

**Université Lumière Lyon2**  
Année 2003

Thèse  
pour obtenir le grade de  
Docteur  
en  
Informatique

présentée et soutenue publiquement par

**Pierre-Emmanuel JOUVE**  
le 10 décembre 2003

# **Apprentissage Non Supervisé et Extraction de Connaissances à partir de Données**

préparée au sein du laboratoire ERIC  
Equipe de Recherche en Ingénierie des Connaissances

sous la direction de  
Nicolas Nicoloyannis

devant le jury, composé de:

Jean-Paul Rasson, Rapporteur	Professeur, Facultés Universitaires N.D. de la Paix, Namur
Gilles Venturini, Rapporteur	Professeur, Université de Tours
Mohand-Saïd Hacid, Examineur	Professeur, Université Claude Bernard-Lyon 1
Michel Lamure, Examineur	Professeur, Université Claude Bernard-Lyon 1
Gilbert Ritschard, Examineur	Professeur, Université de Genève
Nicolas Nicoloyannis, Directeur de thèse	Professeur, Université Lumière-Lyon 2

# Table des matières

<b>1</b>	<b>Introduction, Préambule</b>	<b>1</b>
<b>2</b>	<b>Concepts, Notions et Notations Utiles</b>	<b>7</b>
2.1	Données Catégorielles . . . . .	7
2.1.1	Domaines et Attributs Catégoriels . . . . .	8
2.1.2	Objets Catégoriels . . . . .	9
2.1.2.1	Similarités, Dissimilarités entre Objets Catégoriels . . . . .	10
2.1.3	Ensemble d'Objets Catégoriels . . . . .	11
2.1.3.1	Mode d'un Ensemble d'Objets Catégoriels . . . . .	11
2.1.3.2	Similarités et Dissimilarités entre Ensembles d'Objets Catégoriels . . . . .	12
2.1.3.3	Similarités et Dissimilarités au sein d'un Ensemble d'Objets Catégoriels . . . . .	12
2.1.3.4	Voisinage d'une Partition d'un Ensemble d'Objets Catégoriels . . . . .	13
2.2	Le Nouveau Critère de Condorcet . . . . .	13
<b>3</b>	<b>Classification Non Supervisée</b>	<b>15</b>
3.1	Introduction . . . . .	15
3.1.1	Méthodologie Générale de la Classification Non Supervisée . . . . .	16
3.1.2	Applications de la Classification Non Supervisée . . . . .	16
3.1.3	Taxonomies des Méthodes de Classification Non Supervisée . . . . .	17
3.1.4	Méthodes de Classification Non Supervisée pour Données Catégorielles . . . . .	19
3.1.5	Challenges Actuels en Classification Non Supervisée . . . . .	22
3.2	Une Nouvelle Méthode de Classification Non Supervisée "Orientée Utilisateur" . . . . .	24
3.2.1	Critère d'Évaluation de l'Aspect Naturel d'une Partition d'Objets . . . . .	24
3.2.2	La Méthode de Classification Non Supervisée "Orientée Utilisateur" . . . . .	26
3.2.2.1	Travaux Liés et Spécificités du Travail . . . . .	26
3.2.2.2	L'Algorithme de Classification Non Supervisée . . . . .	27

