

Apprentissage d'ensembles de règles

C. Rouveirol
MASTER MICR

Sommaire

- Approche générer et tester descendante
 - Exemple : ID3 (déjà vu au cours 1 et 2)
 - Apprentissage de règles (DNF)
 - Apprentissage de règles par couverture
 - Algorithme de couverture
 - Boucle externe : Algorithme générique
 - Boucle interne : Algorithme générique
 - Exemples détaillés d'algorithmes : CN2 (Clark 90)
- Algorithme ascendant « guidé par les données »
- Algorithme descendant « guidé par les données »

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

2

Sources

- Transparents de J. Fürnkranz sur Inductive Rule Learning (article Roc'n'Rule)
- Transparents A. Osmani, Apprentissage Symbolique, M2 MICR 2006-2007
- Transparents E. Alphonse, M2 EID, 2006-2007

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

3

Apprentissage d'ensembles de règles (DNF)

Pourquoi des ensembles de règles ?

- arbres de décision souvent complexes et difficiles à interpréter
- problème de réplification de sous-arbres avec arbres de décision, ce qui partitionne inutilement l'espace des instances et pose des difficultés à la recherche heuristique
- les ensembles de règles ordonnés à k littéraux par règles plus expressif que les arbres de profondeur k
- se généralise bien à la logique d'ordre un

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

4

Comment les apprendre ?

- Apprendre directement l'ensemble de règles
 - intérêt théorique mais très peu utilisé
- Apprendre une règle à la fois : stratégie de couverture

L'origine de la stratégie de couverture est due à Michalski (1969) (Algorithme AQ). Le terme Separate-and-conquer dû à Pagallo et Haussler (1990) est simplement une interprétation de la stratégie de couverture.

- **Couverture séquentielle ou simultanée?**
 - Séquentielle: apprend une hypothèse à la fois
 - Simultanée: apprend un arbre (plusieurs règles en parallèle)
- **Ascendant ou descendant?**
- **Générer et tester ou guidé par les données/exemples ?**
 - Générer et tester: explore un espace d'hypothèses bien formées (ordonnées)
 - Guidé par les données: utilise les exemples pour modifier son hypothèse courante
- **Post-pruning des règles**
- **Quelles heuristiques ?**
 - Entropie, prédiction, laplacien...

Recherche descendante dans l'espace des hypothèses

Deux types d'opérateurs :

- «générer-tester» : fondés sur la structure de H seulement

$$\rho(h) = h' \text{ avec } h' \in H \text{ et } h' \leq_h h$$
- «dirigés par les données» : utilisation d'un exemple négatif pour spécialiser

$$\rho(h, e^-) = h' \text{ avec } h' \in H \text{ et } h' \leq_h h \text{ et } \text{not couvre}(h', e^-)$$
- On appelle *couverture inférieure* de h l'ensemble d'hypothèses:

$$CI(h) = \{h' \in H \mid h' \leq_h h \text{ et il n'existe pas de } h'' \in H \text{ t.q. } h' \leq_h h'' \leq_h h\}$$

Apprentissage de règles par couverture

Tous les algorithmes d'apprentissage de règles dits « *separate-and-conquer* » ou « *covering* » partagent le même principe :

- Chercher une règle qui couvre une partie des exemples positifs
- Enlever les exemples positifs couverts de la base d'apprentissage initiale
- Recommencer récursivement le processus jusqu'à ce qu'il n'y ait plus d'exemples positifs à couvrir

A la fin de ce processus, chaque exemple de la base d'apprentissage est couvert par au moins une règle

Algorithme générique de couverture (boucle externe)

Procédure Couverture(Exemples) % renvoie une hypothèse cohérente

Hypothèse = T % élément universel

Tantque Positif(Exemples) $\neq \emptyset$

Règle = **ApprendreMeilleureRègle**(Exemples)

Hypothèse = Hypothèse \cup Règle

Exemples = Exemples \setminus ExemplesCouverts(Règle, Exemples)

Renvoyer Hypothèse

Les algorithmes d'apprentissage par couverture diffèrent au niveau de la méthode implantée pour **ApprendreMeilleureRègle** (Exemples).

