

Apprentissage à partir
d'exemples

09/04/2015 1

Plan

- Introduction à l'apprentissage
- Apprentissage inductif: Arbres de décision
- Apprentissage d'ensembles d'hypothèses
- Apprentissage à partir d'instances : K -voisins les plus proches
- Réseaux de Neurones

09/04/2015 2

Apprentissage par observation

- Les perceptions ne devraient pas uniquement être utilisées pour choisir une action, mais aussi pour **améliorer la capacité de l'agent à agir dans le futur.**
- L'apprentissage est essentiel pour des environnements inconnus, où le concepteur ne peut pas tout savoir à l'avance.
- L'apprentissage modifie les mécanismes de décision de l'agent pour **améliorer sa performance.**

09/04/2015 3

Apprentissage des composantes de l'agent

- Un agent peut apprendre à l'aide d'un **retour approprié**
 - Ex: **conducteur de taxi**.
 - À chaque fois que l'instructeur crie « Freine! », l'agent peut apprendre une règle condition-action.
 - En montrant à l'agent plusieurs images de caméra contenant des autobus, alors l'agent peut apprendre à reconnaître des autobus.
 - Il peut apprendre les effets de ses actions en les essayant. Ex: freiner sec sur une chaussée enneigée.
 - Si l'agent ne reçoit pas de pourboire pour un trajet où les passagers ont été brassés, alors l'agent pourrait apprendre une partie de sa fonction d'utilité globale.

09/04/2015

4

Type de retour pour le critique

- Trois types d'apprentissage
 - **Apprentissage supervisé**
 - Dans ce cas, la base de données d'apprentissage est un ensemble de couples entrée-sortie (x_n, y_n) avec $x_n \in X$ et $y_n \in Y$, que l'on considère être tirées selon une loi sur $X \times Y$ inconnue, par exemple x_n suit une loi uniforme et $y_n = f(x_n) + w_n$ où w_n est un bruit centré.
 - **Apprentissage non supervisé**
 - L'agent apprend à partir des relations entre les perceptions. Il apprend à prédire les perceptions à partir de celles du passé.
 - **Apprentissage par renforcement**
 - L'évaluation de l'action est faite par récompense ou punition (attachée à la tâche).

09/04/2015

5

Apprentissage inductif

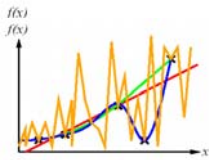
- Forme la plus simple d'apprentissage (**supervisé**) qui **consiste à apprendre une fonction à partir d'exemples**.
 - On veut apprendre une fonction f ;
 - On a des exemples de la forme $(x, f(x))$;
 - On veut trouver une fonction h (**hypothèse**) qui se rapproche le plus possible de f .
- C'est un modèle très simplifié d'apprentissage
 - Ignore les connaissances a priori;
 - Suppose un environnement observable et déterministe;
 - Suppose que des exemples sont donnés.

09/04/2015

6

Méthode d'apprentissage inductif

- h est **cohérente** si elle approxime f sur **tous les exemples**.
- Comment choisir alors parmi plusieurs hypothèses cohérentes?
 - Ex: Ajustement de courbes



- Laquelle choisir ?
- « Ockham's razor »: Choisir la plus simple parmi les hypothèses compatibles.

Si il n'y a pas d'hypothèse cohérente, il faut faire un compromis entre la cohérence et la simplicité.

09/04/2015

Espace des hypothèses

- L'**espace des hypothèses** \mathcal{H} est très important.
 - On doit choisir un espace d'hypothèses qui permet de représenter les données.
 - Par exemple, si on choisit des polynômes et que les données d'apprentissage proviennent d'une fonction sinusoïdale, ~~on~~ n'y arrivera pas.
 - Il ne faut pas prendre trop grand, car l'apprentissage serait trop lent.

\mathcal{H}

09/04/2015

Espace des hypothèses (suite)

- Problème d'apprentissage est dit **réalisable** : Si l'espace des hypothèses \mathcal{H} contient la vraie fonction.
- L'hypothèse la plus probable

$$h^* = \arg \max_{h \in \mathcal{H}} P(h|data)$$

$$h^* = \arg \max_{h \in \mathcal{H}} P(data|h)P(h)$$

- Dans ce cas, on peut favoriser via $P(h)$ les degré de polynôme 1, 2 ou autres.

09/04/2015

Apprentissage d'arbres de décision

- Une des formes les plus simples d'apprentissage, **mais tout de même une de celles qui connaissent le plus de succès.**
- À partir d'exemples, le but est d'apprendre des structures d'arbres permettant de prendre des décisions.
- Chaque nœud représente un test à faire.
- Chaque branche représente une valeur possible résultant du test.
- Une feuille correspond à une prise de décision.

09/04/2015

10

Exemple du restaurant

Faut-il ou pas attendre (*WillWait*) pour une table, dans un restaurant ?

Attributs

1. *Alternate*: whether there is a suitable alternative restaurant nearby.
2. *Bar*: whether the restaurant has a comfortable bar area to wait in.
3. *Fri/Sat*: true on Fridays and Saturdays.
4. *Hungry*: whether we are hungry.
5. *Patrons*: how many people are in the restaurant (values are *None*, *Some*, and *Full*).
6. *Price*: the restaurant's price range (\$, \$\$, \$\$\$).
7. *Raining*: whether it is raining outside.
8. *Reservation*: whether we made a reservation.
9. *Type*: the kind of restaurant (French, Italian, Thai, or burger).
10. *WaitEstimate*: the wait estimated by the host (0-10 minutes, 10-30, 30-60, or >60).

