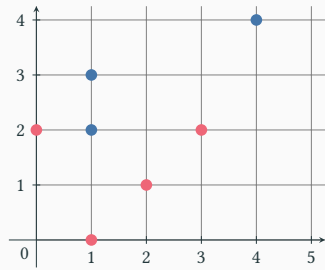


Étant donnée la représentation tabulaire suivante d'un corpus, attribuez les classes "?" par l'heuristique des 3-plus proches voisins

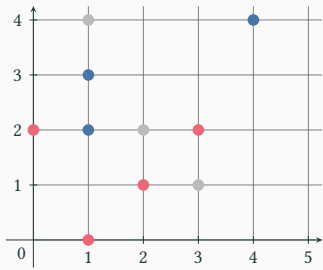
t	wug	neurone	classe
1	3	2	linguistique
2	1	0	linguistique
3	0	2	linguistique
4	2	1	linguistique
5	4	4	informatique
6	1	3	informatique
7	1	2	informatique
8	2	2	?
9	1	4	?
10	3	1	?

3



Représentation du corpus dans le plan

4



Représentation du corpus dans le plan

4

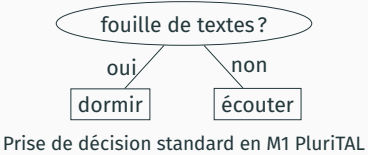
Arbres de décision

Arbres de décision

Arbre de décision

**Dans Weka trees (> J48)**  
**Espace de recherche** L'ensemble des arbres de recherche pour les attributs choisis  
**Techniques de recherche** Plusieurs algorithmes, le plus connu étant C4.5, appelé J48 dans Weka

Modèle de prise de décision déterministe



Prise de décision standard en M1 PluriTAL

5

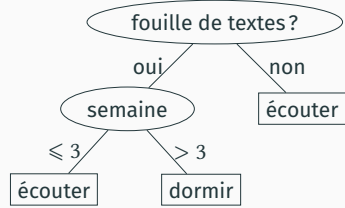
6

Arbre de décision

Arbre de décision

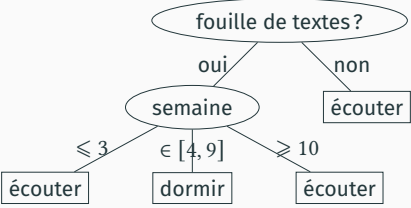
Exercice

Modèle de prise de décision déterministe par succession de choix



Prise de décision standard en M1 PluriTAL

Modèle de prise de décision déterministe par succession de choix (pas nécessairement binaires)



Prise de décision standard en M1 PluriTAL

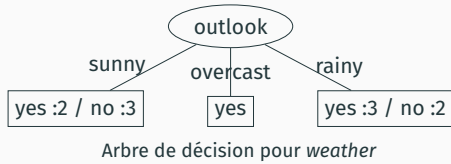
Construire un arbre de décision pour les données suivantes

No	outlook	humidity	windy	play
1	sunny	high	false	no
2	sunny	high	true	no
3	overcast	high	false	yes
4	rainy	high	false	yes
5	rainy	high	false	yes
6	rainy	low	true	no
7	overcast	low	true	yes
8	sunny	high	false	no
9	sunny	low	false	yes
10	rainy	high	false	yes
11	sunny	low	true	yes
12	overcast	high	true	yes
13	overcast	low	false	yes
14	rainy	high	true	no