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

9

Apprentissage de règles par couverture

Beaucoup d'algorithmes de cette famille sont développés pour les différents langages d'hypothèses. Parmi lesquels :

DNF : AQ et AQ15 (Michalski69 et 86), PRISM (Cendrowska87), PFOIL (Mooney95), JoJo(Fensel&Wiese93), **RIPPER** (Cohen95), ...

CNF : PFOIL-CNF (Mooney1995), ICL (de Raedt95)

Decisionlists : CN2 (Clark & Niblett90), PREPEND (Webb&Brkie93), FOILDL (Mooney95) ...

Programmes logiques : INDUCE (Michalski80), FOIL(Quinlan90), PROGOL (Muggleton95), FOCL(Pazzani92)

Règles de régression : RULE (Weiss93), FORS (Karalie95) ...

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

10

Apprentissage d'une règle (boucle interne)

Procédure **ApprendreMeilleureRègle**(Exemples)

MeilleureRègle = T

ListeOuvert = {T}

Tantque ListeOuvert $\neq \emptyset$

tmpListeOuvert = \emptyset

Pour tout $h \in$ ListeOuvert

ListeOuvert = ListeOuvert \setminus h

Pour chaque $\rho(h) = h', h' \in H$

Si h' est correcte alors % solution potentielle

Si $\text{évaluation}(h') > \text{évaluation}(\text{MeilleureRègle})$

alors MeilleureRègle = h'

Sinon Ajouter h' dans tmpListOuvert

Élaguer tmpListOuvert % heuristique

ListeOuvert = tmpListOuvert

FinTantque

Retourner MeilleureRègle

29/10/2007

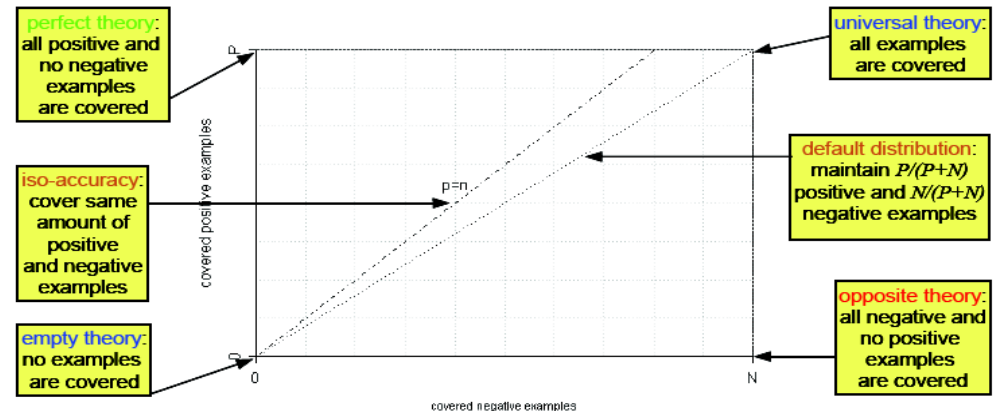
C. Rouveirol

Cours Apprentissage symbolique

11

Coverage Spaces

- good tools for visualizing properties of covering algorithms
 - each point is a theory covering p positive and n negative examples



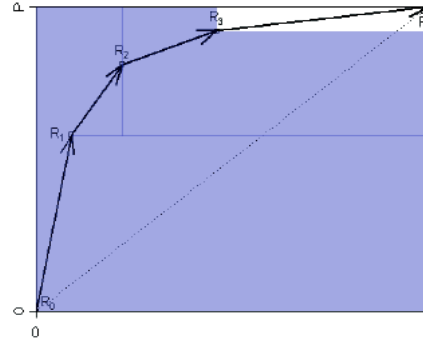
covered negative examples

19

© J. Fürnkranz

Covering Strategy

- Covering or Separate-and-Conquer rule learning algorithms learn one rule at a time
- This corresponds to a path in coverage space:
 - The **empty theory** R_0 (no rules) corresponds to $(0,0)$
 - Adding one rule **never decreases p or n** because adding a rule covers *more* examples (generalization)
 - The **universal theory** R_+ (all examples are positive) corresponds to (N,P)



