

Fouille de données

Christel Vrain

{Christel.Vrain@lifo.univ-orleans.fr}.

Contraintes et Apprentissage

LIFO (FRE 2490)

Université d'Orléans



Plan de l'exposé

1. Présentation de la fouille de données

2. Les travaux menés au LIFO

(a) Programmation Logique Inductive et Classification

(b) Recherche de règles d'association

(c) Recherche de règles de caractérisation

(d) Découverte de classes

(e) Fouille de textes

**Extraction de Connaissances dans
les Bases de Données
(ECD)**

**Knowledge Discovery in Databases
(KDD)**

ECD

extraction non triviale de connaissances implicites, inconnues au préalable, intéressantes, utiles à partir d'informations stockées dans des bases de données

ECD

extraction non triviale de connaissances implicites, inconnues au préalable, intéressantes, utiles à partir d'informations stockées dans des bases de données

- Quels types de connaissances ?
 - lois numériques
 - règles d'association
 - règles de classification
 - ...
- ← dépend de la tâche à résoudre

ECD

extraction non triviale de connaissances implicites, inconnues au préalable, intéressantes, utiles à partir d'informations stockées dans des bases de données

- Quels types de connaissances ?
- Champ de recherche multidisciplinaire
 - Bases de Données (BD)
 - Statistiques
 - Apprentissage Statistique / Symbolique
 - Visualisation

Pourquoi cette émergence ?

un univers de données en plein expansion, mais peu de connaissances

- de nombreuses bases de données, de plus en plus volumineuses
 - connexion des BDs par les réseaux
 - Internet
- besoin “vital” d’outils d’aide à la décision pour interpréter les données
- développement des techniques d’apprentissage automatique

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts

- caractérisation

- discrimination

→ apprentissage à partir d'exemples positifs

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts
- Classification

→ apprentissage à partir d'exemples positifs et négatifs

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts
- Classification
- Régression

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts
- Classification
- Régression
- Prédiction

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts
- Classification
- Régression
- Prédiction

Apprentissage non supervisé

- Analyse d'associations

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts
- Classification
- Régression
- Prédiction

Apprentissage non supervisé

- Analyse d'associations
- Découverte de classes

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts
- Classification
- Régression
- Prédiction

Apprentissage non supervisé

- Analyse d'associations
- Découverte de classes
- Organisation en hiérarchies

Différentes tâches

Apprentissage supervisé

- Description de classes ou concepts
- Classification
- Régression
- Prédiction

Apprentissage non supervisé

- Analyse d'associations
- Découverte de classes
- Organisation en hiérarchies
- Recherche des anomalies

Des exemples d'applications

- Analyse du panier de la ménagère
- Fidélisation des clients
- Publicité ciblée
- Détection des fraudes

Le processus d'ECD

- Prétraitement des données
 - Nettoyage
 - Sélection des attributs pertinents
 - Construction d'attributs synthétisés
 - Aplatissement de la base de données / *Prise en compte du caractère relationnel* de la base
- Fouille de données
 - Apprentissage Automatique
 - Apprentissage Statistique
 - Analyse de Données
- Interprétation et Validation des résultats

Fondements formels : Programmation Logique Inductive

Programmation Logique Inductive

- À la frontière entre l'Apprentissage Symbolique Automatique et la Programmation Logique

- représentation des connaissances sous une forme *compréhensible* pour l'expert

⇒ formalisme à base de règles

$\hat{a}ge = moyen \wedge ens_sup \rightarrow utilise_internet$

- représentation en logique du 1er ordre

$grand_p\grave{e}re(X, Y) \leftarrow p\grave{e}re(X, Z), p\grave{e}re(Z, Y)$

Programmation Logique Inductive

- À la frontière entre l'Apprentissage Symbolique Automatique et la Programmation Logique

