-dans-le-machine-learning-2019432.php>

Des compétences recherchées

- Septembre 2016 : Machine Learning : Apple s'offre TupleJump

<http://www.zdnet.fr/actualites/machine-learning-apple-s-offre-tuplejump-39842324.htm>

TupleJump fait du traitement des données et a développé la technologie FiloDB développé par TupleJump visant à associer machine learning et traitement d'importantes quantité de données.

Rappelons aussi que Apple recrute depuis 2015 une centaine de professionnels de l'intelligence artificielle pour enrichir ses logiciels et ses smartphones en les rendant plus intelligents grâce à des systèmes prédictifs.

Des compétences recherchées

- Mai 2017 : Apple acquiert Lattice Data qui a développé une solution d'apprentissage automatique pour pouvoir exploiter les dark data (données non structurées inutilisables sous leur forme brute)

<http://www.zdnet.fr/actualites/dark-data-apple-met-la-main-sur-lattice-data-39852434.htm>

Des compétences recherchées

- 2 juin 2015 : Facebook ouvre un centre de recherche dédié au Machine Learning à Paris
<http://www.lemondeinformatique.fr/actualites/lire-facebook-ouvre-un-centre-de-recherche-dedie-au-machine-learning-a-paris-61342.html>
- Janvier 2017 : Microsoft acquiert la jeune pousse canadienne Maluuba, spécialisée dans l'apprentissage automatique pour le traitement du langage naturel
<http://www.zdnet.fr/actualites/intelligence-artificielle-microsoft-realise-sa-premiere-acquisition-de-2017-avec-maluuba-39847160.htm>
- Sans oublier... Amazon, IBM, ...

Arbres de décision



Les diapositives suivantes se basent en grande partie sur le cours et le livre de Tom Mitchell. Machine Learning, McGraw Hill, 1997

Ensemble d'apprentissage

Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	Chaud	Elevée	Faible	Non
J2	Soleil	Chaud	Elevée	Fort	Non
J3	Couvert	Chaud	Elevée	Faible	Oui
J4	Pluie	Doux	Elevée	Faible	Oui
J5	Pluie	Froid	Normale	Faible	Oui
J6	Pluie	Froid	Normale	Fort	Non
J7	Couvert	Froid	Normale	Fort	Oui
J8	Soleil	Doux	Elevée	Faible	Non
J9	Soleil	Froid	Normale	Faible	Oui
J10	Pluie	Doux	Normale	Faible	Oui
J11	Soleil	Doux	Normale	Fort	Oui
J12	Couvert	Doux	Elevée	Fort	Oui
J13	Couvert	Chaud	Normale	Faible	Oui
J14	Pluie	Doux	Elevée	Fort	Non

Apprentissage supervisé

- Les données sont étiquetées : on se donne une valeur cible (la classe)
- On cherche à apprendre une fonction qui permet de prédire la valeur de la classe à partir des valeurs des autres attributs

Ensemble d'apprentissage

classe

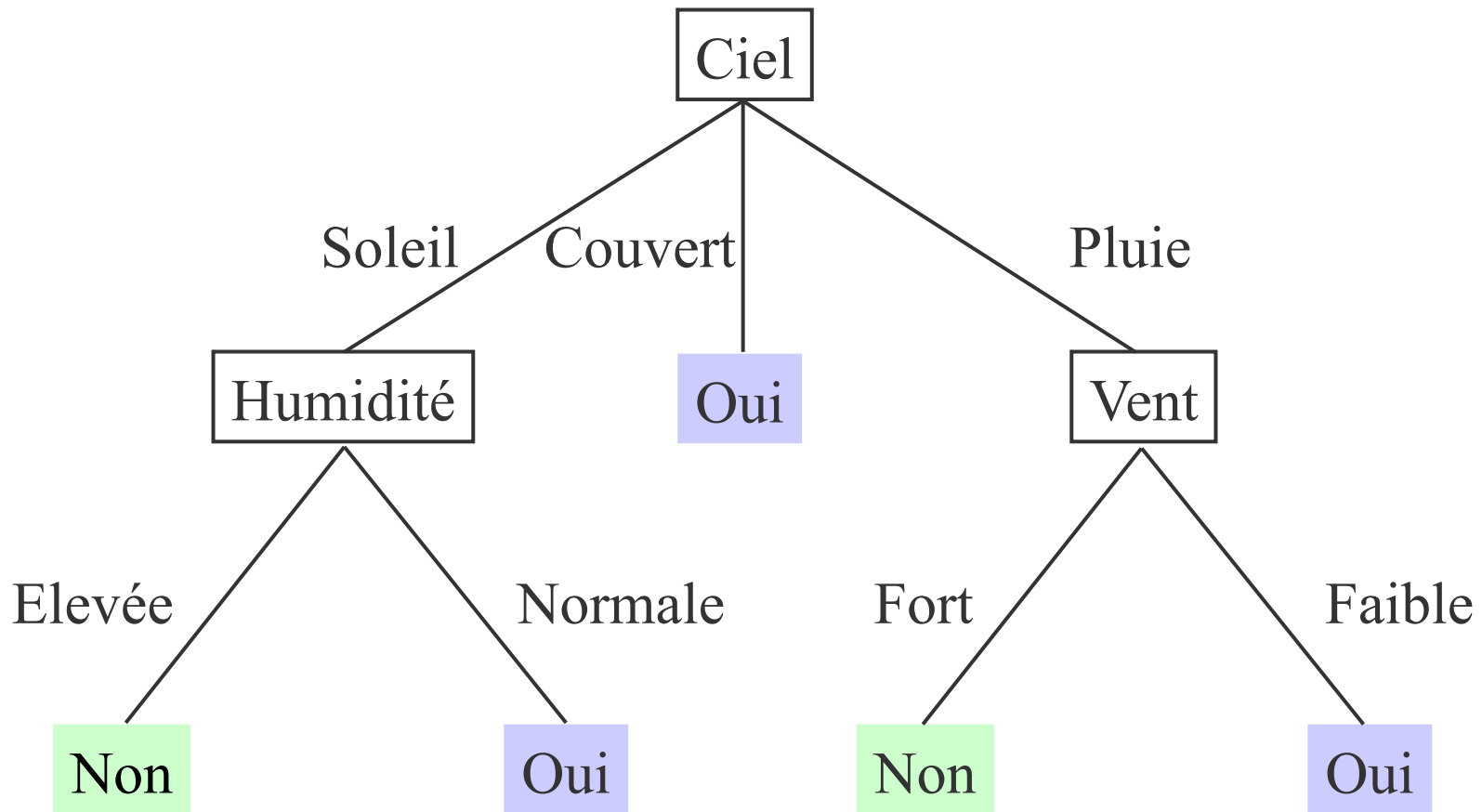
Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	Chaud	Elevée	Faible	Non
J2	Soleil	Chaud	Elevée	Fort	Non
J3	Couvert	Chaud	Elevée	Faible	Oui
J4	Pluie	Doux	Elevée	Faible	Oui
J5	Pluie	Froid	Normale	Faible	Oui
J6	Pluie	Froid	Normale	Fort	Non
J7	Couvert	Froid	Normale	Fort	Oui
J8	Soleil	Doux	Elevée	Faible	Non
J9	Soleil	Froid	Normale	Faible	Oui
J10	Pluie	Doux	Normale	Faible	Oui
J11	Soleil	Doux	Normale	Fort	Oui
J12	Couvert	Doux	Elevée	Fort	Oui
J13	Couvert	Chaud	Normale	Faible	Oui
J14	Pluie	Doux	Elevée	Fort	Non

Ensemble d'apprentissage

Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	Chaud	Elevée	Faible	Non
J2	Soleil	Chaud	Elevée	Fort	Non
J3	Couvert	Chaud	Elevée	Faible	Oui
J4	Pluie	Doux	Elevée	Faible	Oui
J5	Pluie	Froid	Normale	Faible	Oui
J6	Pluie	Froid	Normale	Fort	Non
J7	Couvert	Froid	Normale	Fort	Oui
J8	Soleil	Doux	Elevée	Faible	Non
J9	Soleil	Froid	Normale	Faible	Oui
J10	Pluie	Doux	Normale	Faible	Oui
J11	Soleil	Doux	Normale	Fort	Oui
J12	Couvert	Doux	Elevée	Fort	Oui
J13	Couvert	Chaud	Normale	Faible	Oui
J14	Pluie	Doux	Elevée	Fort	Non

=> apprendre
une fonction qui
permet de savoir
dans quelle
situation on
devrait jouer ou
non