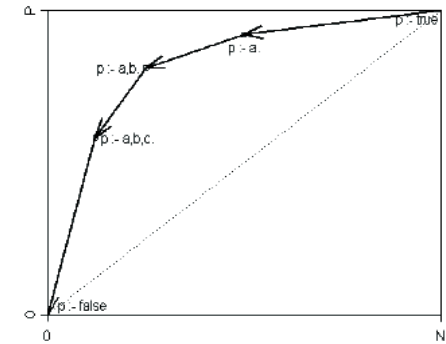
23

© J. Fürnkranz

Top-Down Hill-Climbing

- Successively extends a rule by adding conditions

- This corresponds to a path in coverage space:
 - The rule $p: \text{-true}$ covers all examples (universal theory)
 - Adding a condition never increases p or n (specialization)
 - The rule $p: \text{-false}$ covers no examples (empty theory)



- which conditions are selected depends on a *heuristic function* that estimates the quality of the rule

25

© J. Fürnkranz

Rule Learning Heuristics

- Adding a rule should
 - increase the number of covered negative examples as little as possible (do not decrease *consistency*)
 - increase the number of covered positive examples as much as possible (increase *completeness*)
- An evaluation heuristic should therefore trade off these two extremes

- Example: **Laplace heuristic**

$$h_{Lap} = \frac{p+1}{p+n+2}$$

- grows with $p \rightarrow \infty$
- grows with $n \rightarrow 0$
- Note: Precision is not a good heuristic. Why?

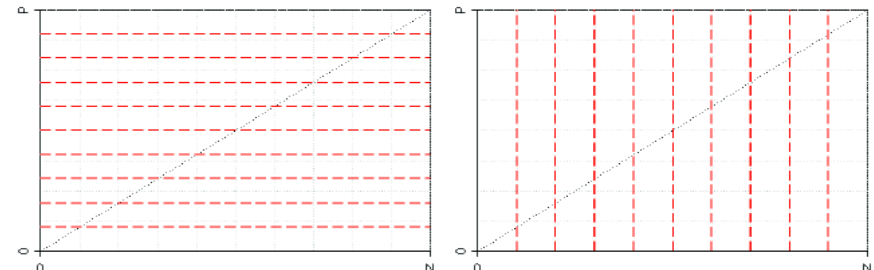
$$h_{Prec} = \frac{p}{p+n}$$

26

© J. Fürnkranz

Isometrics in Coverage Space

- Isometrics are lines that connect points for which a function in p and n has equal values
 - Examples: Isometrics for heuristics $h_p = p$ and $h_n = -n$



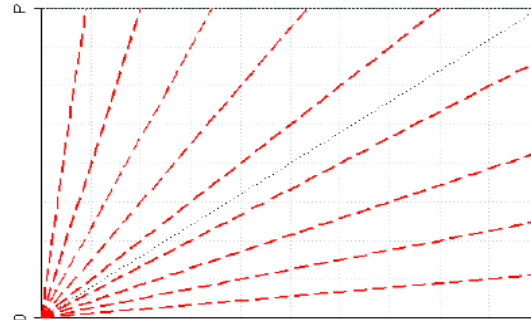
28

© J. Fürnkranz

Precision (Confidence)

$$h_{Prec} = \frac{p}{p+n}$$

- **basic idea:** percentage of positive examples among covered examples
- **effects:**
 - rotation around origin (0,0)
 - all rules with same angle equivalent
 - in particular, all rules on P/N axes are equivalent