09/04/2015

11

Exemple du restaurant

- **Exemple:** En arrivant au restaurant, est-ce que ça vaut la peine d'attendre ?

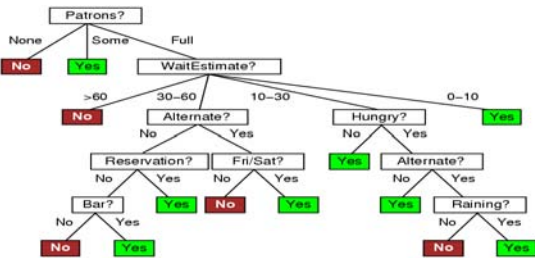
Example	Attributes										Goal <i>WillWait</i>
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	
<i>X</i> ₁	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0-10	Yes
<i>X</i> ₂	Yes	No	No	Yes	Full	\$	No	No	Thai	30-60	No
<i>X</i> ₃	No	Yes	No	No	Some	\$	No	No	Burger	0-10	Yes
<i>X</i> ₄	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10-30	Yes
<i>X</i> ₅	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	No
<i>X</i> ₆	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0-10	Yes
<i>X</i> ₇	No	Yes	No	No	None	\$	Yes	No	Burger	0-10	No
<i>X</i> ₈	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0-10	Yes
<i>X</i> ₉	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	No
<i>X</i> ₁₀	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10-30	No
<i>X</i> ₁₁	No	No	No	No	None	\$	No	No	Thai	0-10	No
<i>X</i> ₁₂	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30-60	Yes

09/04/2015

12

Exemple du restaurant

- **Exemple:** En arrivant au restaurant, est-ce que ça vaut la peine d'attendre ?



09/04/2015

13

Apprentissage d'arbres de décision

- Les arbres de décision sont des arbres qui permettent, la plupart du temps, de donner une réponse binaire du type oui ou non.
- Les instances sont représentés par des pairs (attribut, valeur). Les valeurs peuvent être discrètes ou continues.
- La fonction visée est généralement binaire, mais il est facilement possible d'étendre les arbres de décision pour qu'ils puissent avoir plusieurs valeurs en sorties.
- Ils pourraient donner en sortie des valeurs réelles, mais ce n'est pas beaucoup utilisé.

09/04/2015

14

Apprentissage d'arbres de décision

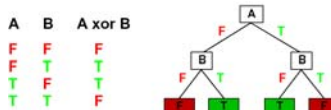
- Les exemples d'entraînement peuvent contenir des erreurs.
 - Les erreurs dans les exemples d'entraînement sont souvent appelés des **bruits**.
 - Les arbres de décision sont robustes aux erreurs de classification et aux erreurs dans les valeurs des attributs.
 - **Bien sûr, il ne faut pas qu'il y en ait trop.**
- Il peut manquer des valeurs pour certains attributs dans les exemples d'entraînement.
- Les arbres de décision sont bons pour les problèmes de classification.

09/04/2015

15

Expressivité

- Toutes les fonctions booléennes peuvent être représentées à l'aide d'un arbre de décision.
- On peut faire cela tout simplement en faisant correspondre chaque ligne de la table de vérité à un chemin dans l'arbre.
- Ceci va donner un arbre de grandeur exponentiel, mais on peut souvent faire beaucoup mieux.



09/04/2015

16

Expressivité

- Les arbres de décisions ne sont pas efficaces (grandeur exponentielle) pour représenter certaines fonctions, comme:
 - **Fonction de parité**: retourne 1 s'il y a un nombre pair de 1 dans les entrées.
 - **Fonction de majorité**: retourne 1 s'il y a plus de la moitié des entrées à 1.
- Combien d'arbres possibles pour n attributs booléens ?
 - Chaque réponse colonne peut être vu comme un nombre de 2^n
 - le nombre de table de vérité distinctes de 2^n rangées = 2^{2^n}

Avec 10 attributs booléens du restaurant, on a 2^{1024} soit environ 10^{308} fonctions à choisir !

09/04/2015

17

Construire un arbre de décision

- L'apprentissage d'un arbre de décision est fait à partir d'exemples de valeurs d'attributs et de la valeur résultante du prédicat à apprendre.
- La valeur du prédicat à apprendre est appelée la **classification** de l'exemple (ex: Vrai / Faux).
- L'ensemble des exemples est appelé l'**ensemble d'entraînement**.

09/04/2015

18

Construire un arbre de décision

- Le but est de trouver le plus petit arbre qui respecte l'ensemble d'entraînement.
- Il ne s'agit pas uniquement de mémoriser les observations: **il faut trouver un arbre qui est capable d'extrapoler pour les exemples qu'il n'a pas déjà vus.**
- L'arbre doit extraire des tendances ou des comportements à partir des exemples.

09/04/2015

19

Algorithme

- Il construit les arbres de décision de haut en bas.
- **Il place à la racine l'attribut le plus important, c'est-à-dire celui qui sépare le mieux les exemples positifs et négatifs.**
- Par la suite, il y a un nouveau nœud pour chacune des valeurs possibles de cet attribut.
- Pour chacun de ces nœuds, on recommence le test avec le sous-ensemble des exemples d'entraînement qui ont été classés dans ce nœud.