3.2.2.3	Complexité de l'Algorithme . . . . .	29
3.2.2.4	Qualités de la Méthode pour l'Utilisateur . . .	30
3.2.2.5	Illustration du Fonctionnement de l'Algorithme	30
3.2.3	Evaluation de l'Algorithme de Classification non Super-	
	visée . . . . .	31
3.2.3.1	Evaluation de la Validité des Classifications . .	31
3.2.3.2	Evaluation de la Stabilité . . . . .	37
3.2.3.3	Evaluation de l'Efficacité Algorithmique . . . .	40
3.2.4	Eléments Additionnels . . . . .	42
3.2.4.1	Valeurs Spécifiques pour le Domaine des Va-	
	riables Catégorielles . . . . .	42
3.2.4.2	Gestion des Valeurs Manquantes : . . . . .	44
3.2.4.3	Introduction de Contraintes : . . . . .	44
3.2.4.4	De l'Apprentissage Non Supervisé à l'Appren-	
	tissage Supervisé : l'Apprentissage Non Super-	
	visé sous Contraintes . . . . .	50
3.3	Conclusion . . . . .	54
<b>4</b>	<b>Validité en Apprentissage Non Supervisé</b>	<b>57</b>
4.1	Validité d'une Classification Non Supervisée :	
	Définition et Evaluation . . . . .	58
4.1.1	Mode d'Evaluation par Critères Externes . . . . .	59
4.1.1.1	Méthode de Monte Carlo . . . . .	59
4.1.1.2	Mesures Statistiques . . . . .	60
4.1.2	Mode d'Evaluation par Critères Internes . . . . .	61
4.1.3	Modes d'Evaluation Relatifs . . . . .	63
4.1.3.1	Cas 1 : Le nombre final de classes, $nc$ , n'est pas	
	contenu dans $P_{alg}$ . . . . .	63
4.1.3.2	Cas 2 : Le nombre final de classes, $nc$ , est contenu	
	dans $P_{alg}$ . . . . .	64
4.1.3.3	Indices . . . . .	64
4.1.4	Autres Modes d'Evaluation . . . . .	67
4.2	Nouveaux Indices et Nouvelle Méthodologie pour l'Evaluation	
	et la Comparaison de la Validité de Classifications Non Super-	
	visées . . . . .	68
4.2.1	Concepts et Formalismes Introductifs . . . . .	69
4.2.1.1	Evaluation de l'homogénéité interne des classes	
	d'une $cns$ . . . . .	71
4.2.1.2	Evaluation de la séparation entre classes d'une	
	$cns$ (ou hétérogénéité entre classes)	
	72	
4.2.1.3	Notions Additionnelles . . . . .	73
4.2.1.4	Remarques importantes concernant l'aspect cal-	
	culatoire . . . . .	73

4.2.2	La nouvelle méthodologie pour l'évaluation et la comparaison de validité de cns . . . . .	75
4.2.2.1	Caractérisation statistique des valeurs de: <i>LM</i> et <i>NLD</i> . . . . .	76
4.2.2.2	Méthodologie . . . . .	77
4.2.2.3	Expérimentations . . . . .	82
4.2.2.4	Expérimentations sur le jeu de données Small Soybean Disease . . . . .	82
4.2.3	Expériences sur le jeu de données Mushrooms . . . . .	92
4.2.3.1	Description . . . . .	92
4.2.3.2	Analyse des Résultats . . . . .	95
4.2.4	Résumé et Informations Supplémentaires . . . . .	96
<b>5</b>	<b>Sélection de Variables, Contributions pour l'apprentissage supervisé et non supervisé</b>	<b>105</b>
5.1	Sélection de Variables pour l'Apprentissage Supervisé . . . . .	107
5.1.1	Caractéristiques de la Sélection de Variables . . . . .	107
5.1.2	Les Types de Méthodes . . . . .	107
5.1.3	Directions de Recherche . . . . .	108
5.1.3.1	Forward Selection (FS) (Ajout de variables) . . . . .	108
5.1.3.2	Backward Elimination (BE) (Suppression de variables) . . . . .	109
5.1.3.3	Méthodes Bidirectionnelles . . . . .	109
5.1.4	Stratégie de Recherche . . . . .	109
5.1.5	Fonction d'Evaluation . . . . .	110
5.1.6	Critère d'Arrêt . . . . .	111
5.1.7	Approches Filtres . . . . .	111
5.1.8	Approches Enveloppes . . . . .	114
5.1.9	Autres Approches . . . . .	115
5.2	Contribution à la Sélection de Variables pour l'Apprentissage Supervisé: Une Nouvelle Méthode Efficace et Rapide . . . . .	118
5.2.1	Hypothèses et Idées Fondamentales . . . . .	118
5.2.2	Evaluation de la Validité d'une Partition dans un Sous-Espace de l'ERD . . . . .	119
5.2.3	La Nouvelle Méthode de Sélections de Variables . . . . .	120
5.2.3.1	La Méthode de Base: une Méthode Exhaustive . . . . .	121
5.2.3.2	Réduction de la Complexité par Introduction d'un AG . . . . .	124
5.2.4	Evaluation Expérimentale . . . . .	126
5.2.4.1	Présentation de l'Evaluation Expérimentale . . . . .	126
5.2.4.2	Analyse de l'Evaluation Expérimentale . . . . .	127
5.2.5	Conclusion . . . . .	131
5.3	Contribution à la Sélection de Variables pour l'Apprentissage Non Supervisé: Une Nouvelle Méthode Efficace et Rapide . . . . .	143