Exemple d'arbre de décision

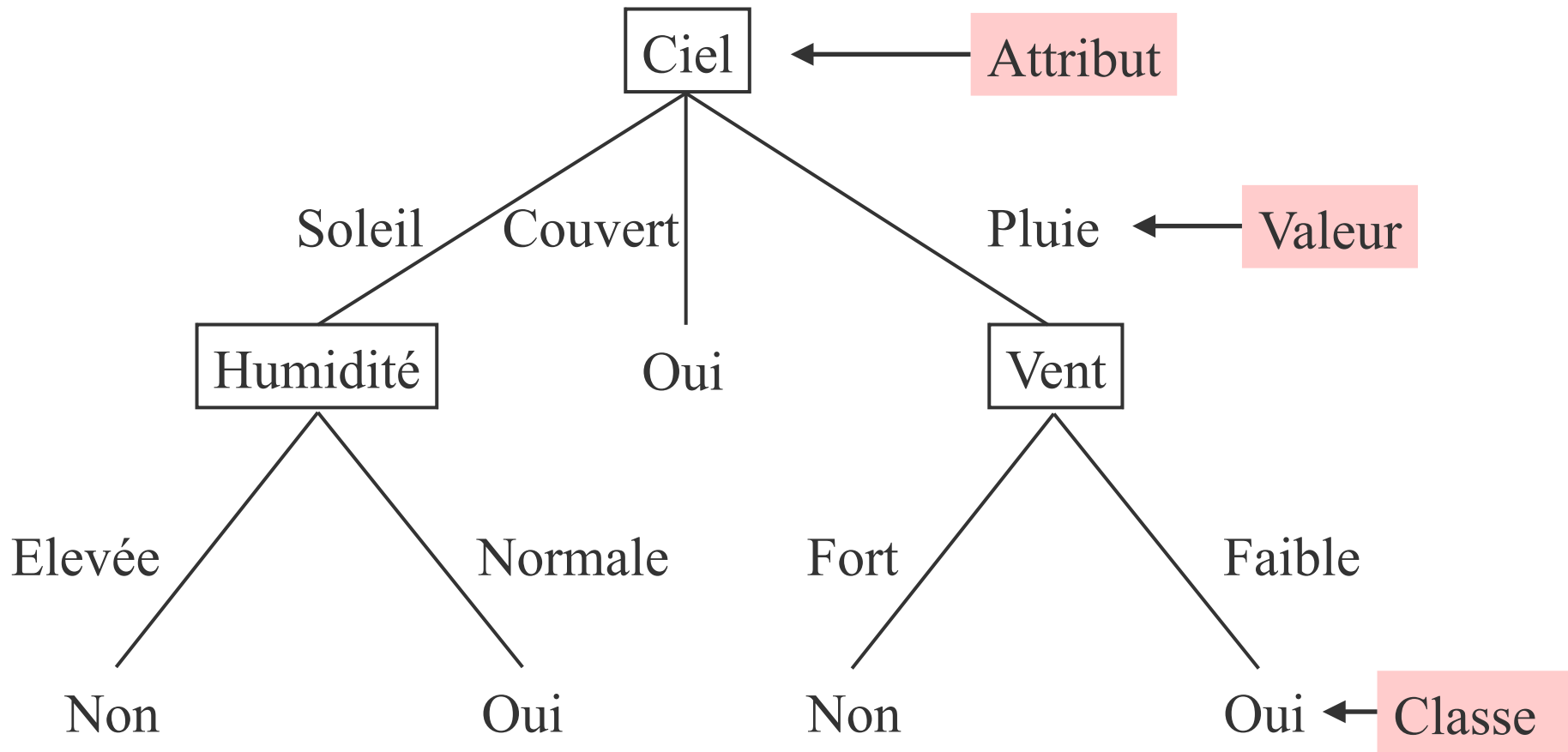


Plan

- Représentation
- ID3 (Iterative Dichotomiser 3) [Quin 86]
- Espace des hypothèses et biais inductif
- Extensions

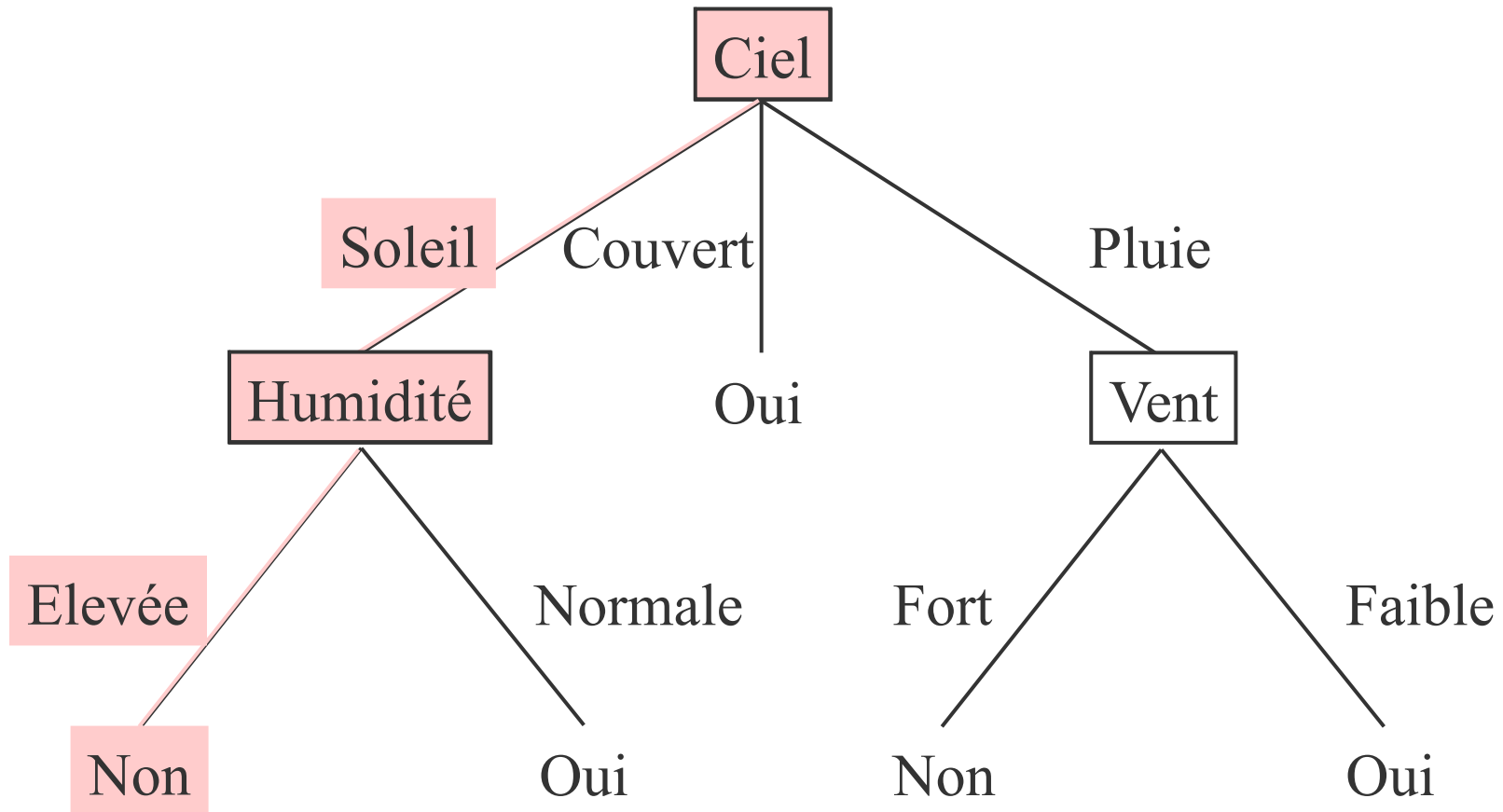
[Quin 86] Quinlan, J. R. 1986. Induction of Decision Trees. Mach. Learn. 1, 1 (Mar. 1986), 81-106

Représentation



Classification d'une instance

<Ciel = Soleil, Température = Chaud, Humidité = Elevée, Vent = Fort>



Expression sous forme logique

Un arbre de décision représente un ensemble de règles :

Si (Ciel = Soleil et Humidité = Normale)
ou (Ciel = Couvert)
ou (Ciel = Pluie et Vent = Faible)
Alors Jouer = oui

Intérêts des arbres de décision

- Expressivité
 - approximation de fonctions à valeurs catégorielles
 - capable d'apprendre des expressions disjonctives
- Lisibilité
 - peut être traduit sous la forme de règles
- Beaucoup d'applications

Domaines d'application

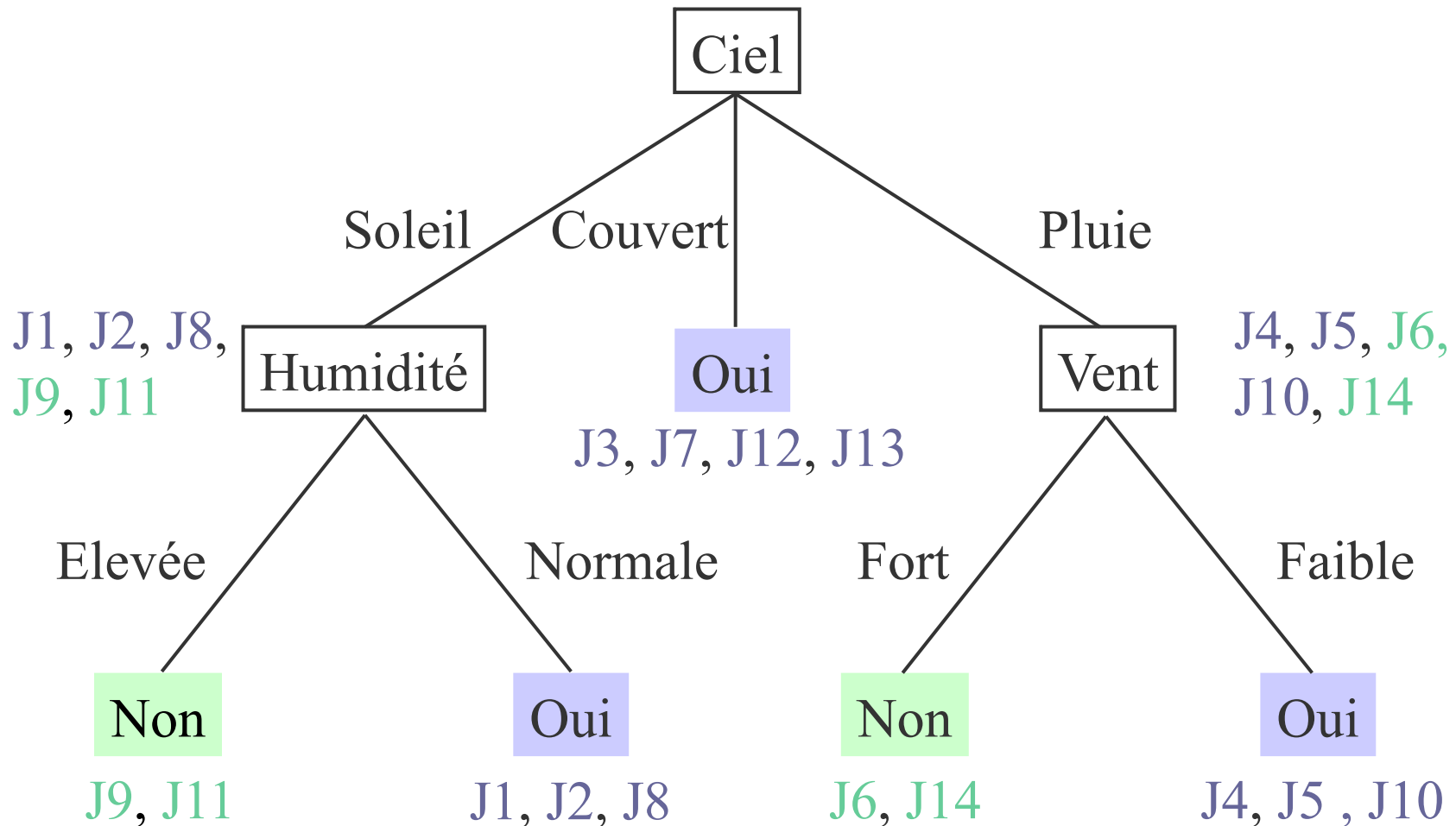
- Instances représentées par des couples attribut-valeur
- Fonction cible à valeurs catégorielles
- Expression disjonctive vraisemblable
- Erreurs possibles dans les exemples
- Valeurs manquantes

Médical, financier,...