09/04/2015

20

Algorithme

```
function DECISION-TREE-LEARNING(examples, attributes, default) returns a decision tree
  inputs: examples, set of examples
         attributes, set of attributes
         default, default value for the goal predicate

  if examples is empty then return default
  else if all examples have the same classification then return the classification
  else if attributes is empty then return MAJORITY-VALUE(examples)
  else
    best ← CHOOSE-ATTRIBUTE(attributes, examples)
    tree ← a new decision tree with root test best
    m ← MAJORITY-VALUE(examplesi)
    for each value vi of best do
      examplesi ← {elements of examples with best = vi}
      subtree ← DECISION-TREE-LEARNING(examplesi, attributes - best, m)
      add a branch to tree with label vi and subtree subtree
    end
  return tree
```

09/04/2015

21

Entropie pour le choix des attributs (Wiki)

- L'**entropie de Shannon**, due à [Claude Shannon](#), est une fonction mathématique qui, intuitivement, correspond à la quantité d'**information** contenue ou délivrée par une source d'information.
- Cette source peut être un texte écrit dans une langue donnée, un **signal électrique** ou encore un **fichier informatique** quelconque (collection d'octets).

09/04/2015

22

Entropie pour le choix des attributs (Wiki)

Pour une source (qui peut être une **variable aléatoire** discrète), X comportant n symboles, chacun ayant une probabilité P_i d'apparaître, l'**entropie** H de la source X est définie comme

$$H_b(X) = -\mathbb{E}[\log_b P(X = x_i)] = -\sum_{i=1}^n P_i \log_b \left(\frac{1}{P_i} \right) = -\sum_{i=1}^n P_i \log_b P_i.$$

où \mathbb{E} désigne l'**espérance mathématique**.

On utilise en général un logarithme à base 2 car l'entropie possède alors les unités de bits/symbole. Les symboles représentent les réalisations possibles de la variable aléatoire X .

$$H(X) = H_2(X) = -\sum_{i=1}^n P_i \log_2 P_i.$$

09/04/2015

23

Choix de l'attribut

- On choisit l'attribut ayant le meilleur gain d'information:

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

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

S : les exemples d'entraînement.

A : l'attribut à tester.

$V(A)$: les valeurs possibles de l'attribut A .

S_v : le sous-ensemble de S qui contient les exemples qui ont la valeur v pour l'attribut A .

c : le nombre de valeurs possibles pour la fonction visée.

p_i : la proportion des exemples dans S qui ont i comme valeur pour la fonction visée.

09/04/2015

24

Exemple

- Le calcul du gain d'information pour l'attribut **Ciel** va donc donner:

$$\begin{aligned}
 \text{Gain}(S, \text{Ciel}) &= \text{Entropie}(S) - \sum_{v \in V(\text{Ciel})} \frac{|S_v|}{|S|} \text{Entropie}(S_v) \\
 &= 0.94 - ((5/14) \times \text{Entropie}(S_{\text{ensoleillé}}) + \\
 &\quad (4/14) \times \text{Entropie}(S_{\text{nuageux}}) + \\
 &\quad (5/14) \times \text{Entropie}(S_{\text{pluvieux}})) \\
 &= 0.94 - ((5/14) \times 0.971 + (4/14) \times 0 + (5/14) \times 0.971) \\
 &= 0.94 - 0.694 \\
 &= 0.246
 \end{aligned}$$

09/04/2015

26

Exemple

- On calcul le gain de la même manière pour les trois autres attributs:

$$\begin{aligned}
 \text{Gain}(S, \text{Ciel}) &= 0.246 \\
 \text{Gain}(S, \text{Humidité}) &= 0.151 \\
 \text{Gain}(S, \text{Vent}) &= 0.048 \\
 \text{Gain}(S, \text{Température}) &= 0.029
 \end{aligned}$$

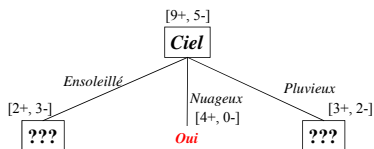
- L'attribut qui a la plus grand gain d'information est l'attribut **Ciel**, donc se sera la racine de l'arbre de décision.

09/04/2015

27

Exemple

- En séparant les exemples selon les valeurs de l'attribut **Ciel**, on obtient l'arbre partiel:



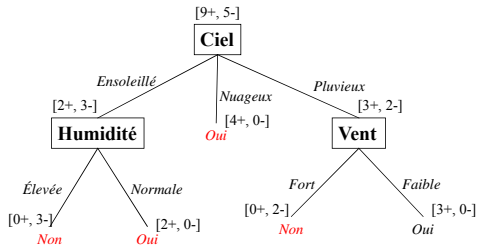
- On peut voir que lorsque le ciel est nuageux, il reste uniquement des exemples positifs, donc ce nœud devient une feuille avec une valeur de **Oui** pour la fonction visée.
- Pour les deux autres nœuds, il y a encore des exemples positifs et négatifs, alors il faut recommencer le même calcul du gain d'information, mais avec les sous-ensembles restants.

09/04/2015

28

Exemple

- En effectuant les calculs restants, on obtient:



09/04/2015

31

Procédure générale d'apprentissage

- Faire la collecte d'un grand ensemble d'exemples.
- Diviser les exemples en deux ensembles: un d'entraînement et l'autre de test (2/3 et 1/3).
- Utiliser l'ensemble d'entraînement comme exemples et générer l'hypothèse h .
- Mesurer le pourcentage d'exemples de l'ensemble de test qui sont correctement identifiés par h .
- Répéter les étapes 1 à 4 pour différentes tailles d'ensembles

09/04/2015

32

Gérer les erreurs dans les exemples

- Si les exemples d'entraînement contiennent des erreurs, il se peut qu'on trouve un arbre plus grand, mais qui ne classifie pas bien les instances autres que celles utilisées dans les exemples d'entraînement (**surapprentissage**).
- Pour vérifier notre arbre, on utilise un **ensemble de validation** (ou de test).
 - L'ensemble de validation est généralement le tiers des exemples disponibles pour l'apprentissage.