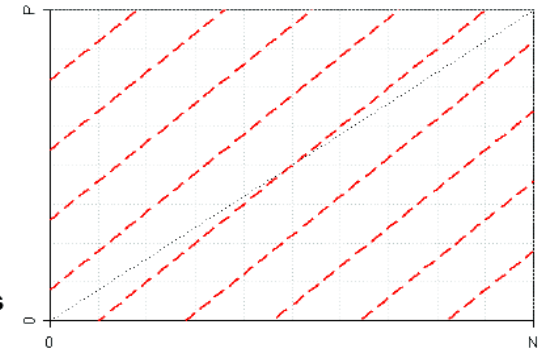
29

© J. Fürnkranz

Accuracy

$$h_{Acc} = \frac{p+(N-n)}{P+N} \approx p-n$$

- **basic idea:** percentage of correct classifications (covered positives plus uncovered negatives)
- **effects:**
 - isometrics are parallel to 45° line
 - covering one positive example is as good as not covering one negative example



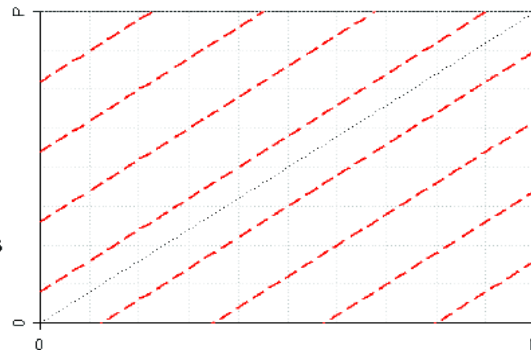
31

© J. Fürnkranz

Weighted Relative Accuracy

$$h_{Acc} = \frac{p+n}{P+N} \left(\frac{p}{p+n} - \frac{P}{P+N} \right) \approx \frac{p}{P} - \frac{n}{N}$$

- **basic idea:** normalize accuracy with the class distribution
- **effects:**
 - isometrics are parallel to diagonal
 - covering x% of the positive examples is as good as not covering x% of the negative examples



32

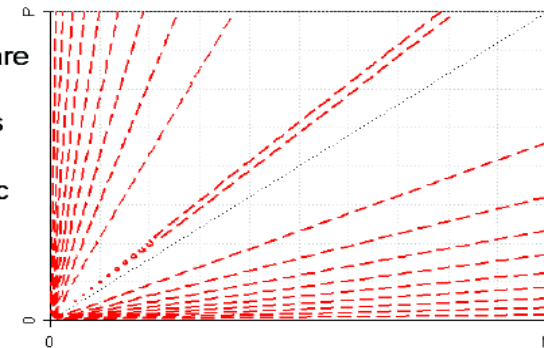
© J. Fürnkranz

Entropy and Gini Index

$$h_{Ent} = -\left(\frac{p}{p+n} \log_2 \frac{p}{p+n} + \frac{n}{p+n} \log_2 \frac{n}{p+n} \right)$$

$$h_{Gini} = 1 - \left(\frac{p}{p+n} \right)^2 - \left(\frac{n}{p+n} \right)^2 \approx \frac{pn}{(p+n)^2}$$

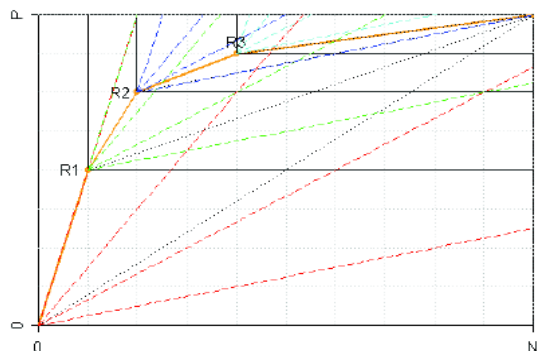
- **effects:**
 - entropy and Gini index are equivalent
 - like precision, isometrics rotate around (0,0)
 - isometrics are symmetric around 45° line
 - a rule that only covers negative examples is as good as a rule that only covers positives