- représentation des connaissances sous une forme *compréhensible* pour l'expert

⇒ formalisme à base de règles

$$\hat{a}ge = moyen \wedge ens_sup \rightarrow utilise_internet$$

- représentation en logique du 1er ordre

$$grand_p\grave{e}re(X, Y) \leftarrow p\grave{e}re(X, Z), p\grave{e}re(Z, Y)$$

- Induction de programmes logiques / Apprentissage de connaissances “expertes”

- traitement de données numériques et bruitées v
- volume de données

Tâche de classification - Exemple 1

Bases de données - Exemples positifs et négatifs

<i>Personne</i>			<i>Père</i>		<i>positif</i>		<i>négatif</i>	
<i>Id</i>	<i>Nom</i>	<i>Prénom</i>	<i>Id₁</i>	<i>Id₂</i>	<i>Id₁</i>	<i>Id₂</i>	<i>Id₁</i>	<i>Id₂</i>
1	<i>Dupont</i>	<i>Jean</i>	1	2	1	3	3	1
2	<i>Dupont</i>	<i>Paul</i>	2	3	1	4	1	2
3	<i>Dupont</i>	<i>Marie</i>	2	4				
4	<i>Dupont</i>	<i>Antoine</i>						

Connaissances apprises

$grand_père(X, Y) \leftarrow père(X, Z), père(Z, Y)$



Tâche de classification - Exemple 2

Connaissances du domaine	Concept cible	
	positifs	négatifs
$zéro(0), succ(0, 1)$ $succ(1, 2), succ(2, 3)$	$pair(0), pair(2)$ $impair(1), impair(3)$	<i>le reste</i>

⇒ trouver un programme logique \mathcal{P} définissant *pair* et *impair*

⇒ programme correct

$$pair(X) \leftarrow succ(X, Y), impair(Y)$$
$$impair(X) \leftarrow \neg pair(X)$$



Principe

recherche d'une "bonne" hypothèse dans l'espace des hypothèses

→ relation de généralité

● Stratégies de recherche

- exploration déterministe : critères statistiques

→ ICN, MULT_ICN [*Martin, Vrain*]

- exploration stochastique

→ programmation génétique [*Martin, Moal, Vrain*]

→ algorithme génétique → GRIL [*Braud, Vrain*]

● Biais syntaxiques

→ modélisation des biais par une grammaire d'arbres
[*Moal*]

Limitations et extensions

- Limitation à des programmes Datalog

- traduction des symboles de fonction en symboles de prédicats

fonction successeur représentée par $\text{succ}(X, Y)$

- spécification d'un ensemble fini de constantes,
⇒ Les domaines numériques doivent être bornées.

Si $D = \{0, 1, 2, 3\}$, que vaut $\text{succ}(3, ?)$

- Les expressions ne sont pas évaluées.

$\text{pair}(s(s(X))) \leftarrow \text{pair}(X)$

⇒ représentation inadéquate des données numériques

Limitations et extensions

- Limitation à des programmes Datalog

- Extensions

→ Programmation Logique avec Contraintes [*Martin, Vrain*]

- un langage pour exprimer les contraintes
- un domaine de calcul

$$taxe(P, Y) \leftarrow Y = 0.186 * X, prix(P, X)$$

Limitations et extensions

- Limitation à des programmes Datalog
- Extensions
 - Programmation Logique avec Contraintes [*Martin, Vrain*]
 - Bases de Données Contraintes [*Turmeaux, Vrain*]

Limitations et extensions

- Limitation à des programmes Datalog
- Extensions
 - Programmation Logique avec Contraintes [*Martin, Vrain*]
 - Bases de Données Contraintes [*Turmeaux, Vrain*]
 - Applications à la Fouille de Données
 - ← capacité à traiter plusieurs relations
 - méta-informations disponibles : types, contraintes d'intégrité, ...
 - taille importante des données : place mémoire, coût du test de couverture, ...