09/04/2015

33

Élagage de l'arbre

- Cette méthode coupe des parties de l'arbre en choisissant un nœud et en enlevant tout son sous-arbre.
 - Ceci fait donc du nœud une feuille et on lui attribut la valeur de classification qui revient le plus souvent.
- Des nœuds sont enlevés seulement si l'arbre résultant n'est pas pire que l'arbre initial sur les exemples de validation.
- On continue tant que l'arbre résultant offre de meilleurs résultats sur les exemples de validation.
- Ceci a pour but de réduire l'arbre en enlevant des branches qui auraient été ajoutées par une erreur dans les exemples.

09/04/2015

34

Exemple

- Est-ce une bonne journée pour jouer au tennis ?

Journée	Ciel	Température	Humidité	Vent	JouerTennis
J1	Ensoleillé	Chaude	Élevée	Faible	Non
J2	Ensoleillé	Chaude	Élevée	Fort	Non
J3	Nuageux	Chaude	Élevée	Faible	Oui
J4	Pluvieux	Tempérée	Élevée	Faible	Oui
J5	Pluvieux	Froide	Normal	Faible	Oui
J6	Pluvieux	Froide	Normal	Fort	Non
J7	Nuageux	Froide	Normal	Fort	Oui
J8	Ensoleillé	Tempérée	Élevée	Faible	Non
J9	Ensoleillé	Froide	Normal	Faible	Oui
J10	Pluvieux	Tempérée	Normal	Faible	Oui
J11	Ensoleillé	Tempérée	Normal	Fort	Oui
J12	Nuageux	Tempérée	Élevée	Fort	Oui
J13	Nuageux	Chaude	Normal	Faible	Oui
J14	Pluvieux	Tempérée	Élevée	Fort	Non

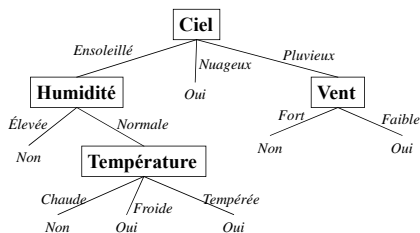
09/04/2015

35

Élagage de l'arbre

- Si on ajoute une erreur dans l'exemple précédent en modifiant le premier exemple, on obtient l'arbre suivant:

J1	Ensoleillé	Chaude	Élevée	Faible	Non
J1	Ensoleillé	Chaude	Normale	Faible	Non



09/04/2015

36

Élagage de l'arbre

- L'élagage pour cet arbre consiste à enlever le nœud *Température* qui vient d'être ajouté avec l'erreur.
- Par la suite, on teste le nouvel arbre élagué sur les exemples de validation.
- Si la performance n'est pas diminuée, alors on conserve l'arbre élagué.

09/04/2015

37

Valeurs d'attributs manquantes

- S'il y a des valeurs pour certains attributs qui ne sont pas disponibles, alors on peut :
 - Donner la valeur moyenne pour cet attribut.
 - On regarde les autres exemples et on calcule la moyenne des valeurs présentes.
 - On utilise cette moyenne pour estimer la valeur manquante dans le calcul du gain d'information.
 - Attribuer une probabilité pour la valeur manquante.
 - On regarde les autres exemples et on calcule la probabilité de chaque valeur possible pour l'attribut.
 - On utilise par la suite ces probabilités pour calculer le gain d'information.

09/04/2015

38

Valeurs d'attributs manquantes

- Pour la première stratégie, les calculs ne changent pas. On ne fait qu'utiliser la valeur moyenne pour remplacer la valeur manquante.
- Pour la deuxième stratégie, les calculs sont modifiés pour utiliser les probabilités.
 - On commence par calculer la probabilité de chacune des valeurs possibles pour l'attribut manquant.
 - Par exemple, supposons qu'il nous manquerait une valeur pour l'attribut *Vent*.
 - La probabilité que le vent soit faible, selon nos exemples d'entraînement est de $8/14 * 100 = 57\%$.
 - La probabilité que le vent soit fort est de $6/14 * 100 = 43\%$.
 - Ceci nous donne des fractions d'exemples dans nos calculs.

09/04/2015

39

Exemple

- Supposons que l'on a un exemple d'entraînement positif de plus, mais que la valeur pour l'attribut *Vent* est absente.

$$\begin{aligned} \text{Entropie}(S) &= \sum_{i=1}^c -p_i \log_2 p_i \\ &= -10/15 \log_2 10/15 + -5/15 \log_2 5/15 \\ &= 0.918 \end{aligned}$$

$$\text{Entropie}(S_{\text{fort}}) = (-4/6.43) \log_2 4/6.43 + (-3/6.43) \log_2 3/6.43 = 0.939$$

$$\text{Entropie}(S_{\text{faible}}) = (-7/8.57) \log_2 7/8.57 + (-2/8.57) \log_2 2/8.57 = 0.728$$

$$\begin{aligned} \text{Gain}(S, \text{Vent}) &= \text{Entropie}(S) - \sum_{v \in V(\text{Vent})} \frac{|S_v|}{|S|} \text{Entropie}(S_v) \\ &= 0.918 - ((6.43/15) \times \text{Entropie}(S_{\text{fort}}) + \\ &\quad (8.57/15) \times \text{Entropie}(S_{\text{faible}})) \\ &= 0.918 - ((6.43/15) \times 0.939 + (8.57/15) \times 0.728) \\ &= 0.18 - 0.818 \\ &= 0.1 \end{aligned}$$

09/04/2015

40

Attributs multivalués

- Il y a problème avec la fonction de gain d'information.
 - Lorsque les attributs ont beaucoup de valeurs possibles, comme par exemple un attribut date, leur gain est très élevé, car il classifie parfaitement les exemples.
 - Par contre, ils vont générer un arbre de décision d'une profondeur de 1 qui ne sera pas très bon pour les instances futures.
- Solution: on peut utiliser une fonction qui se nomme *GainRatio* qui pénalise les attributs qui ont trop de valeurs possibles.