Principes

- Induction descendante d'arbres de décision
- Approche diviser pour régner (identifier des sous-problèmes, leur trouver une solution et combiner ces solutions pour résoudre le problème général)
- À chaque nœud : meilleure façon de séparer en classes le sous-ensemble d'apprentissage

Exemple d'arbre de décision



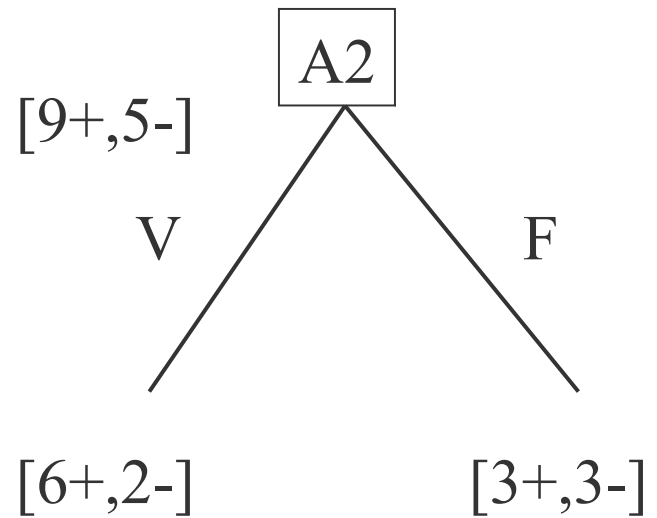
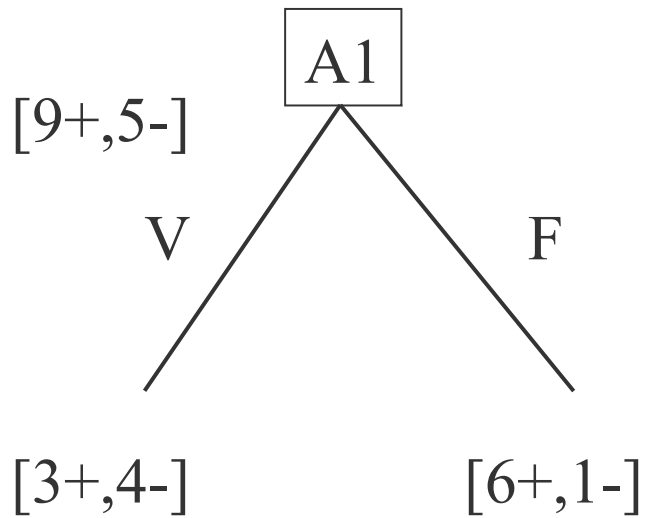
Algorithme de base

- $A = \text{MeilleurAttribut}(\text{Exemples})$
- Affecter A à la racine
- Pour chaque valeur de A , créer un nouveau nœud fils de la racine
- Classer les exemples dans les nœuds fils
- Si tous les exemples d'un nœud fils sont homogènes, affecter leur classe au nœud, sinon recommencer à partir de ce nœud

Ensemble d'apprentissage

Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	Chaud	Elevée	Faible	Non
J2	Soleil	Chaud	Elevée	Fort	Non
J3	Couvert	Chaud	Elevée	Faible	Oui
J4	Pluie	Doux	Elevée	Faible	Oui
J5	Pluie	Froid	Normale	Faible	Oui
J6	Pluie	Froid	Normale	Fort	Non
J7	Couvert	Froid	Normale	Fort	Oui
J8	Soleil	Doux	Elevée	Faible	Non
J9	Soleil	Froid	Normale	Faible	Oui
J10	Pluie	Doux	Normale	Faible	Oui
J11	Soleil	Doux	Normale	Fort	Oui
J12	Couvert	Doux	Elevée	Fort	Oui
J13	Couvert	Chaud	Normale	Faible	Oui
J14	Pluie	Doux	Elevée	Fort	Non

Choix de l'attribut



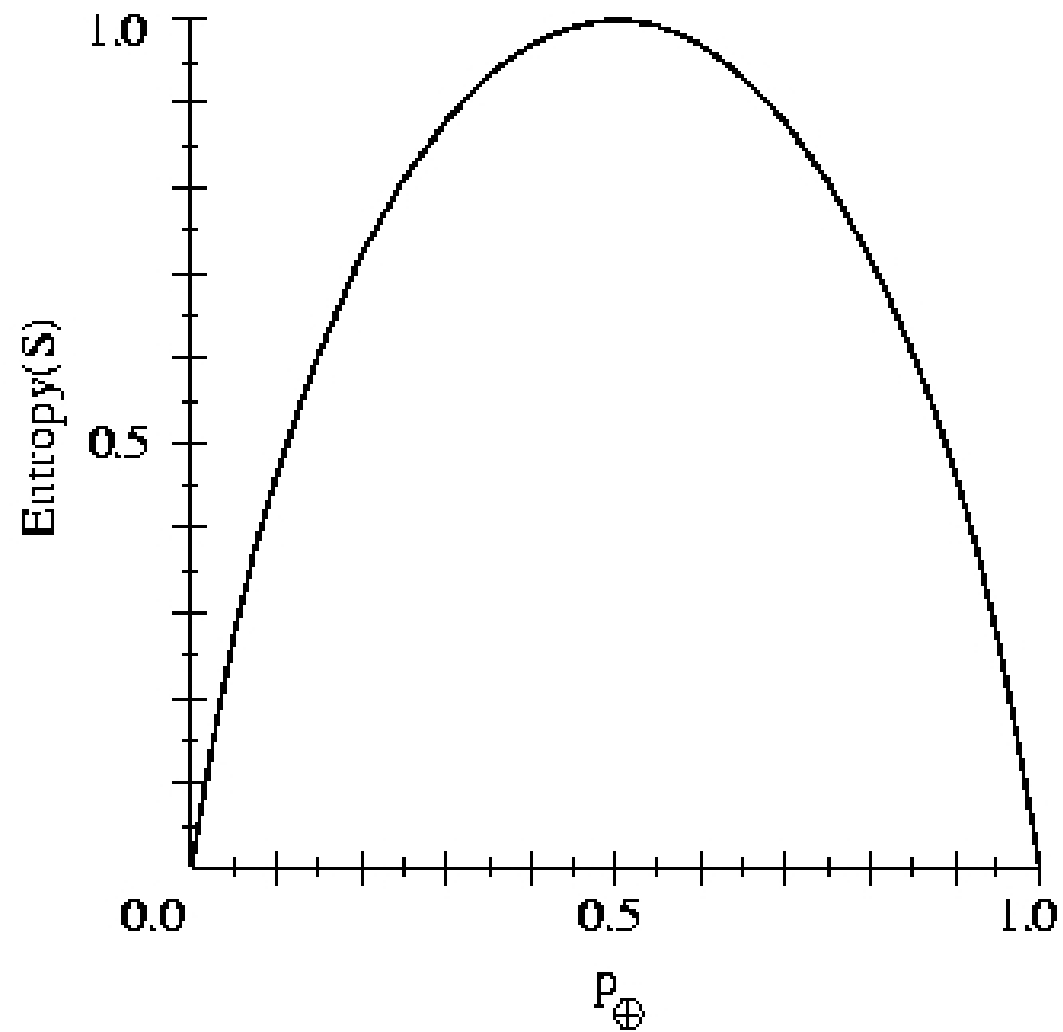
Entropie (cas d'une classe à deux valeurs)

$$\textit{Entropie}(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-$$

- S est un ensemble d'exemples
- p_+ est la proportion d'exemples positifs
- p_- est la proportion d'exemples négatifs
- Mesure l'homogénéité des exemples

$$\textit{Entropie}([9+, 5-]) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940$$

Entropy



Interprétation de l'entropie

- Nombre minimum de bits nécessaires pour coder la classe d'un élément quelconque de S
- Théorie de l'information : un code de longueur optimale utilise $-\log_2 p$ bits à un message de probabilité p .

$$\textit{Entropie}(S) \equiv p_+(-\log_2 p_+) + p_-(-\log_2 p_-)$$

Entropie : cas général

$$\textit{Entropie} (S) \equiv \sum_{i=1}^c - p_i (\log_2 p_i)$$

Gain d'information

- $\text{Gain}(S, A) = \text{Réduction d'entropie due à un tri suivant les valeurs de } A$

$$\text{Gain}(S, A) \equiv \text{Entropie}(S) - \sum_{v \in \text{Valeurs}(A)} \frac{|S_v|}{|S|} \text{Entropie}(S_v)$$

Choix de l'attribut

[9+,5-]

E=0,940

Humidité

Elevée

Normale

[3+,4-]

E=0,985

[6+,1-]