33

Optimizing Precision

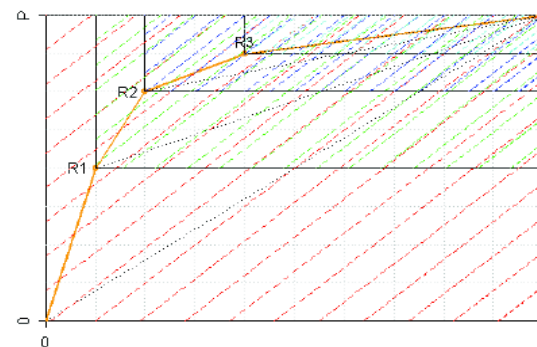
- Precision tries to pick the steepest continuation of the curve
 - does not assume any costs



© J. Fürnkranz

Optimizing Accuracy

- Accuracy assumes the same costs in all subspaces
 - a local optimum in a sub-space is also a global optimum in the entire space



© J. Fürnkranz

RIPPER (Cohen, 95)

Généralisation du gain d'information utilisé dans ID3 et C4.5 à l'apprentissage de règles :

- Recherche : gloutonne
- Opérateur de spécialisation: ajout de descripteurs
 - $A = v, A < v, A \geq v$
- Fonction d'évaluation : Gain d'information

R: règle à spécialiser

L: littéral candidat à ajouter à R

p_0, n_0 : positifs (négatifs) couverts par R

p_1, n_1 : positifs (négatifs) couverts par R+L

p : positifs couverts par R et R+L

- Elagage de l'hypothèse

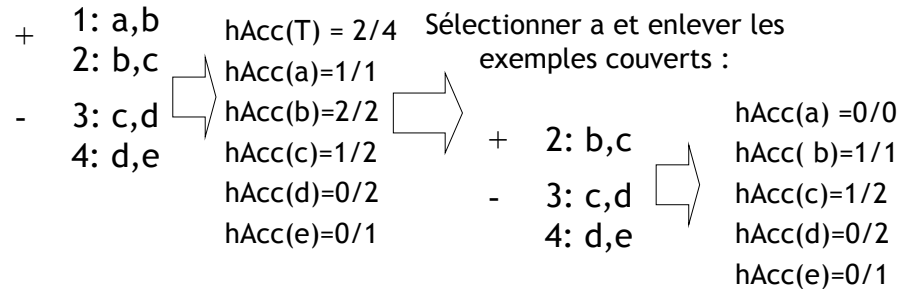
$$Gain(L,R) \equiv p \left(\lg \left(\frac{p_1}{p_1 + n_1} \right) - \lg \left(\frac{p_0}{p_0 + n_0} \right) \right)$$

Autres algorithmes

- Recherche : gloutonne, faisceau, meilleur d'abord, ...
- Opérateur de spécialisation : dépend du langage des hypothèses
- Fonction d'évaluation : précision, Laplacien, m-estimate, ...

Exemple avec heuristique de précision

$$hAcc(r\grave{e}gle, exemples) = \frac{p(r\grave{e}gle)}{p(r\grave{e}gle) + n(r\grave{e}gle)}$$



On sélectionne la DNF : $a \vee b$

CN2 (Clark&Boswell 1991)

- Apprend une liste de décision (Rivest, 1987)
- Une liste de décision: conjonction (bornée) associée à une classe. Les listes sont ordonnées. La dernière conjonction, true, représente la règle par défaut de la liste de décision.
- Exemple: la liste de décision suivante représente $x1 \text{ XOR } x2$
 - $(x1x2, 0)$
 - $(x1, \text{not}(x2)), 1)$
 - $(\text{not}(x1)\text{not}(x2), 0)$
 - $(\text{true}, 1)$

CN2 (Clark&Boswell 1991)