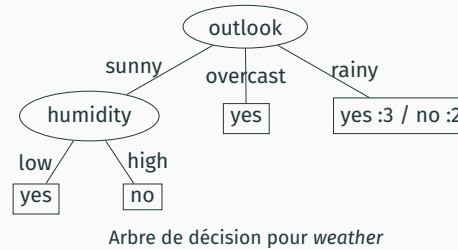
7

8

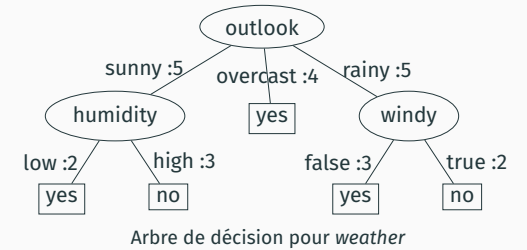
9



10



10



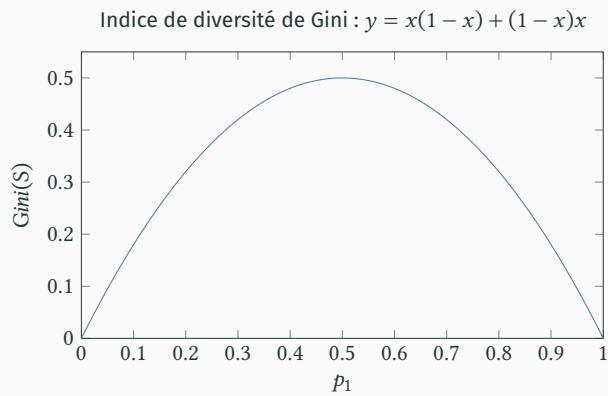
10

Arbre de décision	Technique d'apprentissage	Technique d'apprentissage
<pre>def classify(outlook, humidity, windy):     if outlook == "sunny":         if humidity == "low":             return "yes"         else:             return "no"     elif outlook == "overcast":         return "yes"     elif outlook == "rainy":         if windy:             return "no"         else:             return "yes"</pre> <p>Code Python correspondant à l'arbre précédent (voir icampus)</p>	<p>C'est assez facile de construire un arbre parfait pour l'ensemble d'entraînement :</p> <ul style="list-style-type: none"> <li>Il suffit d'énumérer tous les attributs jusqu'à avoir généré toutes les combinaisons existantes</li> </ul>	<p>C'est assez facile de construire un arbre parfait pour l'ensemble d'entraînement :</p> <ul style="list-style-type: none"> <li>Il suffit d'énumérer tous les attributs jusqu'à avoir généré toutes les combinaisons existantes</li> </ul> <p>→ Surapprentissage!</p>
<p>Technique d'apprentissage</p> <p>C'est assez facile de construire un arbre parfait pour l'ensemble d'entraînement :</p> <ul style="list-style-type: none"> <li>Il suffit d'énumérer tous les attributs jusqu'à avoir généré toutes les combinaisons existantes</li> </ul> <p>→ Surapprentissage!</p> <p>Comment faire pour avoir un arbre bon, mais pas trop profond ?</p>	<p>Technique d'apprentissage</p> <p>C'est assez facile de construire un arbre parfait pour l'ensemble d'entraînement :</p> <ul style="list-style-type: none"> <li>Il suffit d'énumérer tous les attributs jusqu'à avoir généré toutes les combinaisons existantes</li> </ul> <p>→ Surapprentissage!</p> <p>Comment faire pour avoir un arbre bon, mais pas trop profond ?</p> <ul style="list-style-type: none"> <li>Faire en sorte de trier vite et bien</li> </ul> <p>→ En choisissant les attributs les plus discriminants</p> <p>L'idée est de construire l'arbre progressivement, prenant à chaque étape le test le plus <b>discriminant</b>, reste à savoir comment on le détermine.</p>	<p>Indice de diversité de Gini</p> <p><b>Définition</b> On appelle <i>indice de diversité de Gini</i> d'une partition <math>S = \sqcup_{1 \leq i \leq n} c_i</math></p> $Gini(S) = \sum_{1 \leq i \leq n} p_i(1 - p_i)$ <p>avec <math>p_i = \frac{\#c_i}{\sum_{j=1}^n \#c_j}</math></p> <p>Autrement dit, l'indice de diversité de Gini est la probabilité qu'un exemple choisi au hasard et classé au hasard soit mal classé.</p> <p>Cet indice est d'autant plus élevé que la partition sépare bien les éléments.</p>

12

12

13



Indice de diversité de Gini pour un problème à deux classes

On peut en déduire une valuation de « être un attribut discriminant » :  
« Générer une partition équilibrée »

On peut en déduire une valuation de « être un attribut discriminant » :  
« Générer une partition équilibrée »

- Soit un attribut  $a$  à valeurs discrètes
- Pour toute valeur  $v$  prise par  $a$ , on note  $S_{a=v}$  l'ensemble des éléments de  $S$  pour lesquels  $a$  vaut  $v$ .

On peut en déduire une valuation de « être un attribut discriminant » :  
« Générer une partition équilibrée »

- Soit un attribut  $a$  à valeurs discrètes
- Pour toute valeur  $v$  prise par  $a$ , on note  $S_{a=v}$  l'ensemble des éléments de  $S$  pour lesquels  $a$  vaut  $v$ .