E=0,592

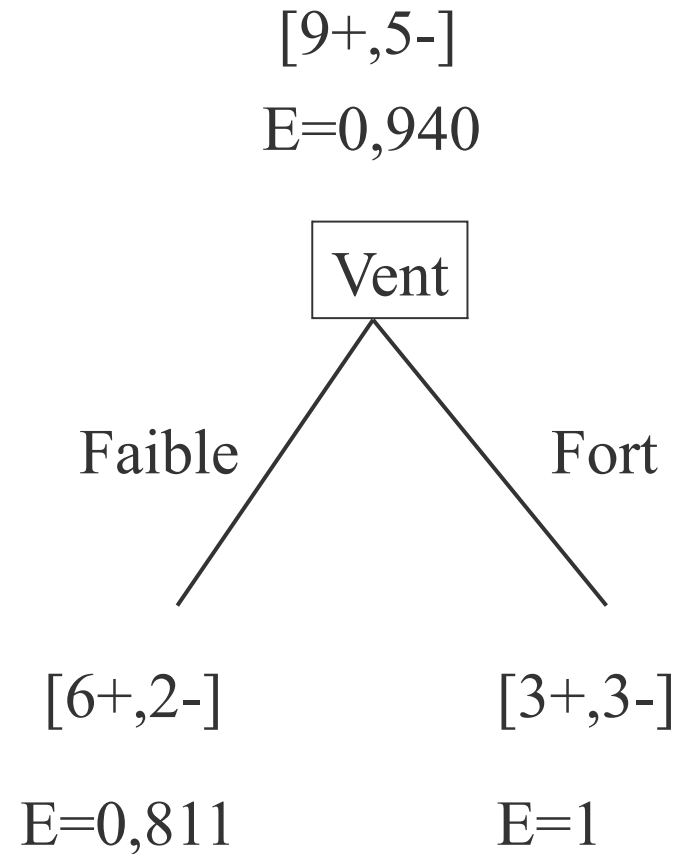
Gain(S, Humidité)

= 0,940-(7/14)0,985-(7/14)0,592

= 0,151

Choix de l'attribut

$$\begin{aligned} \text{Gain}(S, \text{Vent}) \\ &= 0,940 - (8/14)0,811 - (6/14) 1 \\ &= 0,048 \end{aligned}$$



Choix de l'attribut

$[9+, 5-]$

$E=0,940$

Ciel

Soleil

Couvert

Pluie

$[2+, 3-]$

$[4+, 0-]$

$[3+, 2-]$

$E=0,971$

$E=0$

$E=0,971$

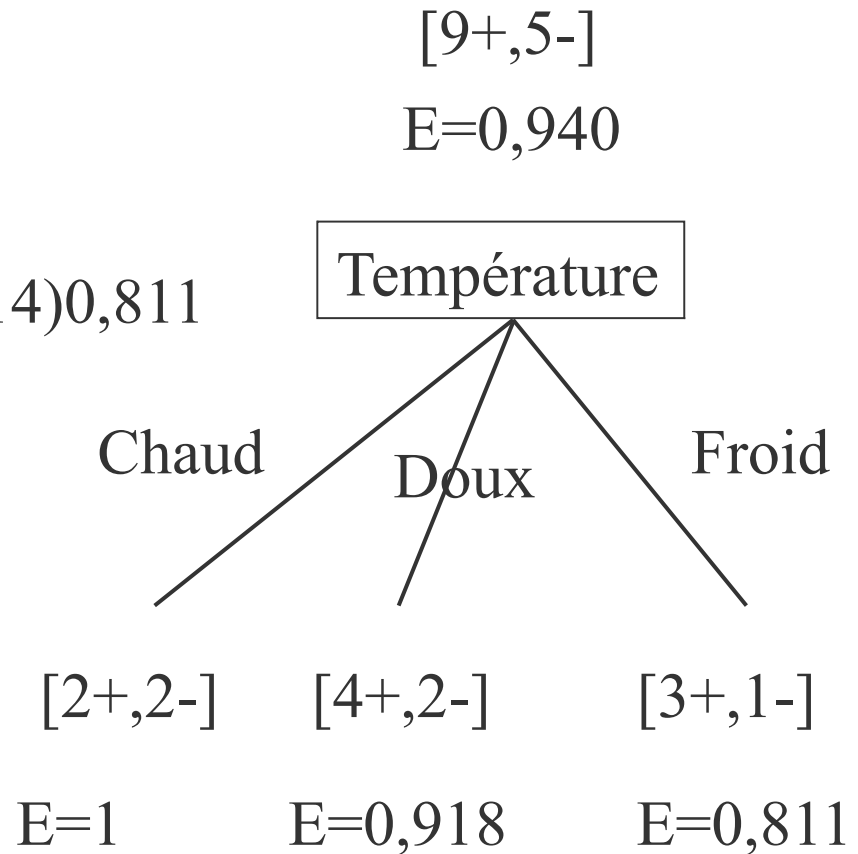
$\text{Gain}(S, \text{Ciel})$

$= 0,940 - (5/14)0,971 - (5/14)0,971 - 0$

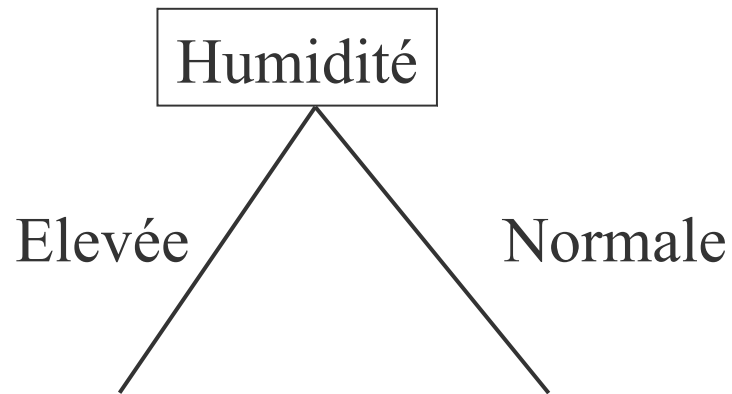
$= 0,246$

Choix de l'attribut

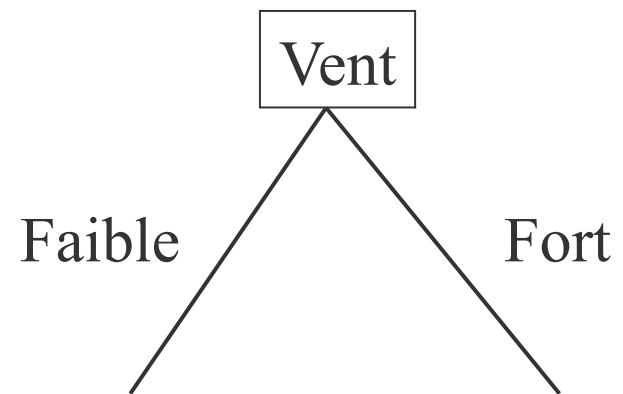
$$\begin{aligned} &\text{Gain}(S, \text{Température}) \\ &= 0,940 - (4/14)1 - (6/14)0,918 - (4/14)0,811 \\ &= 0,029 \end{aligned}$$