Limitations et extensions

- Limitation à des programmes Datalog
- Extensions
 - Programmation Logique avec Contraintes [*Martin, Vrain*]
 - Bases de Données Contraintes [*Turmeaux, Vrain*]
 - Applications à la Fouille de Données
 - Domaines visés :
 - Bases de Données Relationnelles
 - Systèmes d'Information Géographiques

Travaux Actuels

Personnes

- Maître de conférences
 - Sylvie Billot
 - Matthieu Exbrayat
 - Lionel Martin
 - Frédéric Moal
- ATER
 - Agnès Braud
 - Ansaf Salleb
- Doctorants
 - Andrei Letchnenko
 - Guillaume Cleuziou
 - Teddy Turmeaux

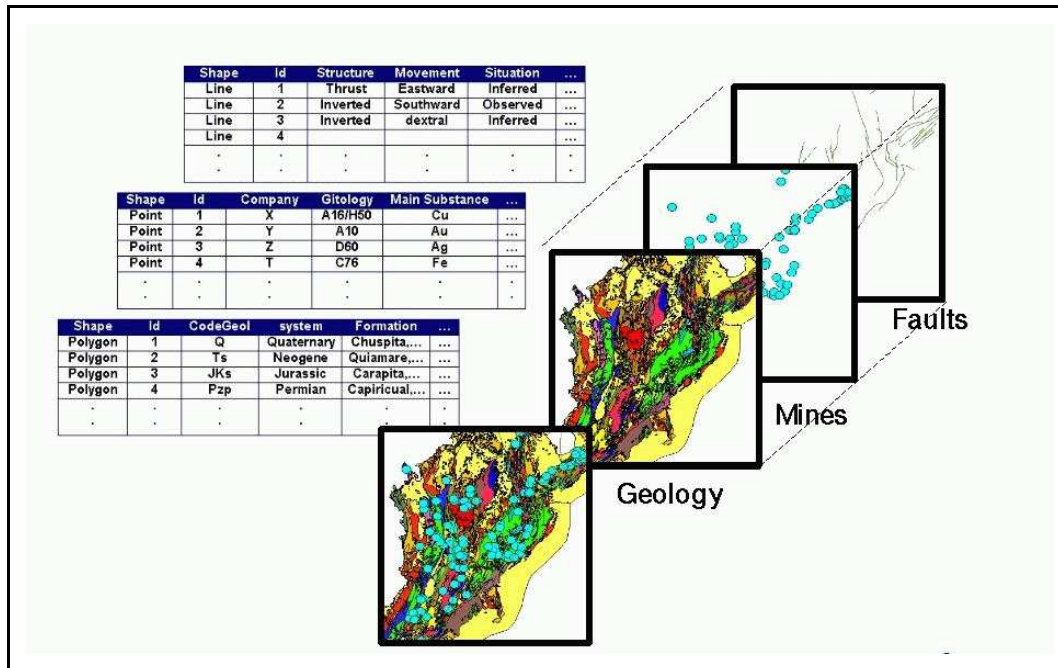
Travaux actuels

- apprentissage et BD relationnelle :
complexité [*Exbrayat, Braud, Turmeaux*]
- approche à base de distances [*Martin, Moal*]
- Apprentissage supervisé
 - Règles de classification
 - Règles de caractérisation [*Turmeaux, Salleb, Vrain*]
- Apprentissage non supervisé
 - Règles d'association [*Salleb, Vrain*]
 - Regroupement (Clustering) [*Cleuziou, Martin, Vrain*]
→ recherche de classes non disjointes
 - BD transactionnelles et chargement en mémoire
[*Maazouzi, Salleb, Vrain*]

Applications

- collaboration avec le BRGM [*D. Cassard*]
- données athérosclérose (PKDD challenge)
- apprentissage de solveurs [*Letchnenko & al.*]
- applications aux textes [*Billot, Cleuziou, Martin, Vrain*]

Collaboration avec le BRGM



- rechercher des règles d'association entre les dépôts de minerais, les mines, les failles, les volcans ...
- caractériser des dépôts de minerais □

Fouille de textes

- ACI BIOTIM : IRD, INRIA (ATOLL ET IMEDIA), CEDRIC, INRA
 - ⇒ méthodes génériques d'exploration de masses de données contenant textes et images pour acquérir la sur-couche sémantique commune
 - ⇒ développer des méthodes génériques d'interrogation pluri-modale des données.
- plate-forme CORAL - LIFO [Billot, Clavier & al.]