5.3.1	Evaluation de l'Adéquation entre deux Ensembles de Variables . . . . .	144
5.3.2	Remarques Importantes Concernant l'Aspect Calculatoire 145	
5.3.3	Evaluation de l'adéquation entre $\mathbf{EV}$ un Ensemble de Variables et $\mathbf{EV}_*$ un Sous Ensemble de $\mathbf{EV}$ ( $\mathbf{EV}_* \subseteq \mathbf{EV}$ )	146
5.3.4	Evaluation/Comparaison de l'Adéquation entre un Ensemble de Variables ( $\mathbf{EV}$ ) et des Sous Ensembles de $\mathbf{EV}$	148
5.3.5	La Nouvelle Méthode de Sélection de Variables . . . . .	148
5.3.6	Evaluations Expérimentales . . . . .	149
5.3.6.1	<b>Expérience #1</b> : Evaluation expérimentale sur jeux de données synthétiques . . . . .	149
5.3.6.2	<b>Expérience #2</b> : Evaluation Expérimentale sur Jeux de Données de l'UCI . . . . .	154
5.3.7	Conclusion . . . . .	162
<b>6</b>	<b>Agrégation de Classifications Non Supervisées : La Problématique "Cluster Ensembles"</b>	<b>165</b>
6.1	Introduction . . . . .	165
6.1.1	Illustration de la Problématique "Cluster Ensembles" . .	167
6.1.2	Motivations, Objectifs de la Problématique "Cluster Ensembles" . . . . .	168
6.1.2.1	Réutilisation de Connaissances . . . . .	169
6.1.2.2	Calcul Distribué pour la cns . . . . .	169
6.1.3	Travaux Liés . . . . .	172
6.1.4	Principaux Challenges pour la Problématique "Cluster Ensembles" . . . . .	175
6.2	Mesures d'Adéquation . . . . .	175
6.2.1	Adéquation entre Classifications Non Supervisées . . .	176
6.2.2	Adéquation pour un Couple de Classification Non Supervisée . . . . .	176
6.2.3	Adéquation entre une Classification Non Supervisée et un Ensemble de Classifications Non Supervisées . . . . .	177
6.3	Contribution à la Problématique "Cluster Ensembles" : Trois Méthodes pour l'Agrégation de Classifications Non Supervisées . .	177
6.3.1	Première Méthode pour l'Agrégation de cns: Une Méthode Intuitive . . . . .	178
6.3.2	Seconde Méthode pour l'Agrégation de Classifications Non Supervisées : Utilisation de la Méthode KEROUAC .	179
6.3.2.1	Utilisation de KEROUAC pour la cns en considérant des Méta-Variables . . . . .	179
6.3.2.2	Relation entre $P_*$ and $P_\beta$ . . . . .	180
6.3.2.3	Conclusion . . . . .	181
6.3.2.4	Illustration . . . . .	181

6.3.2.5	Propriétés de la Méthode . . . . .	182
6.3.3	Troisième Méthode pour l'Agrégation de Classifications Non Supervisées : Utilisation de la Méthode K-Modes . .	183
6.3.3.1	Illustration . . . . .	184
6.3.3.2	Propriétés de la Méthode . . . . .	184
6.3.4	Evaluations Expérimentales . . . . .	184
6.3.4.1	Evaluations, Comparaisons et Discussions Pré- liminaires . . . . .	184
6.3.4.2	Evaluations, Comparaisons et Discussions Com- plémentaires . . . . .	191
6.3.4.3	Comportement de la méthode KEROUAC face à des cns à agréger possédant des nombre de classes très différents . . . . .	204
6.4	Conclusion . . . . .	207
<b>7</b>	<b>Conclusion</b>	<b>211</b>
<b>8</b>	<b>Données Utilisées pour les Expérimentations</b>	<b>217</b>
8.1	Jeu de Données ADULT . . . . .	217
8.2	Jeu de Données MUSHROOMS . . . . .	218
8.3	Jeu de Données BREAST CANCER . . . . .	220
8.4	Jeu de Données CAR . . . . .	222
8.5	Jeu de Données : ADULT . . . . .	224
8.6	Jeu de Données Contraceptive Method Choice . . . . .	225
8.7	Jeu de Données FLAGS . . . . .	226
8.8	Jeu de Données GERMAN . . . . .	227
8.9	Jeu de Données HOUSE VOTES 84 . . . . .	229
8.10	Jeu de Données IONOSPHERE . . . . .	230
8.11	Jeu de Données MONKS . . . . .	231
8.12	Jeu de Données NURSERY . . . . .	232
8.13	Jeu de Données PIMA . . . . .	234
8.14	Jeu de Données SICK . . . . .	235
8.15	Jeu de Données SMALL SOYBEAN DISEASES . . . . .	236
8.16	Jeu de Données VEHICLE . . . . .	237
8.17	Jeu de Données WINE . . . . .	240
8.18	Jeu de Données SPAM . . . . .	241
	<b>Bibliographie</b>	<b>243</b>
	<b>Table des figures</b>	<b>254</b>
	<b>Liste des tableaux</b>	<b>257</b>