09/04/2015

41

Attributs multivalués

- On utilise la fonction *GainRatio* au lieu de *Gain* lorsque l'on a des attributs multivalués.

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

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

S_i à S_c sont les c ensembles résultant de la partition de S par les c valeurs de l'attribut A .

09/04/2015

42

Attributs à valeurs continues

- On utilise un **point de coupe** pour obtenir une discrétisation des variables continues.
 - Ex: la variable *Température* est continue et on a les 6 exemples suivants.

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

- On met les valeurs en ordre croissant et on regarde les endroits où la fonction change de valeur.
- À ces endroits, on choisit la médiane comme valeur de coupe.
- On compare toutes les valeurs de coupe et on choisit celle qui apporte le plus grand gain d'information.

09/04/2015

45

Exemple

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

- Il y a deux valeurs de coupe:
 - $(48 + 60)/2 = 54$
 - $(80 + 90)/2 = 85$
- On calcule les gains d'information pour ces deux valeurs de coupes.

$$\begin{aligned}
 Entropie(S) &= \sum_{i=1}^c -p_i \log_2 p_i \\
 &= -3/6 \log_2 3/6 + -3/6 \log_2 3/6 \\
 &= 1
 \end{aligned}$$

09/04/2015

46

Exemple

- Gain pour la coupe 54:

$$\begin{aligned}
 Gain(S, T_{54}) &= Entropie(S) - \sum_{v \in V(T_{54})} \frac{|S_v|}{|S|} Entropie(S_v) \\
 &= 1 - (((2/6) \times (-2/2 \log_2 2/2 + -0/2 \log_2 0/2)) + \\
 &\quad ((4/6) \times (-3/4 \log_2 3/4 + -1/4 \log_2 1/4))) \\
 &= 1 - (0 + 0.541) \\
 &= 0.459
 \end{aligned}$$

- Gain pour la coupe 85:

$$\begin{aligned}
 Gain(S, T_{85}) &= Entropie(S) - \sum_{v \in V(T_{85})} \frac{|S_v|}{|S|} Entropie(S_v) \\
 &= 1 - (((5/6) \times (-2/5 \log_2 2/5 + -3/5 \log_2 3/5)) + \\
 &\quad ((1/6) \times (-1/1 \log_2 1/1 + -0/1 \log_2 0/1))) \\
 &= 1 - (0.809 + 0) \\
 &= 0.191
 \end{aligned}$$

09/04/2015

45

Exemple

- En comparant les gains d'informations, on se rend compte que 54 est la meilleure valeur de coupe.
 - Par conséquent, la température sera utilisée comme un attribut binaire qui ne pourra prendre que deux valeurs: plus petit que 54 ou plus grand ou égale à 54.

09/04/2015

46

Exemples d'applications des arbres de décision

- Conception d'équipement de plate-forme de forage :
 - Un système expert appelé GASOIL (BP) pour la conception de système de séparation d'huile sur des plates-formes en mer.
 - La séparation dépend (entre autres) de : proportion relative de gaz, d'huile et d'eau, débit, pression, densité, viscosité, température, etc.
 - À ce moment, le plus gros système expert au monde: 2500 règles et un estimé de 10 personnes-années pour le construire.
 - Utilisation de l'apprentissage d'arbres de décision à partir d'une base de données d'exemples de conception existante: construit en 100 personnes-jours.
 - Plus performant que les humains et a sauvé des millions.

09/04/2015

47

Exemples d'applications des arbres de décision

- Pilote d'avion: Une alternative aux méthodes formelles pour la construction d'un contrôleur pour des systèmes complexes est d'apprendre à poser les bonnes actions dans chaque situation.
- Pour des systèmes complexes, comme les avions, il peut être impossible de prendre l'approche formelle. Il reste l'autre approche...
 - 3 pilotes expérimentés avec plan de vol, répétés 30 fois: acquisition de données action-situation (90 000 exemples) avec 20 variables.
 - Apprentissage d'un l'arbre de décision.
 - Conversion en C et insertion dans la boucle de contrôle.
 - Résultat: Apprend à voler et mieux que les pilotes humains, car les erreurs humaines ont été éliminées par le mécanisme de généralisation.

09/04/2015

48

Apprentissage d'ensembles d'hypothèses

- Au lieu d'apprendre une seule hypothèse, on en apprend plusieurs et on combine leur prédiction pour retourner la réponse.
 - Ex: si on a 5 arbres de décisions
 - Pour évaluer une instance, on exécute les 5 arbres
 - On retourne la réponse qui revient le plus souvent:
Vote par majorité
 - Ex: (2 vrais, 5 faux), donc on retourne faux.

09/04/2015

49

La méthode « boosting »

- Cette méthode est utilisée pour générer un ensemble d'hypothèses.
- Attribue un poids à tous les exemples.
- Plus un exemple a un grand poids, plus il sera considéré important par l'algorithme d'apprentissage.

09/04/2015

50

La méthode « boosting »

- L'algorithme fonctionne de la manière suivante:
 - Donner un poids de 1 à tous les exemples.**
 - Répéter M fois**
 - Générer une hypothèse à partir des exemples
 - Augmenter le poids des exemples mal classifiés.
 - Diminuer le poids des exemples bien classifiés.
- L'hypothèse finale est la **majorité pondérée des M hypothèses** où le poids d'une hypothèse dépend de sa performance sur l'ensemble d'entraînement.