Conclusion

- Travaux formels : complétude, correction, complexité des algorithmes
- ECD
 - Tâches :
 - Classification
 - Caractérisation
 - Recherche de règles d'associations
 - Découverte de classes non disjointes
 - Types de bases de données :
 - relationnelles, géographiques
 - textes
 - transactionnelles
 - Applications : BRGM, IRD, ...

Exemple de Bases de Données Contraintes

<i>Ville</i>			<i>Forêt</i>
<i>Nom</i>	<i>Surface</i>	<i>Population</i>	<i>Surface</i>
<i>Paris</i>	$3x - 4y \leq 2$ $x \leq 10$ $x - y \geq -3$ $x + 3y \leq 37$ $x \geq 6$ $3x + 4y \geq 46$	2000000	$2x \leq 19$ $-3x - 4y \leq 2$ $2x - 6y \geq -17$
<i>Rocquencourt</i>	$y \geq 0$ $x \geq 8$ $2x + 2y \leq 17$	3877	$2x - 6y \geq -17$ $2y \geq 11$ $2x - 4y \leq 5$
<i>Orsay</i>	$2x - 2y \geq -13$ $2x + 2y \leq 17$ $2x \geq 1$	15000	$x \geq 9$ $9x - y \leq 85$ $x - y \leq 10$



Opérations algébriques

- *jointure naturelle* de deux relations R_1 et R_2 :

Exemple : $Ville \bowtie Forêt$

- *sélection* de R sur une contrainte c :

Exemple : $\sigma_{Population \geq 20000}(Ville)$

- *projection* de R sur \tilde{Z} , ($\tilde{Z} \subseteq Var(R)$)

Exemple : $\Pi_{X_o}(Ville)$

- *union* de deux relations : $R_1 \cup R_2 = \{t | t \in R_1 \text{ or } t \in R_2\}$



Exemple d'apprentissage

<i>Nom</i>	<i>Classe</i>	<i>Définitions possibles</i>
<i>Orsay</i>	(+)	$\Pi_{Name}(Ville \bowtie Forêt)$
<i>Paris</i>	(+)	$\Pi_{Name}(\sigma_{Population > 3877}(Ville))$
<i>Rocquencourt</i>	(-)	



Intérêts

- un cadre unifié pour représenter les domaines numériques et symboliques
 - permet de représenter des relations avec des tuples infinis
 - permet de représenter des ensembles infinis d'exemples
- ⇒ extension de l'ILP
- implémentation ouverte :
définition de nouvelles stratégies (modèle, stochastique ...)
 - Problèmes : données tests ?



Règles d'association

- Recherche des propriétés fréquentes
 - étape coûteuse en temps de calcul et en espace mémoire
 - élagage
- Recherche des règles (seuil de confiance)
 - génération de beaucoup de règles

Règles d'association

- Règles statistiques

$$Mine(x) \wedge Gitologie(x, A) \rightarrow Gitologie(x, A_1)(92, 12\%)$$

- Règles de contrôle

$$\begin{aligned} Mine(x) \wedge Gitologie(x, H12) \\ \rightarrow \\ Substance_principale(x, Au)(89, 32\%) \end{aligned}$$

- Nouvelles règles

$$\begin{aligned} Mine(x) \wedge Faille(z) \wedge Gitologie(x, C5) \wedge Proche_de(x, z) \\ \rightarrow \\ Structure(x, Strike_slip)(43, 75\%) \end{aligned}$$

Caractérisation

tâche descriptive de fouille de données

- ciblé sur un ensemble d'exemples positifs
- ne nécessite pas d'exemples négatifs

⇒ un *cadre général* pour la caractérisation d'un ensemble d'objets, *ensemble cible*, à partir