On définit lors le gain associé à  $a$  par

$$g(S, a) = Gini(S) - \sum_{v \in a} \frac{\#S_{a=v}}{\#S} Gini(S_{a=v})$$

- Un attribut est d'autant plus discriminant que son gain est élevé.
- On peut procéder de même avec  $H$  (ou un autre indice de diversité)

**Définition**  
On appelle *entropie* d'une partition  $S = \bigsqcup_{1 \leq i \leq n} c_i$

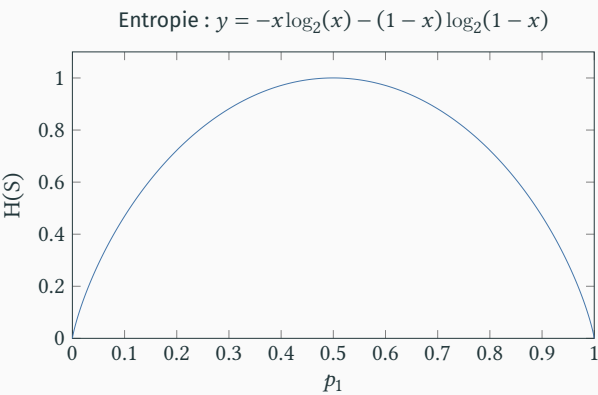
$$H(S) = - \sum_{1 \leq i \leq n} p_i \log_2(p_i)$$

avec  $p_i = \frac{\#c_i}{\sum_{j=1}^n \#c_j}$

Intuitivement, si on choisit au hasard et de façon uniforme un exemple  $x$  dans  $S$

- $p_i$  est la probabilité de l'évènement « La classe de  $x$  est  $c_i$  »
- $-\log_2 p_i$  mesure la surprise de l'évènement « La classe de  $x$  est  $c_i$  »

Finalement,  $H$  est donc la surprise moyenne de  $S$ .

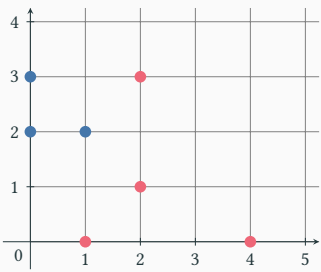


Entropie pour un problème à deux classes

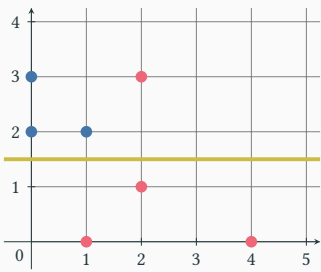
Pour des attributs numériques, on se ramène à un choix discret en utilisant des seuils

- Pour un attribut  $a$  à valeur numérique  $\in [\alpha, \beta]$  et  $r \in [\alpha, \beta]$ , on note  $a_r$  l'attribut booléen  $a(x) \leq r$
- On choisit  $s$  tel que  $g(S, a_r)$  soit maximal

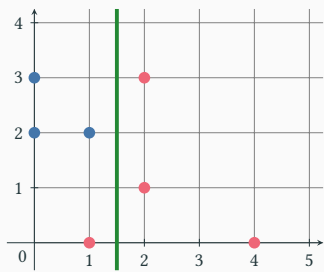
On peut ensuite utiliser  $a_r$  au lieu de  $a$  dans notre choix d'attributs.



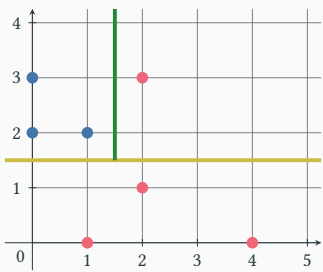
Exemple de classification par seuils



Exemple de classification par seuils



Exemple de classification par seuils



Exemple de classification par seuils

Propriétés

La principale qualité des arbres de décision tient dans leur simplicité

- *White box* : on peut comprendre le résultat
- Le modèle est de petite taille et est efficace même avec peu de données
- La procédure d'apprentissage est assez proche du raisonnement humain conscient
- Facilement combinables

Leur principal défaut est leur instabilité

→ des petites variations dans l'ensemble de test peuvent conduire à des changements importants dans l'arbre

Des extensions plus performantes, mais moins simples existent (*random forest...*)