Choix de l'attribut

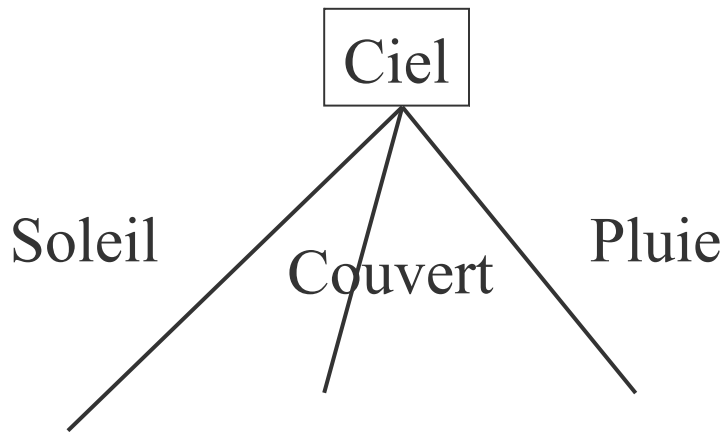


$$\text{Gain}(S, \text{Humidité}) = 0,151$$

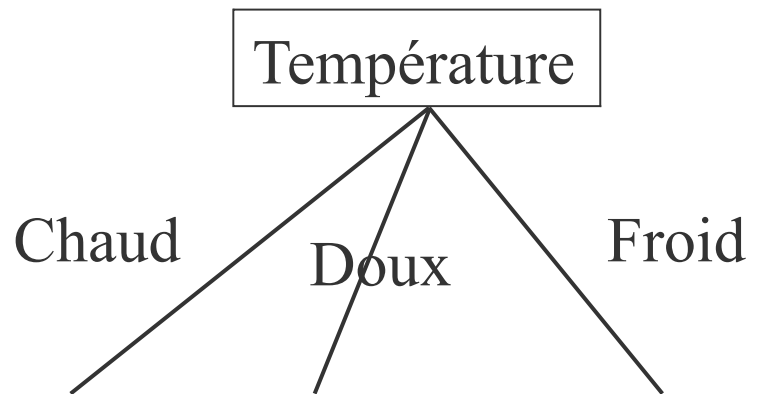


$$\text{Gain}(S, \text{Vent}) = 0,048$$

Choix de l'attribut

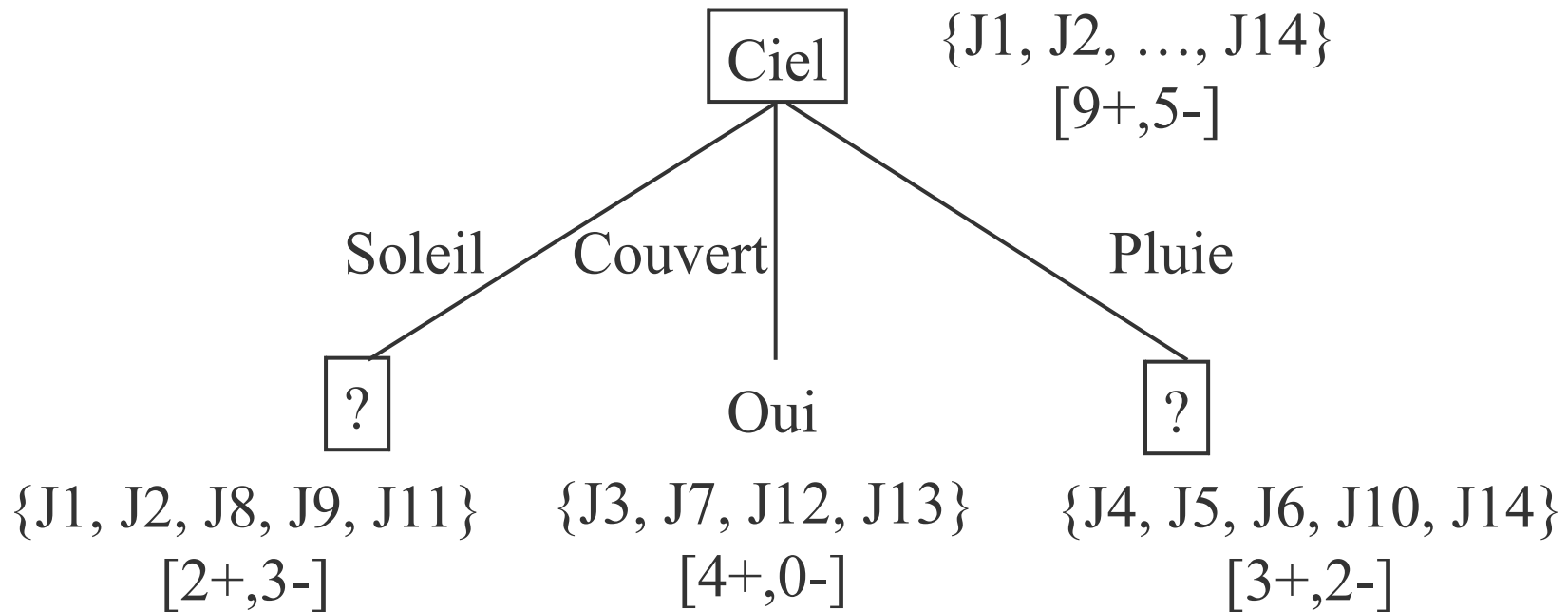


$$\text{Gain}(S, \text{Ciel}) = 0,246$$



$$\text{Gain}(S, \text{Température}) = 0,029$$

Choix du prochain attribut

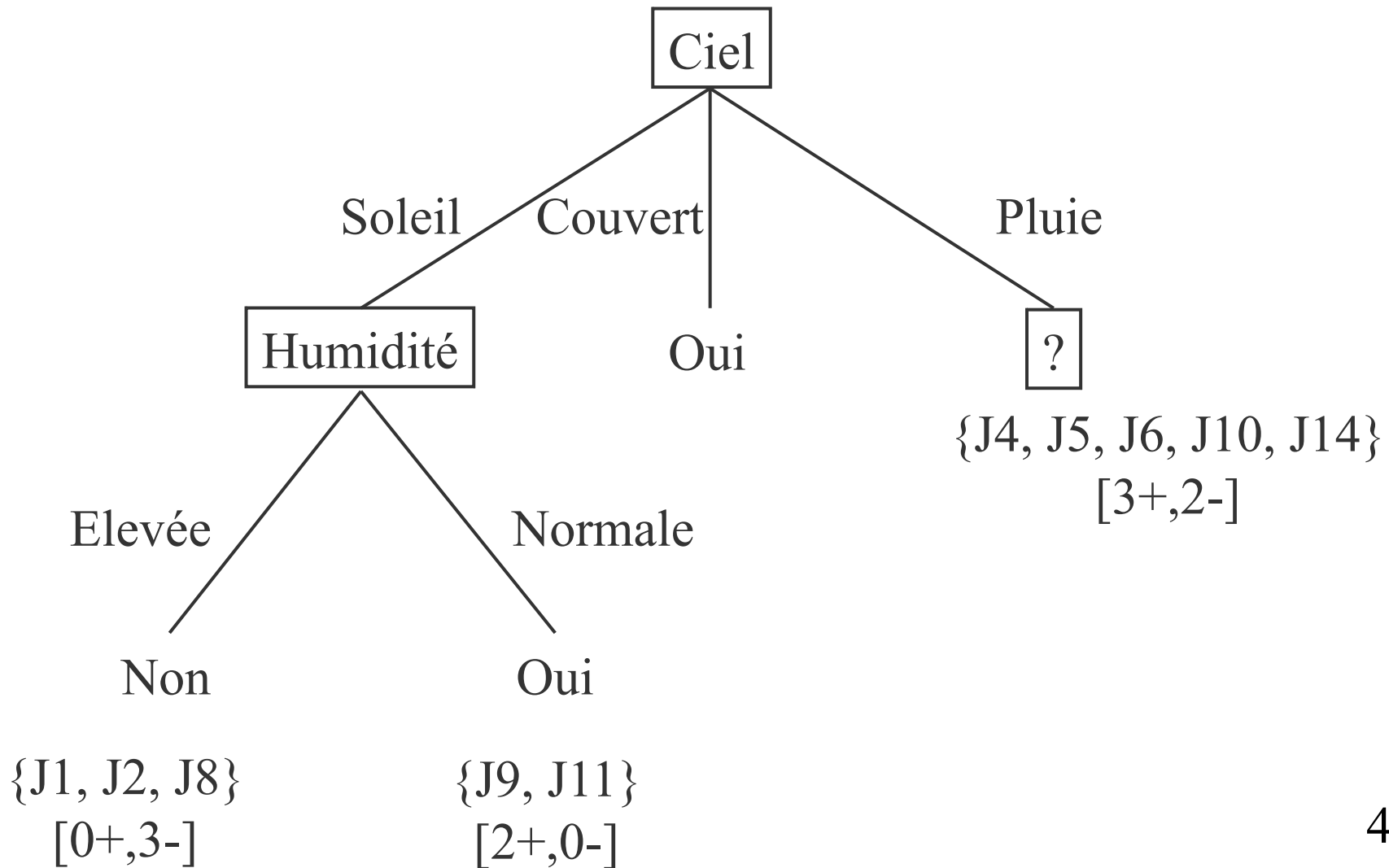


$$\text{Gain}(S_{\text{Soleil}}, \text{Humidité}) = 0,970 - (3/5) 0 - (2/5) 0 = 0,970$$

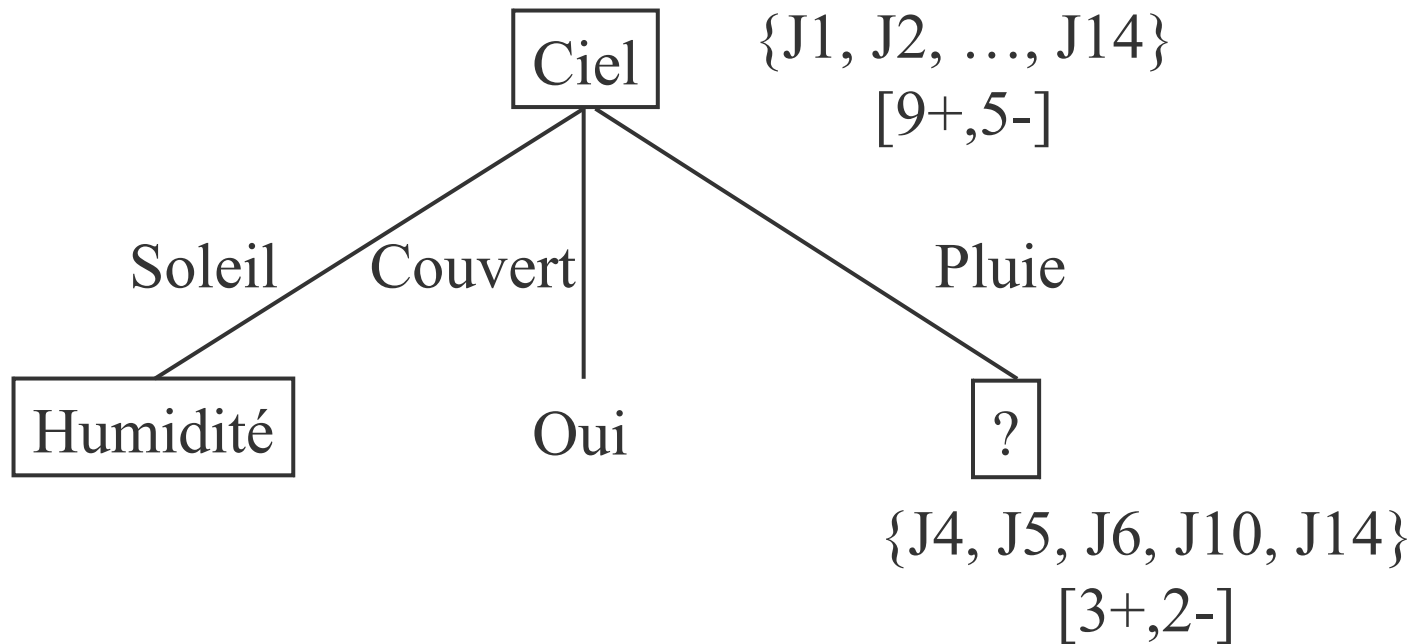
$$\text{Gain}(S_{\text{Soleil}}, \text{Température}) = 0,970 - (2/5) 0 - (2/5) 1 - (1/5) 0 = 0,570$$

$$\text{Gain}(S_{\text{Soleil}}, \text{Vent}) = 0,970 - (2/5) 1 - (3/5) 0,918 = 0,019$$