- des propriétés des objets cibles
- des propriétés des objets liés

Applications

- bases de données géographiques
- bases de données relationnelles

Règle caractéristique

- Règle caractéristique $\delta :: p$: conjonction d'un schéma caractéristique δ et d'une propriété p

$\forall M \text{ Profondeur_Benioff}(M) \in [75..150]$

$\forall M \exists G \text{ Mine}(M) \wedge \text{Geologie}(G) \wedge \text{Age}(G, \text{tertiaire})$

- Utilisation d'agrégats : $\text{Agrégat}_{\text{proche}}(V, \text{count}) \geq 2$
- Recherche pour chaque type de schémas des propriétés p vérifiées par au moins ϵ , ϵ seuil donné

Elagage

- *Règle caractéristique intéressante*
⇒ caractère contrastant de p entre \mathcal{E}_{target} et $\mathcal{E} - \mathcal{E}_{target}$
- *Relation de généralité entre règles*

$$\begin{aligned}\forall M \forall_{3K_m} F \succeq \forall M \forall_{5K_m} F \succeq \forall M \forall_{10K_m} F \\ \forall M \exists_{10K_m} F \succeq \forall M \exists_{5K_m} F \succeq \forall M \exists_{3K_m} F\end{aligned}$$

Propriété :

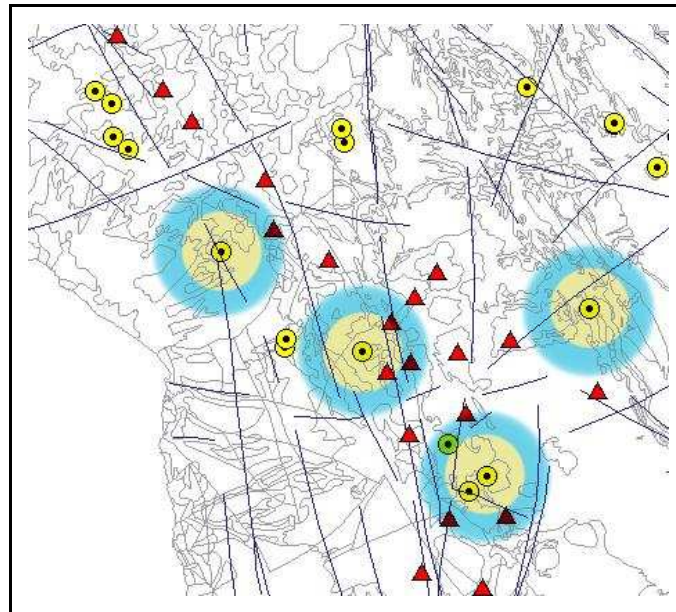
$$\begin{aligned}couverture(\delta_2, p, \mathcal{E}_{target}) \geq \epsilon \\ \Rightarrow \\ couverture(\delta_1, p, \mathcal{E}_{target}) \geq \epsilon\end{aligned}$$

Application

Entrées :

- mines, géologie, volcans, failles, séismes
- relations de distance entre objets

But : caractériser les mines d'or



⇒ construction de *buffers croissants* autour des cibles

Résultat

Un exemple de règle obtenue couvrant près de 60% des mines d'or et rejetant la majeure partie des autres mines.

$$\forall M \exists_{10km} G :: \text{Mine}(M) \wedge \text{Geologie}(G) \wedge$$
$$\text{Substance}(M, \text{or}) \wedge$$
$$\text{Profondeur_Benioff}(M) \in [75..150] \wedge$$
$$\text{Distance_Benioff}(M) \in [170..275] \wedge$$
$$\text{Pente}(M) \in [8^\circ..16^\circ] \wedge$$
$$\text{Age}(G, \text{tertiaire}) \wedge$$
$$\text{Lithologie}(M, \text{volcanique}) \wedge$$
$$\text{Gitologie}(M, \text{épithermale}) \wedge$$
$$\text{Morphologie}(M, \text{veines})$$