09/04/2015

51

Apprentissage à base d'instances

09/04/2015

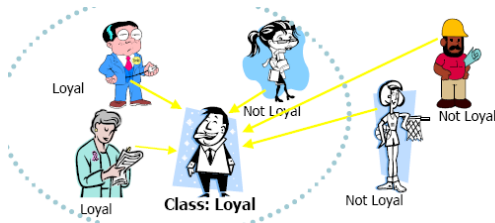
Algorithme

- **Paramètre** : le nombre k de voisins
- **Donnée** : un échantillon de m exemples et leurs classes
 - La classe d'un exemple X est $c(X)$
- **Entrée** : un enregistrement Y
- 1. Déterminer les k plus proches exemples de Y en calculant les distances
- 2. Combiner les classes de ces k exemples en une classe c
- **Sortie** : la classe de Y est $c(Y) = c$

53

Exemple: Client loyal ou non

$K = 3$



54

Autres distances

Summation :

$$D(X,Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Distance euclidienne ponderée :

$$D(X,Y) = \sqrt{\sum_{i=1}^n w_i (x_i - y_i)^2}$$

55

Exemple (1)

Customer	Age	Income	No. credit cards	Loyal
John	35	35K	3	No
Rachel	22	50K	2	Yes
Hannah	63	200K	1	No
Tom	59	170K	1	No
Nellie	25	40K	4	Yes
David	37	50K	2	?

56

Exemple (2)

K = 3

Customer	Age	Income	No. credit cards	Loyal	Distance from David
John	35	35K	3	No	$\sqrt{[(35-37)^2 + (35-50)^2 + (3-2)^2]} = 15.16$
Rachel	22	50K	2	Yes	$\sqrt{[(22-37)^2 + (50-50)^2 + (2-2)^2]} = 15$
Hannah	63	200K	1	No	$\sqrt{[(63-37)^2 + (200-50)^2 + (1-2)^2]} = 152.23$
Tom	59	170K	1	No	$\sqrt{[(59-37)^2 + (170-50)^2 + (1-2)^2]} = 122$
Nellie	25	40K	4	Yes	$\sqrt{[(25-37)^2 + (40-50)^2 + (4-2)^2]} = 15.74$
David	37	50K	2	Yes	

57

Apprentissage à base d'instances

- Ne construit pas de fonctions à l'aide des exemples d'entraînement.
- Ne fait qu'emmagasiner les exemples.
- Le calcul n'est fait que lorsque le programme a une nouvelle instance à classer.
- **Avantage:** estimation locale pour chaque instance à classer.
- **Désavantage:** coût de classification élevé.

09/04/2015

k -voisins les plus proches

- Les instances sont des points dans un espace à n -dimensions où n est le nombre d'attributs.
- On utilise la distance Euclidienne pour déterminer les voisins les plus proches.
- Une instance x est définie par son vecteur d'attributs:
 $\langle a_1(x), a_2(x), \dots, a_n(x) \rangle$
- La distance entre deux instances x_i et x_j est :

$$d(x_i, x_j) \equiv \sqrt{\sum_{r=1}^n (a_r(x_i) - a_r(x_j))^2}$$

09/04/2015

Fonction visée à valeurs discrètes

- L'ensemble V contient les valeurs possibles.
- L'estimation de la fonction f est tout simplement la valeur qui revient le plus souvent dans les k voisins.

$$\hat{f}(x_q) \leftarrow \max_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$$

où

$$\delta(a, b) = \begin{cases} 1 & \text{si } a = b \\ 0 & \text{sinon} \end{cases}$$

09/04/2015

k-voisins les plus proches avec distances pondérées

- Une amélioration à l'algorithme des *k*-voisins le plus proche est de **pondérer la contribution** de chacun des *k* voisins par leur distance par rapport à l'instance à classer.
- **Les K-voisins sont robustes au bruit**
- Pour les valeurs discrètes, on obtient:

$$\hat{f}(x_q) \leftarrow \max_{v \in V} \sum_{i=1}^k w_i \delta(v, f(x_i)) \quad \text{où} \quad w_i \equiv \frac{1}{d(x_q, x_i)^2}$$