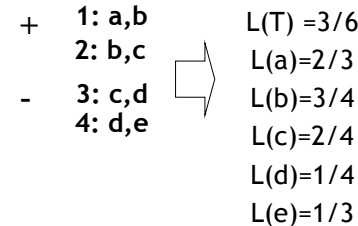
- Recherche : faisceau (taille 5)
- Opérateur de spécialisation : $A=v, A$
- Recherche heuristique : dans la version originale

$$h_{ent} = -\left(\frac{p}{p+n} \log_2 \frac{p}{p+n} + \frac{n}{p+n} \log_2 \frac{n}{p+n}\right)$$

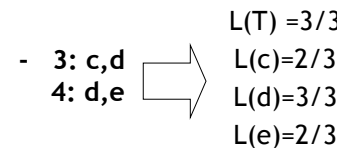
- Puis, évolution vers :

$$LaPlace(r\grave{e}gle, exemples) = \frac{p(r\grave{e}gle) + 1}{p(r\grave{e}gle) + n(r\grave{e}gle) + NumClasses}$$

Exemple CN2



On sélectionne le faisceau=2. On garde a, b. On continue la recherche $a \wedge b, a \wedge c, b \wedge c$. La meilleure règle obtenue est: si b alors classe = +



La meilleure règle obtenue est: si d alors classe = -.

Stratégie d'apprentissage descendante «dirigée par les données»

- Un système emblématique de la stratégie DDD : AQ (Michalski et al. 78, 80, 83)

Procédure :

1. Choisir un exemple dans E^+ comme graine
 2. Rejeter les exemples négatifs un à un en produisant uniquement les raffinements qui couvrent la graine
- $\rho(h, e^-)$ dépend de la graine $\Rightarrow \rho_g(h, e^-)$

29/10/2007

C. Rouveïrol

Cours Apprentissage symbolique

29

Opérateurs de spécialisation «dirigés par les données»

Soit $g = (a_1, \dots, a_i, \dots, a_n)$ la graine, $e^- = (b_1, \dots, b_i, \dots, b_n)$ le négatif à rejeter :

- opérateur nominal :

$$- a_i = v \text{ et } b_i \neq v \Rightarrow a_i = v$$

$$- a_i = v \text{ et } b_i = v' \Rightarrow a_i \neq v$$

- opérateur numérique :

$$- a_i = v \text{ et } b_i = v' \Rightarrow a_i \in]v', + [\text{ si } v' > v \text{ et sinon } a_i \in]v, + [$$

Des opérateurs définis pour les hiérarchies, les ensembles (Michalski, 83)

29/10/2007

C. Rouveïrol

Cours Apprentissage symbolique

30

Le système DDD AQ

- Recherche : faisceau
- Opérateur de spécialisation : opérateurs DDD
- Fonction d'évaluation : LEX (fonction lexicographique)
 - Maximiser le nombre d'exemples couverts par la règle
 - Minimiser le nombre d'exemples négatifs couverts
 - Minimiser le nombre d'attributs de la règle
 - ...

29/10/2007

C. Rouveïrol

Cours Apprentissage symbolique

31

Stratégie d'apprentissage descendante «dirigées par les données»

- Remarque fondamentale : le nombre de raffinement possibles d'une hypothèse dépend du nombre de valeurs d'attributs en commun entre la graine et l'exemple à rejeter
- Définition d'un ordre total sur les négatifs : on traite les plus proches d'abord
- Si le négatif et la graine ne diffèrent que d'un attribut alors G ne fragmente pas : une nuance critique (ou *near-miss* de Winston, 75)

29/10/2007

C. Rouveïrol

Cours Apprentissage symbolique

32

A Pathology for Top-Down Learning

- Parity problems (e.g. XOR)
 - r relevant binary attributes
 - s irrelevant binary attributes
 - each of the $n = r + s$ attributes has values 0/1 with probability $\frac{1}{2}$
 - an example is positive if the number of 1's in the relevant attributes is even, negative otherwise
- Problem for top-down learning:
 - by construction, each condition of the form $a_i = 0$ or $a_i = 1$ covers approximately 50% positive and 50% negative examples
 - irrespective of whether a_i is a relevant or an irrelevant attribute
 - top-down hill-climbing cannot learn this type of concept
- Typical recommendation:
 - use *bottom-up learning* for such problems