Exemple d'arbre de décision



Choix du prochain attribut

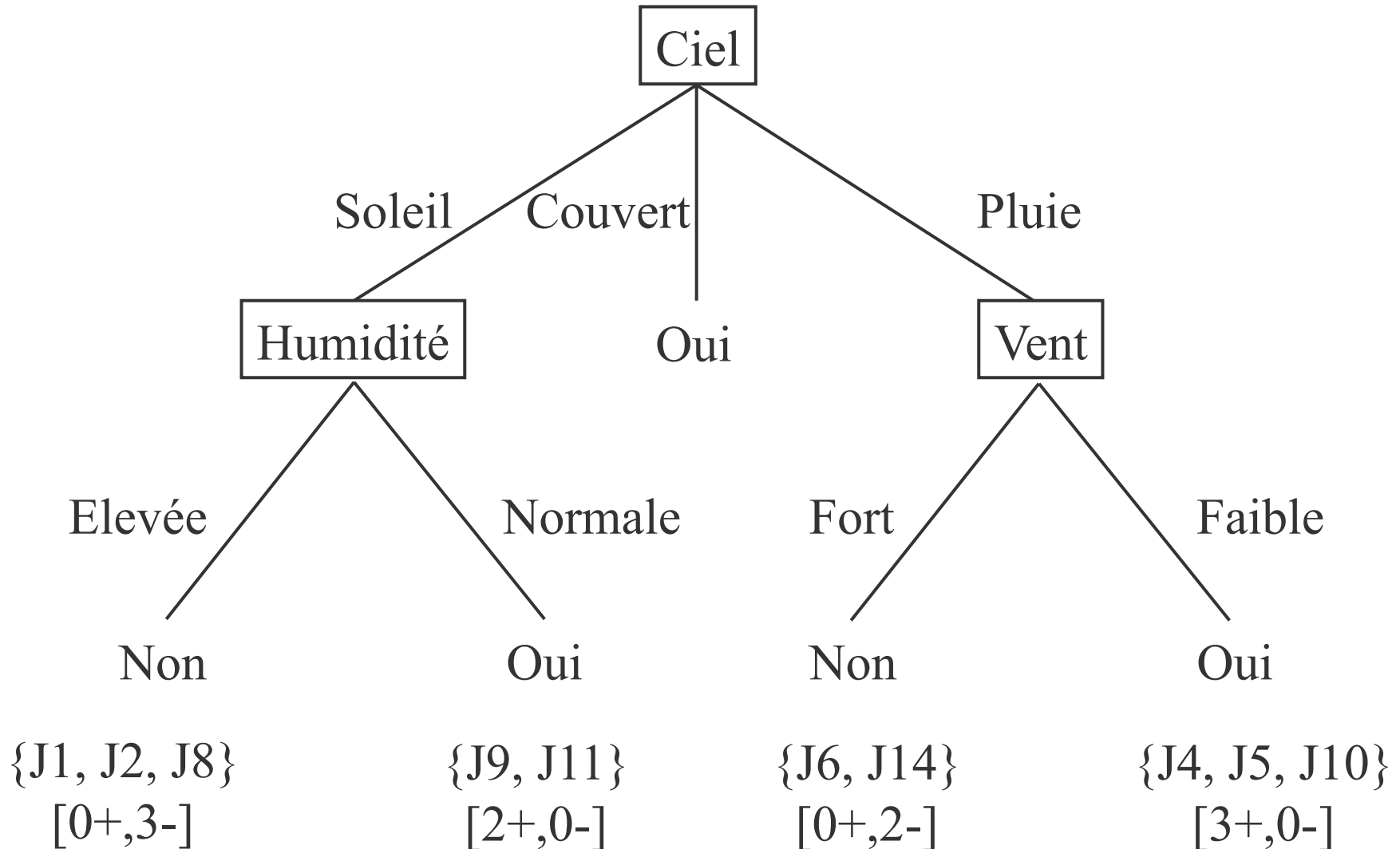


$$\text{Gain}(S_{\text{Pluie}}, \text{Humidité}) = 0,970 - (2/5) 1 - (3/5) 0,918 = 0,019$$

$$\text{Gain}(S_{\text{Pluie}}, \text{Température}) = 0,970 - (0/5) - (3/5) 0,918 - (2/5) 1 = 0,019$$

$$\text{Gain}(S_{\text{Pluie}}, \text{Vent}) = 0,970 - (2/5) 0 - (3/5) 0 = 0,970$$

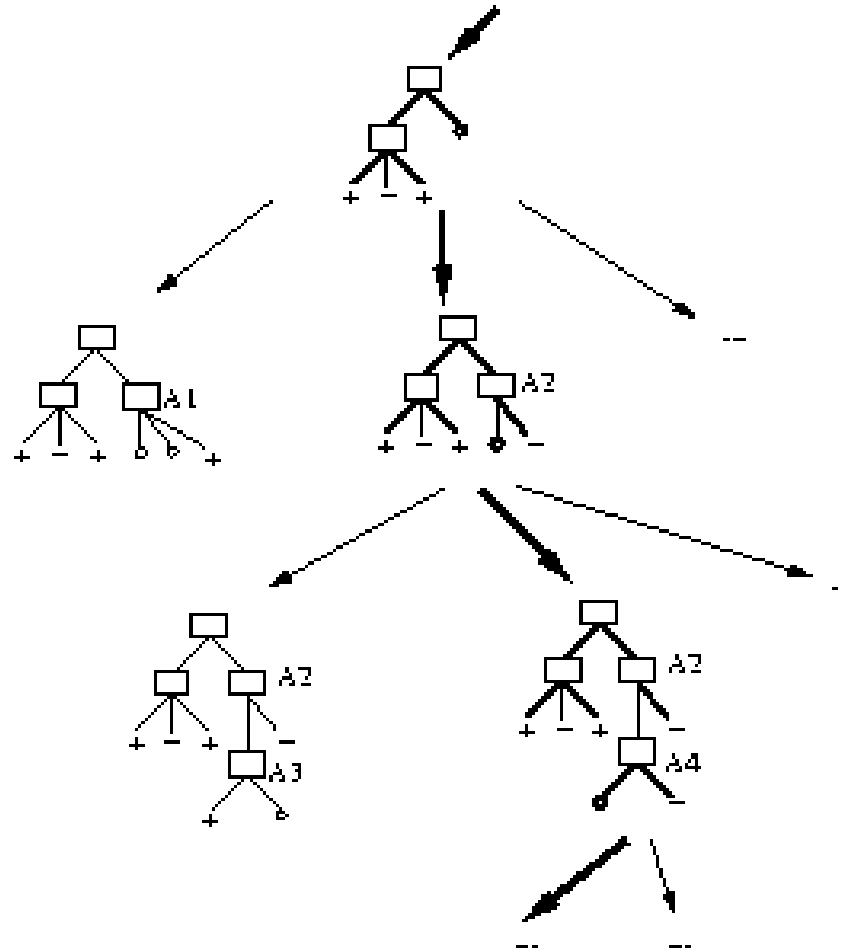
Exemple d'arbre de décision



Espace des hypothèses d'ID3

- Apprentissage vu comme une recherche dans un espace d'hypothèses
- « Hill-climbing » à partir de l'arbre vide, guidé par le gain d'information

Espace des hypothèses d'ID3



Espace des hypothèses d'ID3

- Espace des hypothèses est complet
- Rend une seule solution, pas toutes...
- Pas de retour en arrière
- Choix faits sur des critères statistiques

Biais inductif d'ID3

- « préfère les arbres les plus courts »
- ceux qui placent les attributs de meilleurs gains d'information près de la racine
- Approche heuristique d'une recherche en largeur d'abord

Types de biais

- Biais de restriction (a priori de l'espace de recherche)
 - biais de langage
- Biais de préférence
 - biais de recherche

Pourquoi préférer les hypothèses les plus courtes ?

- Rasoir d'Occam (biais de préférence)
 - ➔ Préférer les hypothèses les plus simples qui expliquent les données
- Plus générale, plus de chances d'être réfutée
- Taille de l'hypothèse dépend du langage

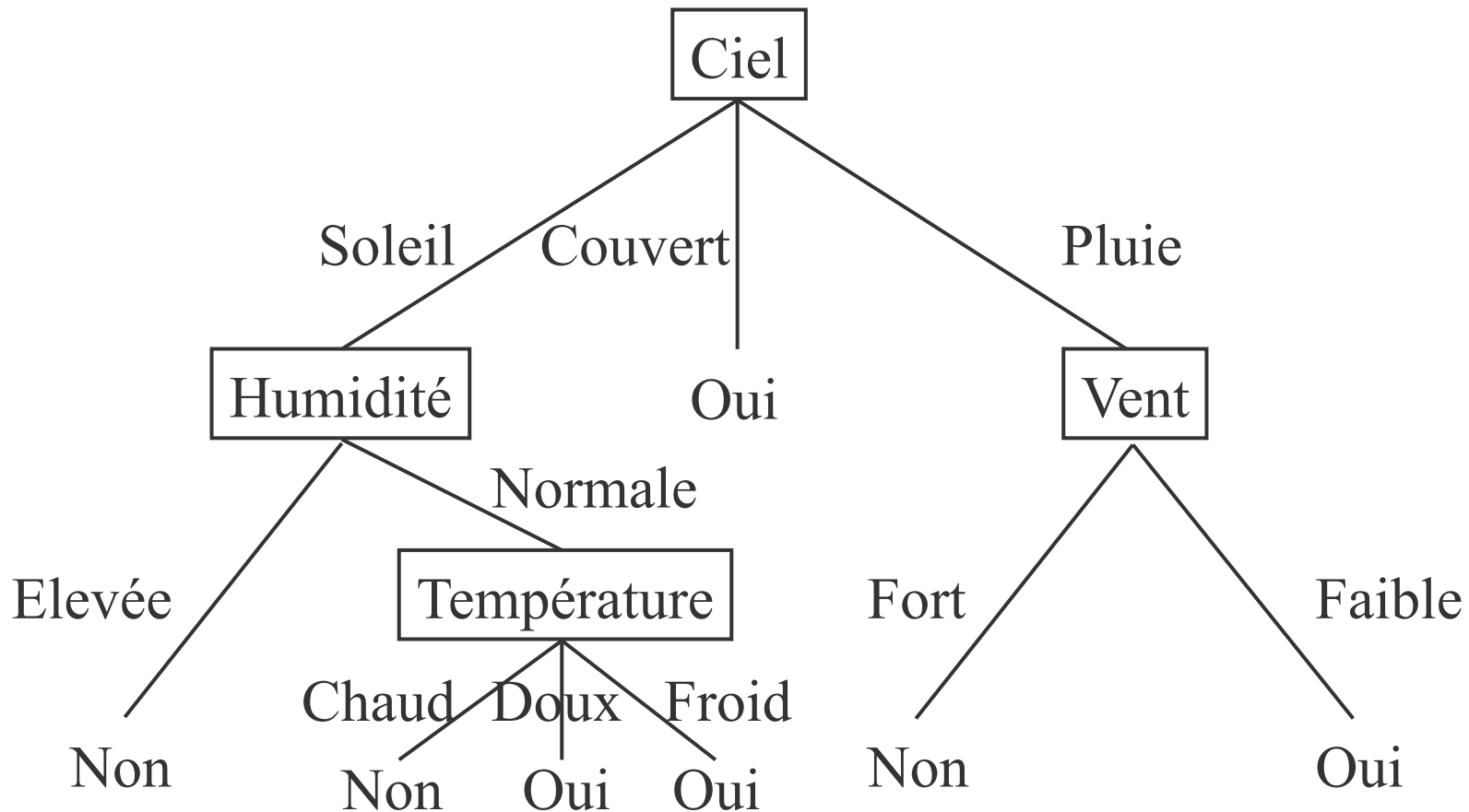
Extensions

- Sur-apprentissage
- Attributs continus
- Ratio du gain d'information
- Valeurs manquantes
- Coût des attributs