09/04/2015

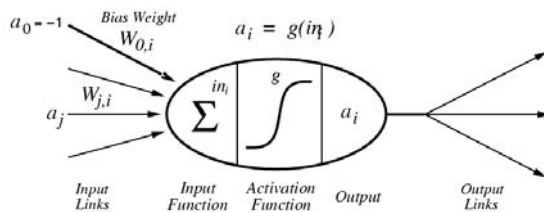
Réseaux de Neurones

09/04/2015

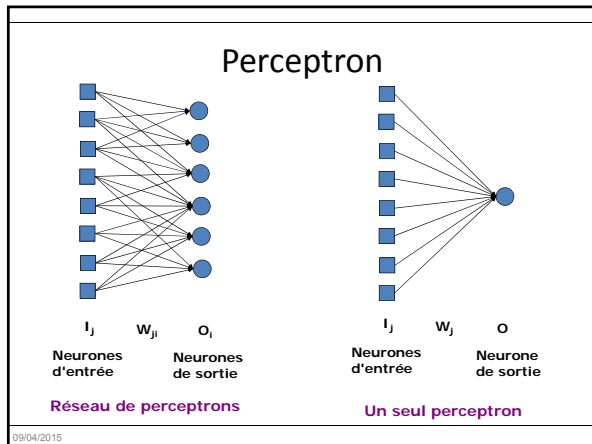
102

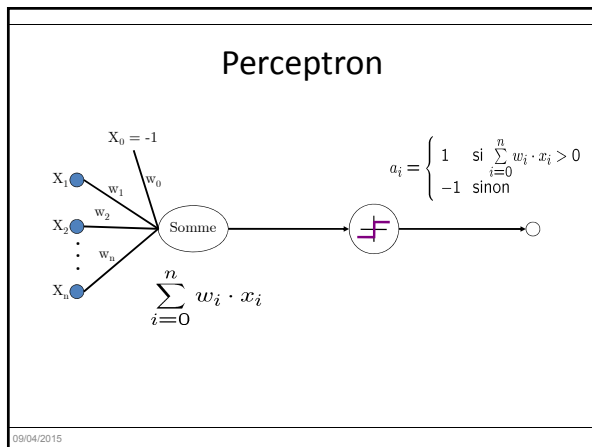
Réseau de neurones

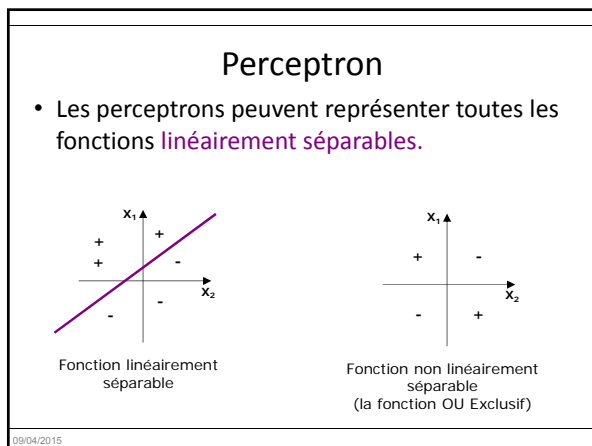
- Plusieurs neurones interconnectés ensemble
- Chaque neurone a la forme suivante:



09/04/2015







Exemple de perceptron

- La fonction ET
 - Les entrées prennent les valeurs 1 (vrai) ou 0 (faux)

$X_0 = -1$

$$\sum_{i=0}^n w_i \cdot x_i$$

$$a_i = \begin{cases} 1 & \text{si } \sum_{i=0}^n w_i \cdot x_i > 0 \\ -1 & \text{sinon} \end{cases}$$

Pour un OU, il suffirait de poser $w_0 = 0.3$

09/04/2015

Règle d'apprentissage du perceptron

- Commencer avec des poids de valeurs aléatoires.
- Traiter tous les exemples d'entraînement jusqu'à ce que le perceptron les classe tous correctement.
- À chaque exemple, les poids sont révisés à l'aide de la règle suivante:
 - $w_i \leftarrow w_i + \Delta w_i$
 - $\Delta w_i = \eta(t - o)x_i$

t : sortie désirée
o : sortie obtenue
η : constante d'apprentissage, généralement petite (ex. 0.1)

09/04/2015

Règle d'apprentissage du perceptron

$w_i \leftarrow w_i + \Delta w_i \quad \Delta w_i = \eta(t - o)x_i$

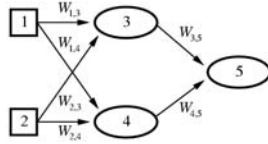
- Pourquoi une telle règle converge-t-elle vers une bonne hypothèse ?

Sortie perceptron	$(t - o)$	Δw_i	w_i
Bonne valeur	0	0	Pareil
-1 au lieu de +1	+	+	Augmente
+1 au lieu de -1	-	-	Diminue

09/04/2015

Ajustement des poids

- L'ajustement des poids dans les réseaux de neurones revient à modifier les paramètres de la fonction représentée par le réseau.
- Par exemple, la sortie du réseau a_5 dépend de la valeur des entrées a_1 et a_2 . Les paramètres de la fonction sont les poids W .



$$a_5 = g(W_{3,5}a_3 + W_{4,5}a_4)$$

$$= g(W_{3,5}g(W_{1,3}a_1 + W_{2,3}a_2) + W_{4,5}g(W_{1,4}a_1 + W_{2,4}a_2))$$

09/04/2015

Règle Delta

- Permet de gérer les exemples qui ne sont pas linéairement séparables
 - Converge vers l'estimation la plus juste possible
 - Utilise l'algorithme de la descente du gradient
 - La règle utilisée pour ajuster les poids est:

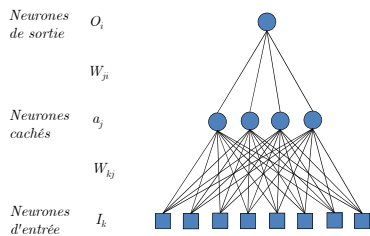
$$w_i \leftarrow w_i + \Delta w_i$$

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{id}$$

D : ensemble de tous les exemples d'entraînement

09/04/2015

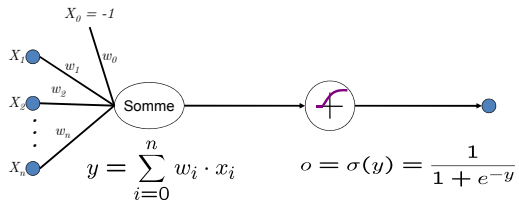
Réseau à plusieurs couches



09/04/2015

Sigmoïd

- Fonction continue entre 0 et 1



09/04/2015

Algorithme de rétropropagation

- Utilise l'algorithme de descente du gradient.
- Tente de diminuer la différence au carré entre les sorties du réseau et les sorties désirées.
- L'erreur E du réseau est:

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{sorties}} (t_{kd} - o_{kd})^2$$

t_{kd} : sortie désirée
 o_{kd} : sortie obtenue
 D : ensemble des exemples d'entraînement
 sorties : ensemble des sorties du réseau

09/04/2015

Algorithme de rétropropagation Texte du Pr. Marc Parizeau

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{sorties}} (t_{kd} - o_{kd})^2$$