41

© J. Fürnkranz

Stratégies d'apprentissage «dirigées par les données»

- Approche ascendante : on parcourt par généralisation à partir de l'élément nul pour trouver un élément de S
 - Algorithme par moindre-généralisé

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

34

Stratégie d'apprentissage ascendante

Opérateurs de raffinement:

- «générer-tester» : fondés sur la structure de H seulement
 - $\delta(h) = h'$ avec $h' \in H$ et $h \leq_h h'$
- «dirigés par les données» : utilisation d'un exemple positif pour généraliser
 - $\delta(h, e+) = h'$ avec $h' \in H$ et $h \leq_h h'$ et $\text{couvre}(h', e+)$
- On appelle couverture supérieure de h l'ensemble d'hypothèses :
 $CS(h) = \{h' \in H \mid h \leq_h h' \text{ et il n'existe pas de } h'' \text{ t.q. } h \leq_h h'' \leq_h h'\}$

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

35

Opérateur de généralisation «dirigé par les données»

- Adaptation de l'opérateur de moindre-généralisé au traitement des exemples d'apprentissage
- => astuce de la représentation unique (*single representation trick*, Dietterich et al, 1982) : on plonge l'espace des instances dans l'espace des hypothèses. Un exemple peut toujours être vu comme la règle la plus spécifique couvrant uniquement cette exemple

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

36

Opérateur de généralisation «dirigés par les données»

- Soient deux hypothèses h_1 et $h_2 \in H$
 - $\delta(h_1, h_2) = h_1 \vee h_2 = h'$ avec $h' \in H$ et $h_1 \leq_h h'$ et $h_2 \leq_h h'$
- Plotkin (70) a le premier montré l'intérêt du calcul de la borne supérieure d'un ensemble pour le problème de l'apprentissage. On parle de moindre-généralisé ou de **lgg** (least-general-generalization) ou de **msg** (most-specific generalization)
- La lgg permet de généraliser par plus petit pas (par définition)
- Son calcul dépend du poset utilisé comme espace d'hypothèses
- Elle est définie uniquement entre hypothèses

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

37

Opérateur de généralisation «dirigé par les données»

Moindre-généralisé en langage attribut-valeur

Soit $e_j = (a_{j1}, \dots, a_{ji}, \dots, a_{jn})$ la description de l'exemple e_j

$\text{lgg}(e_i, e_j) = h = (b_1, \dots, b_k, \dots, b_n)$, t.q. pour tout k allant de 1 à n

Si $a_{ik} = a_{jk}$ alors $b_k = a_{ik}$

Sinon

Si a_{ik} est de type nominal, alors $b_k = ?$ (une valeur indéterminée)

Si a_{ik} est de type continu, alors $b_k =$ le plus petit intervalle contenant a_{ik} et a_{jk}

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

38

Opérateur de généralisation «dirigés par les données»

Moindre-généralisé :

Exemple : $e_1 = (\text{age}=24, \text{sexe}=F, \text{taille}=1.8, \text{yeux}=M)$

$e_2 = (\text{age}=26, \text{sexe}=H, \text{taille}=1.8, \text{yeux}=M)$

$\text{lgg}(e_1, e_2) = (\text{age}=[24,26], \text{sexe}=?, \text{taille}=1.8, \text{yeux}=M) = h$

Le même calcul peut se faire entre deux hypothèses ou entre une hypothèse et un exemple

Exemple : $e_3 = (\text{age}=27, \text{sexe}=F, \text{taille}=1.8, \text{yeux}=B)$

$\text{mg}(h, e_3) = (\text{age}=[24,27], \text{sexe}=?, \text{taille}=1.8, \text{yeux}=?)$