Sur-apprentissage

- Effet des exemples bruités :
 - ➔ J15 <Soleil, Chaud, Normale, Fort, Non>

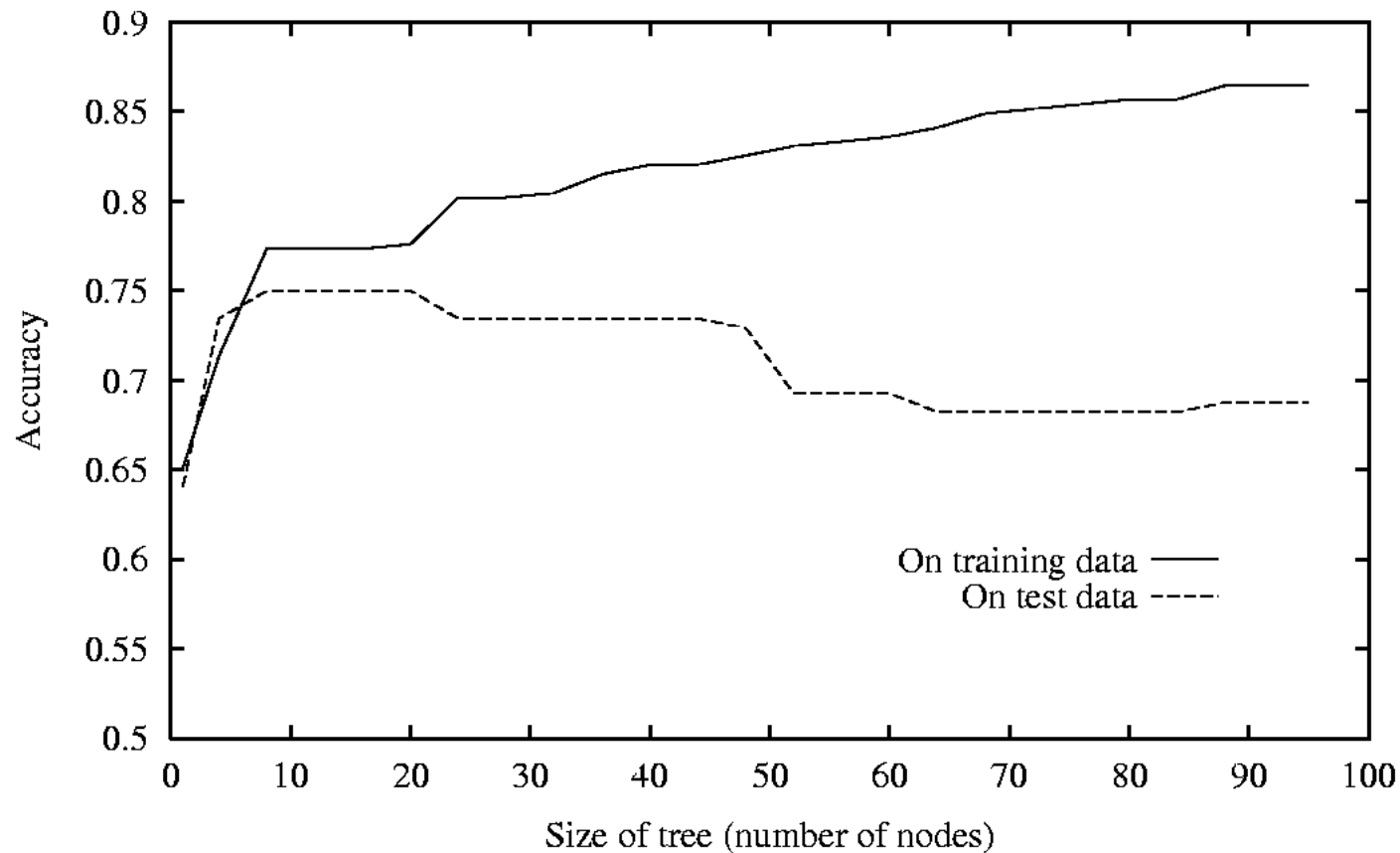
Sur-apprentissage



Sur-apprentissage

- Erreur expérimentale : $\text{erreur}_{\text{exp}}(h)$
- Erreur réelle sur la distribution D des instances : $\text{erreur}_D(h)$
- Sur-apprentissage : il existe h'
 - $\text{erreur}_{\text{exp}}(h) < \text{erreur}_{\text{exp}}(h')$
 - $\text{erreur}_D(h) > \text{erreur}_D(h')$

Exemple de sur-apprentissage



Eviter le sur-apprentissage

- Arrêter la croissance de l'arbre quand la division des données n'est plus statistiquement significative
- Générer l'arbre entier, puis élaguer

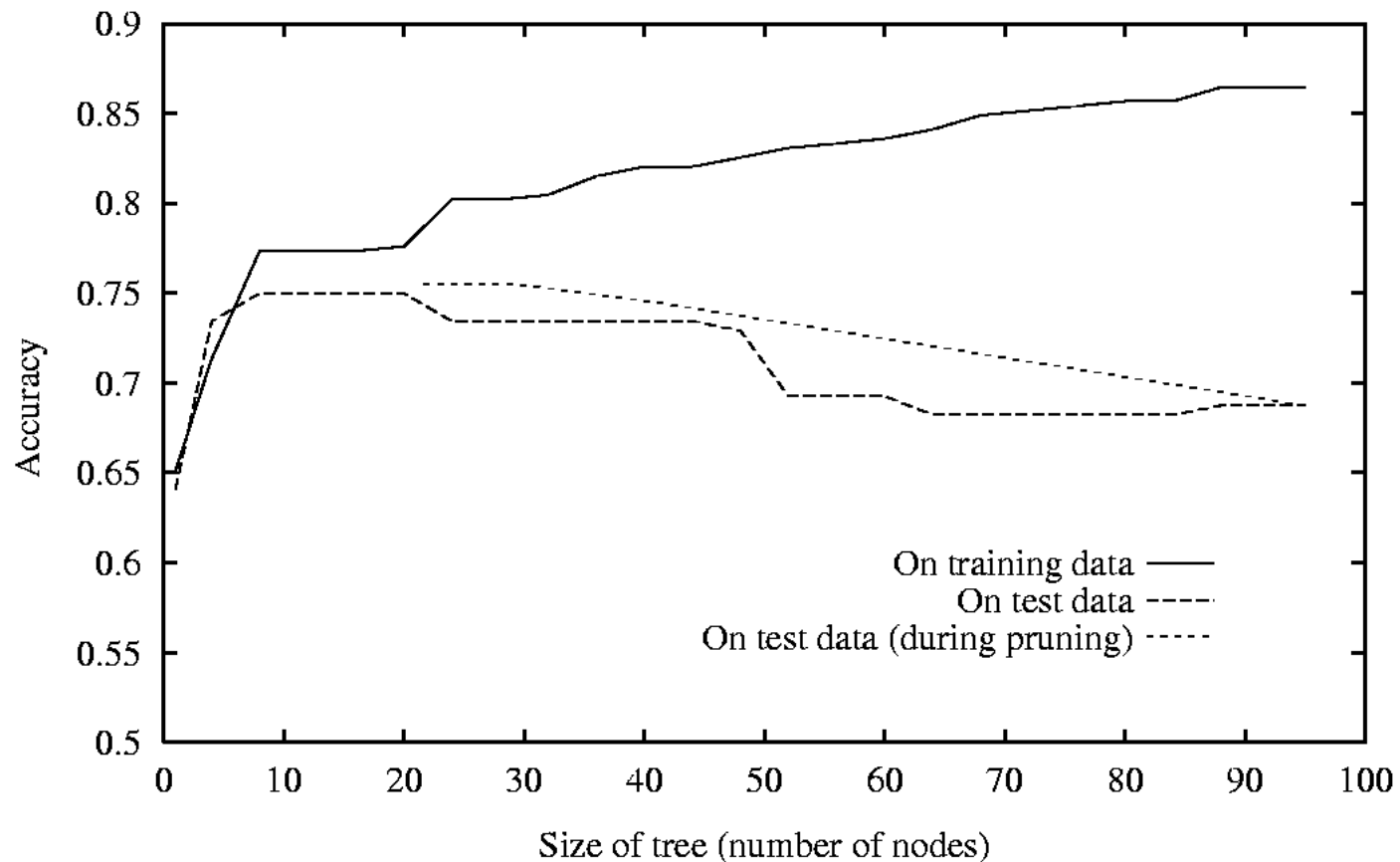
Sélection du « meilleur » arbre

- Mesurer les performances sur un ensemble distinct de données de validation
- Mesurer les performances sur l'ensemble d'apprentissage et effectuer test statistique
- MDL : minimiser $\text{taille}(\text{arbre}) + \text{taille}(\text{erreurs de classification}(\text{arbre}))$

Élagage basé sur l'erreur

- Diviser les données en ensembles d'apprentissage et de validation
- Tant que l'élagage réduit l'erreur
 - ➔ Evaluer sur l'ensemble de validation l'influence d'un élagage à partir de chaque nœud
 - ➔ Effectuer le meilleur élagage

Effet de l'élagage basé sur l'erreur



Post-élagage des règles

- Convertir l'arbre en un ensemble de règles
- Elaguer chaque règle indépendamment
- Ordonner les règles obtenues en fonction de leur précision
- Estimation de la précision d'une règle
 - ensemble de validation
 - estimation pessimiste (C4.5)

Intérêts de la conversion en règles

- Supprime la distinction entre les nœuds
- Plus flexible que l'élagage de l'arbre
- Améliore la lisibilité