## 7 Conclusion

*"On ne fait jamais attention à ce qui a été fait ; on ne voit que ce qui reste à faire."*

- Marie Curie -

Les travaux présentés dans cette thèse participent donc en premier lieu de l'intégration de la classification non supervisée au sein d'un processus d'ECD.

Nous nous sommes penchés sur différentes étapes de ce processus (la sélection de variables au chapitre 5, l'application d'algorithmes de classification non supervisée au chapitre 3, l'évaluation de la validité des résultats au chapitre 4) en nous efforçant de proposer de nouvelles solutions originales permettant d'une part la résolution de problèmes ou l'intégration d'exigences issues de la pratique et d'autre part de bien prendre en compte le rôle central de l'utilisateur dans le processus ECD.

Ainsi, la méthode de sélection de variables proposée (pour la classification non supervisée) semble permettre le traitement de jeux de données plus volumineux en limitant le temps de traitement tout en assurant un bon niveau de qualité des résultats. Cette contribution nous apparaît intéressante d'une part car peu de méthodes permettant l'automatisation de ce processus existent, et d'autre part car certaines des méthodes existantes (telles celles fondées sur l'analyse factorielle) nécessitent d'appliquer l'algorithme de classification non supervisée sur un espace de représentation des données modifié (les axes principaux de l'analyse factorielle par exemple), rendant ainsi plus complexe l'interprétation des résultats.

La mise au point de la méthode de classification non supervisée (pour données catégorielles) KEROUAC a, quant à elle, été motivée par le désir de disposer d'une méthode à l'utilisabilité augmentée : sa mise en œuvre ne nécessite pas une expertise particulière tout comme l'analyse de ses résultats ne suppose pas une interprétation complexe à mener. De plus, cette méthode autorise le traitement de données manquantes ou encore l'intégration de contraintes, de connaissances qui constituent autant d'éléments bien utiles en pratique. Enfin, les éléments classiquement définis comme essentiels à une bonne méthode de classification non supervisée (validité et stabilité des résultats fournis, coût calculatoire "allégé"...) constituent autant de points non négligés par la méthode KEROUAC.

Evaluer aisément la qualité/validité des résultats d'un processus de classification non supervisée correspond là encore à une préoccupation importante tant d'un point de vue pratique qu'académique. Aussi avons nous proposé une méthodologie originale d'évaluation/comparaison de la validité de classification non supervisée, méthodologie s'appuyant sur une analyse graphique relativement intuitive. Cette méthodologie semble de plus se démarquer de la majorité des approches existantes de part son coût calculatoire limité combiné à la possibilité de comparer toutes formes de partitions, et de part l'aspect caractérisation visuelle d'un jeu de données qu'elle propose.

Enfin, de manière à intégrer mieux encore des exigences résultant de la pratique (traitement de données distribuées physiquement, traitement de données hétérogènes, réutilisation de connaissances, nécessité d'accélérer le processus de classification non supervisée...) nous avons introduit la problématique "Cluster Ensembles" et proposé l'utilisation de trois méthodes pour sa résolution. Les intérêts principaux de cette partie sont de souligner de manière simultanée les différents apports de l'agrégation de classification non supervisée et de proposer d'utiliser deux méthodes de classification non supervisée efficaces pour procéder à l'agrégation. Ce dernier point étant particulièrement intéressant car il semble démontrer que la simple réutilisation de méthodes existantes permet d'atteindre des résultats qualitativement aussi intéressants que ceux obtenus avec des méthodes spécifiques (en permettant même une accélération du processus d'agrégation pour la méthode K-Modes, ou la limitation du paramétrage dans le cas de la méthode KEROUAC).

---

Dénominateurs communs à l'ensemble de ces travaux, les concepts de comparaisons par paires, d'agrégation de préférences (concepts sous-jacents au critère de Condorcet) se sont avérés être des outils puissants.