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

39

Stratégie ascendante guidée par les données

```
ApprendMeilleureRègle(E+, E-)
MeilleureRègle = SélectionnerUnExemple(E+);
Sol = MeilleureRègle; ListeOuvert := MeilleureRègle
while ListeOuvert ≠ ∅
  TempListeOuvert := ∅
  for each h ∈ ListeOuvert do
    for each e ∈ E+ - ext(h) do
      TempListeOuvert := TempListeOuvert ∪ lgg(h, e)
    end for
  Candidats := Candidats ∪ TempListeOUverte
end for
Sol := MeilleurEvaluation(Sol ∪ Candidats)
Candidats := MaximalementSpécifiques(ElagageBU(Candidats, Sol))
L := KmeilleursGuidageBU(Candidats, k)
end while
return(Sol)
end
```

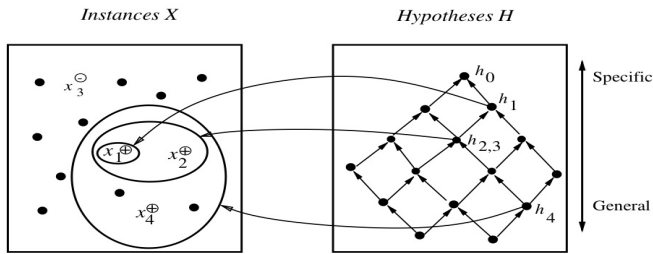
29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

40

Stratégie ADD : FIND-S



X_1 <soleil chaud normal fort chaude stable oui > h_1 <soleil chaud normal fort chaude stable oui >
 X_2 <soleil chaud forte fort chaude stable oui > h_2 <soleil chaud ? fort chaude stable oui >
 X_3 <pluie froid forte fort chaude instable non > h_3 <soleil chaud ? fort chaude stable oui >
 X_4 <soleil chaud forte fort fraîche instable oui > h_4 <soleil chaud ? fort ? ? oui >

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

41

Différents algorithmes ADD

Différenciés par le choix de la graine :

- le premier exemple (FIND-S)
- recherche par faisceau + choix aléatoire de paires d'exemples positifs (GOLEM, Muggleton et al. 92)
- recherche par faisceau + autant de graines que la taille du faisceau (ELENA, Brezellec et al. 93)
- autant de graines que d'exemples positifs + recherche du plus petit sous-ensemble de règles (GLOBO, Torre 99)

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

42

Conclusion sur les approches algorithmiques de l'apprentissage symbolique

- Algorithme de covering lorsque le concept à apprendre est caractérisé par un ensemble de règles.
- Le parcours est guidé par une heuristique (technique de parcours + fonction d'évaluation)
- Le parcours de l'espace de recherche se fait au moyen d'opérateurs de raffinement (ascendant ou descendant, «générer-tester» ou «dirigés par les données»)
- Un système d'apprentissage est caractérisé par sa stratégie d'apprentissage, son langage d'hypothèses, son type d'opérateurs, sa recherche heuristique

29/10/2007

29/10/2007

C. Rouveirol

Cours Apprentissage symbolique

44

References

- Clark, P. and Boswell, R. Rule induction with CN2: some recent improvements. In Machine Learning - EWSL-91. Proceedings of the European Working Session on Learning., pages 151--163, 1991
- Clark, P. and Niblett, T. The CN2 induction algorithm. Machine Learning, 3(4):261--284, 1989
- Dietterich, T. G., London, R. L., Clarkson, K., and Dromey, G. . Learning and inductive inference. Chapter XIV in The Handbook of Artificial Intelligence, Vol. III, 323-512, William Kaufmann, 1982
- Fürnkranz J. and Flach P. A., "ROC'n Rule Learning: Towards a Better Understanding of Covering Algorithms", Machine Learning, 58 (1),39-77, 2005
- Plotkin, G. D. A note on inductive generalisation. Machine Intelligence 6:101–124, 1971.