Estimation pessimiste C4.5

- Ciel=couvert \wedge Humidité=normale \rightarrow Oui
 - $\rightarrow n=40, r=12, \text{erreur}_S = 12/40 = 0,3$
 - $\rightarrow \sigma = [r/n(1-r/n)/n]^{1/2} = [0,3 \times 0,7/40]^{1/2} = 0,07$
 - $\rightarrow \text{erreur}_S + z_N \sigma = 0,3 + 1,96 \times 0,07 = 0,437$
- Humidité=normale \rightarrow Oui
 - $\rightarrow n=160, r=56, \text{erreur}_S = 56/160 = 0,35$
 - $\rightarrow \sigma = [r/n(1-r/n)/n]^{1/2} = [0,35 \times 0,65/160]^{1/2} = 0,04$
 - $\rightarrow \text{erreur}_S + z_N \sigma = 0,35 + 1,96 \times 0,04 = 0,424$

Attributs continus

Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	85	85	Faible	Non
J2	Soleil	80	90	Fort	Non
J3	Couvert	83	86	Faible	Oui
J4	Pluie	70	96	Faible	Oui
J5	Pluie	68	80	Faible	Oui
J6	Pluie	65	70	Fort	Non
J7	Couvert	64	65	Fort	Oui
J8	Soleil	72	75	Faible	Non
J9	Soleil	69	70	Faible	Oui
J10	Pluie	75	80	Faible	Oui
J11	Soleil	75	70	Fort	Oui
J12	Couvert	72	90	Fort	Oui
J13	Couvert	81	75	Faible	Oui
J14	Pluie	71	91	Fort	Non

Attributs continus

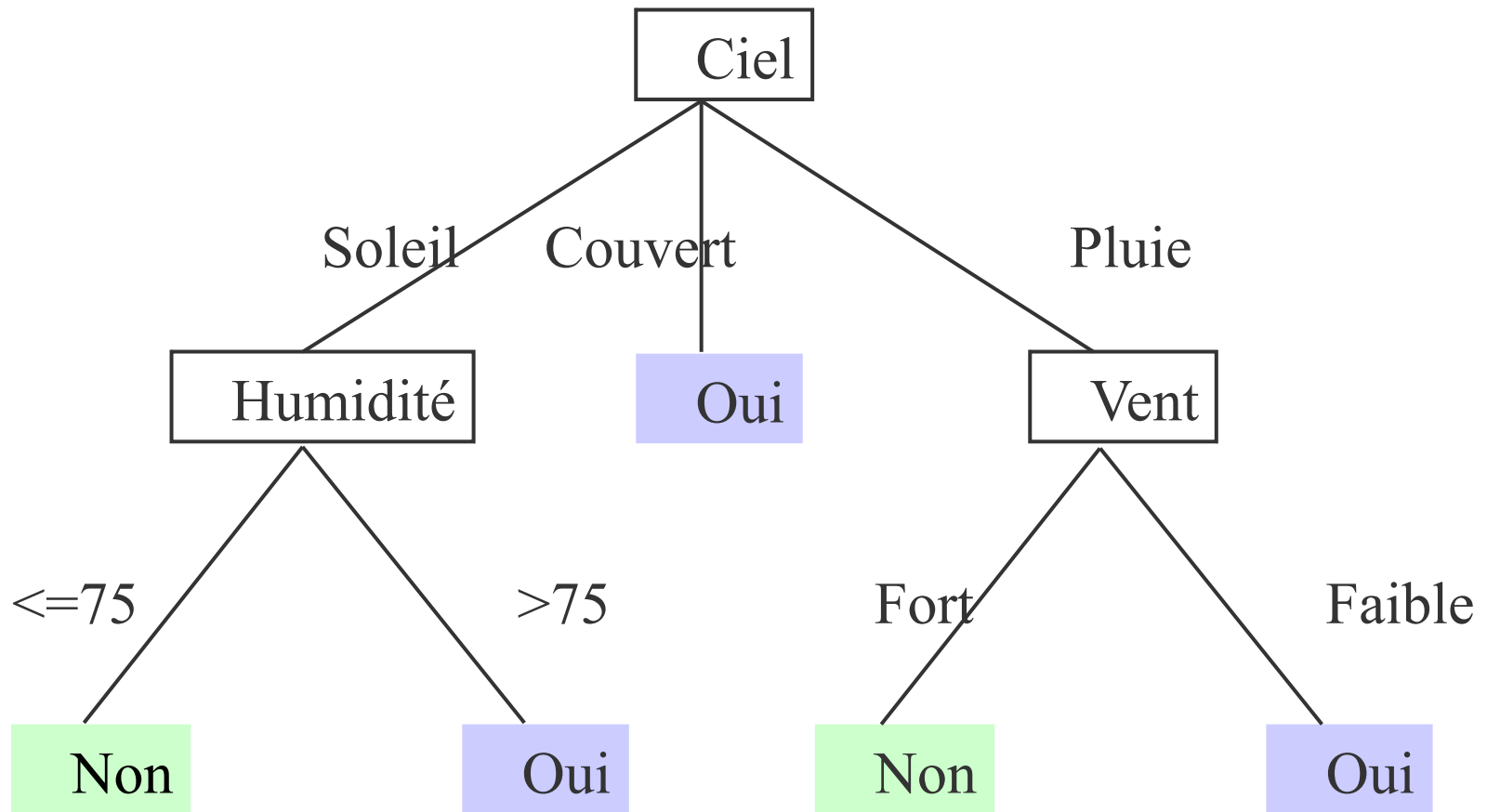
- Créer un attribut discret
→ (Température > 75) = V, F

Température	40	48	60	72	80	90
Jouer	Non	Non	Oui	Oui	Non	Non

Choix du seuil

- Entropie $[2+,4-] = 0,918$
- $\text{Gain}_1 = 0,918 - (1/6)0 - (5/6)0,970 = 0,109$
- $\text{Gain}_2 = 0,918 - (2/6)0 - (4/6)1 = 0,251$
- $\text{Gain}_3 = 0,918 - (3/6)0,918 * 2 = 0 !$
- $\text{Gain}_4 = 0,918 - (4/6)1 - (2/6)0 = 0,251$
- $\text{Gain}_5 = 0,918 - (5/6)0,970 - (1/6)0 = 0,109$

Exemple d'arbre de décision



Ensemble d'apprentissage

Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	85	85	Faible	Non
J2	Soleil	80	90	Fort	Non
J3	Couvert	83	86	Faible	Oui
J4	Pluie	70	96	Faible	Oui
J5	Pluie	68	80	Faible	Oui
J6	Pluie	65	70	Fort	Non
J7	Couvert	64	65	Fort	Oui
J8	Soleil	72	75	Faible	Non
J9	Soleil	69	70	Faible	Oui
J10	Pluie	75	80	Faible	Oui
J11	Soleil	75	70	Fort	Oui
J12	Couvert	72	90	Fort	Oui
J13	Couvert	81	75	Faible	Oui
J14	Pluie	71	91	Fort	Non

?



Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	Chaud	Elevée	Faible	Non
J2	Soleil	Chaud	Elevée	Fort	Non
J3	Couvert	Chaud	Elevée	Faible	Oui
J4	Pluie	Doux	Elevée	Faible	Oui
J5	Pluie	Froid	Normale	Faible	Oui
J6	Pluie	Froid	Normale	Fort	Non
J7	Couvert	Froid	Normale	Fort	Oui
J8	Soleil	Doux	Elevée	Faible	Non
J9	Soleil	Froid	Normale	Faible	Oui
J10	Pluie	Doux	Normale	Faible	Oui
J11	Soleil	Doux	Normale	Fort	Oui
J12	Couvert	Doux	Elevée	Fort	Oui
J13	Couvert	Chaud	Normale	Faible	Oui
J14	Pluie	Doux	Elevée	Fort	Non

Ensemble d'apprentissage

Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	85	85	Faible	Non
J2	Soleil	80	90	Fort	Non
J3	Couvert	83	86	Faible	Oui
J4	Pluie	70	96	Faible	Oui
J5	Pluie	68	80	Faible	Oui
J6	Pluie	65	70	Fort	Non
J7	Couvert	64	65	Fort	Oui
J8	Soleil	72	75	Faible	Non
J9	Soleil	69	70	Faible	Oui
J10	Pluie	75	80	Faible	Oui
J11	Soleil	75	70	Fort	Oui
J12	Couvert	72	90	Fort	Oui
J13	Couvert	81	75	Faible	Oui
J14	Pluie	71	91	Fort	Non

?



Jour	Ciel	Température	Humidité	Vent	Jouer
J1	Soleil	Chaud	Elevée	Faible	Non
J2	Soleil	Chaud	Elevée	Fort	Non
J3	Couvert	Chaud	Elevée	Faible	Oui
J4	Pluie	Doux	Elevée	Faible	Oui
J5	Pluie	Froid	Normale	Faible	Oui
J6	Pluie	Froid	Normale	Fort	Non
J7	Couvert	Froid	Normale	Fort	Oui
J8	Soleil	Doux	Elevée	Faible	Non
J9	Soleil	Froid	Normale	Faible	Oui
J10	Pluie	Doux	Normale	Faible	Oui
J11	Soleil	Doux	Normale	Fort	Oui
J12	Couvert	Doux	Elevée	Fort	Oui
J13	Couvert	Chaud	Normale	Faible	Oui
J14	Pluie	Doux	Elevée	Fort	Non

Température $\geq 80 \Rightarrow$ Chaud

$70 \leq$ Température $< 80 \Rightarrow$ Doux

Température $< 70 \Rightarrow$ Froid

Humidité $> 75 \Rightarrow$ Elevée

Humidité $\leq 75 \Rightarrow$ Normale

Ratio du gain d'information

- Si un attribut a beaucoup de valeurs, le gain d'information le sélectionnera
 - ➔ Plus les ensembles sont petits, plus ils sont purs

$$GainRatio(G, A) \equiv \frac{Gain(S, A)}{SplitInformation(S, A)}$$

$$SplitInformation(S, A) \equiv - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

73

Valeurs manquantes

- Si un exemple n'a pas de valeur pour A
 - ➔ si un nœud teste A, utiliser la valeur la plus commune parmi les exemples de ce nœud
 - ➔ utiliser la valeur la plus fréquente parmi les exemples de la même classe
 - ➔ affecter une probabilité à chaque valeur de A

Coût des attributs

- Exemples
 - diagnostic médical
 - robotique
- Heuristiques

$$\frac{Gain^2(S, A)}{Coût(A)}$$

$$\frac{2^{Gain(S, A)} - 1}{(Coût(A) + 1)^w}, w \in [0; 1]$$