On exprime la variation de poids par

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial w_{ji}}$$

[Texte du Prof Marc Parizeau](#)

09/04/2015

75

Exemple

$\eta = 0.05$ $o_4 = 0.6457$
 $o_5 = 0.6442$

Entrées

Sortie

$\delta_5 = o_5(1-o_5)(t_5-o_5) = 0.6442(1-0.6442)(2-0.6442) = 0.3108$
 $\delta_4 = o_4(1-o_4) \sum_{k \in \{5\}} w_{k4} \delta_k = 0.6457(1-0.6457)(0.3 \times 0.3108) = 0.0213$

09/04/2015

Exemple

$\eta = 0.05$ $\delta_4 = 0.0213 o_4 = 0.6457$
 $\delta_5 = 0.3108 o_5 = 0.6442$

Entrées

Sortie

$\Delta w_{40} = \eta \delta_4 x_{40} = 0.05 \times 0.0213 \times -1 = -0.0011$
 $\Delta w_{41} = \eta \delta_4 x_{41} = 0.05 \times 0.0213 \times 2 = 0.0021$
 $\Delta w_{42} = \eta \delta_4 x_{42} = 0.05 \times 0.0213 \times 5 = 0.0053$
 $\Delta w_{53} = \eta \delta_5 x_{53} = 0.05 \times 0.3108 \times -1 = -0.0155$
 $\Delta w_{54} = \eta \delta_5 x_{54} = 0.05 \times 0.3108 \times 0.6457 = 0.0100$

09/04/2015

Exemple

$\eta = 0.05$ $\Delta w_{40} = -0.0011$
 $\Delta w_{41} = 0.0021$
 $\Delta w_{42} = 0.0053$
 $\Delta w_{53} = -0.0155$
 $\Delta w_{54} = 0.0100$

Entrées

Sortie

$w_{40} = w_{40} + \Delta w_{40} = -0.5 + -0.0011 = -0.5011$
 $w_{41} = w_{41} + \Delta w_{41} = -0.2 + 0.0021 = -0.1979$
 $w_{42} = w_{42} + \Delta w_{42} = 0.1 + 0.0053 = 0.1053$
 $w_{53} = w_{53} + \Delta w_{53} = -0.4 + -0.0155 = -0.4155$
 $w_{54} = w_{54} + \Delta w_{54} = 0.3 + 0.0100 = 0.3100$

09/04/2015

Condition d'arrêt

- Le nombre d'itérations est important car:
 - Si trop faible, l'erreur n'est pas suffisamment réduite.
 - Si trop grand, le réseau devient trop spécifique aux données d'entraînement.
- Il y a plusieurs conditions d'arrêt possible:
 - Après un certain nombre fixe d'itérations.
 - Lorsque l'erreur dans les sorties des exemples d'entraînement descend en dessous d'une certaine borne.
 - Lorsque l'erreur avec les exemples de validation descend en dessous d'une certaine borne.

09/04/2015

Convergence

- La surface qui représente l'erreur contient plusieurs minimums locaux.
- L'algorithme ne garantit pas l'atteinte du minimum global.
- Mais, tout de même **bon en pratique**
 - Poids initialiser à de petites valeurs.
 - La sigmoïd est presque linéaire pour de petites valeurs, donc pas de minimums locaux.
 - Les poids grandissent à une vitesse contrôlée par le coefficient d'apprentissage.
 - Donc, les poids ont le temps de se rapprocher du minimum global avant que la fonction devienne plus complexe avec plusieurs minimums locaux.

09/04/2015

Caractéristiques

- Instances représentées par des paires attribut-valeur. Les attributs peuvent être indépendants ou non. Les valeurs peuvent être réelles.
- Les sorties peuvent être des valeurs discrètes, continues ou vectorielles.
- Les exemples peuvent contenir des erreurs.
- Le temps d'apprentissage peut être long.
- L'évaluation de la fonction apprise est très rapide.
- La fonction visée est souvent très peu compréhensible par un humain.

09/04/2015

Construction du réseau

- Il existe aucune techniques précise
 - La construction se fait souvent par essaies et erreurs
- Il y a quelques algorithmes:
 - Algorithmes génétiques
 - Demande beaucoup de temps de calcul
 - « Optimal brain damage »
 - Commence avec un réseau complètement connecté et on enlève les poids qui sont proches de zéro. On peut aussi enlever les neurones qui ne contribuent pas au résultat.
 - « Tiling algorithm »
 - Commence avec un neurone et ajoute des neurones jusqu'à ce qu'ils couvrent tous les exemples d'entraînement.

09/04/2015

Applications

- L'interprétation d'images.
- La reconnaissance vocale.
- La reconnaissance de mots écrits à la main.
<http://members.aol.com/Trane64/java/JRec.html>
- L'apprentissage de stratégies de contrôle pour les robots.
- Une des meilleurs méthodes connues pour l'interprétation de données provenant de capteurs dans le monde réel.

09/04/2015

Applications

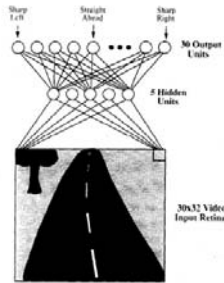
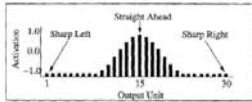
- Le système ALVINN a conduit une voiture à 55 milles/h pendant 90 miles sur une autoroute.



09/04/2015

Applications

- L'architecture d'ALVINN est un réseau de neurones à rétropropagation avec une couche cachée.
- La sortie est une distribution de Gauss centrée autour de la direction appropriée.



09/04/2015