Outre la présentation de solutions concrètes évoquées plus tôt, cette thèse vise également à indiquer que ces concepts véhiculent certainement des éléments de solutions à de nombreux problèmes alors qu'ils apparaissent actuellement sous-exploités.

Assistés des éléments théoriques fournis par la S-Théorie de Michaud [Mic87], [Mic91] (éléments visant notamment à la résolution efficace de problèmes d'optimisation posés dans le cadre précis de l'agrégation de préférences) nous pensons que ces concepts peuvent constituer des bases alternatives pour la mise au point de futures méthodes et méthodologies efficaces dans le cadre de l'ECD.

---

Des travaux additionnels, sortant du domaine du "Non Supervisé" et touchant au domaine du "Supervisé" ont également été abordés, soit en tant que

travaux aboutis (la méthode de sélection de variables pour l'apprentissage supervisé du chapitre 4), soit sous forme de prototype (la méthode d'apprentissage supervisé par apprentissage non supervisé sous contraintes du chapitre 3), ou simplement évoqués comme travaux à venir (mise au point de méthode d'apprentissage généralisé par le biais de KEROUAC (voir chapitre 3), utilisation des méthodes d'agrégation de classifications non supervisées dans le cadre de l'agrégation de classifieurs supervisés (voir chapitre 6)).

Concernant la méthode de sélection de variables pour l'apprentissage supervisé, les expérimentations menées la place au niveau des méthodes de référence actuelles tant pour la qualité des modèles d'apprentissage qu'elle permet de bâtir, que pour la réduction de l'espace de représentation des données qu'elle implique et le temps d'exécution qui lui est associé. Enfin, nous avons vu que cette méthode était fondée sur un certain nombre d'hypothèses dont certaines apparaissent trop strictes et dont la relaxation (envisagée dans des travaux futures) devrait (peut être ?) permettre l'amélioration de ces performances.

La méthode d'apprentissage non supervisé sous contraintes (qui n'est actuellement qu'un prototype) semble, quant à elle, fournir des modèles de qualité relativement similaire à celle des approches classiques, et son développement devrait permettre de répondre à un certain nombre de questions intéressantes touchant notamment à l'intérêt de l'apprentissage semi-supervisé.

Ces travaux (participant non plus de l'intégration de la classification non supervisée au sein du processus d'ECD mais de l'utilisation de la classification non supervisée pour d'autres éléments constitutifs d'un processus d'ECD) présentent donc un intérêt en tant que tel, mais, leur intégration dans ce document est également motivée par la volonté de souligner l'importance (et peut être la prééminence ?) de la structuration de l'information pour le processus d'ECD (la structuration de l'information étant abordée ici par le biais de la classification non supervisée).

En effet, si la relation entre classification non supervisée et apprentissage supervisé par les méthodes de type plus proche voisin est évidente (le concept "qui se ressemblent s'assemblent" sous-jacent à ce type de méthode d'apprentissage supervisé régit également la classification non supervisée) les méthodes d'apprentissage semi-supervisé mettant à profit la structure interne de l'information dans le cadre d'apprentissage supervisé sont encore peu nombreuses et peu exploitées. Ce constat peut être généralisé à l'ensemble des méthodes d'ECD : un faible nombre d'entre elles exploitent vraiment la structuration de l'information.

Le prototype de méthode d'apprentissage non supervisé sous contraintes, la méthode de sélection de variables pour l'apprentissage supervisé (présentée comme une méthode visant à déterminer l'espace de représentation des données tel que sa structure interne reflète au mieux la réalité à apprendre) doivent ainsi être, en partie, considérés comme autant d'éléments contribuant à montrer l'intérêt que peut revêtir l'apport de la structuration dans le cadre

d'un processus d'ECD. Notons enfin que, à l'opposé, il nous semble parfois utile d'adopter une perspective apprentissage supervisé afin de "féconder" des travaux touchant à la classification non supervisée comme le montre la méthode de sélection de variables pour l'apprentissage non supervisé qui peut être envisagée comme une généralisation de la méthode de sélection de variables proposée dans le cadre supervisé (plus spécifiquement une méthode de sélection de variables dans un cadre supervisé où plusieurs variables endogènes cohabitent).

Nous synthétisons finalement les contributions de cette thèse, les relations unissant les univers "Supervisé" et "Non Supervisé" apparaissant dans cette dissertation, ainsi que les travaux à venir sur les figures 7.1, 7.2.

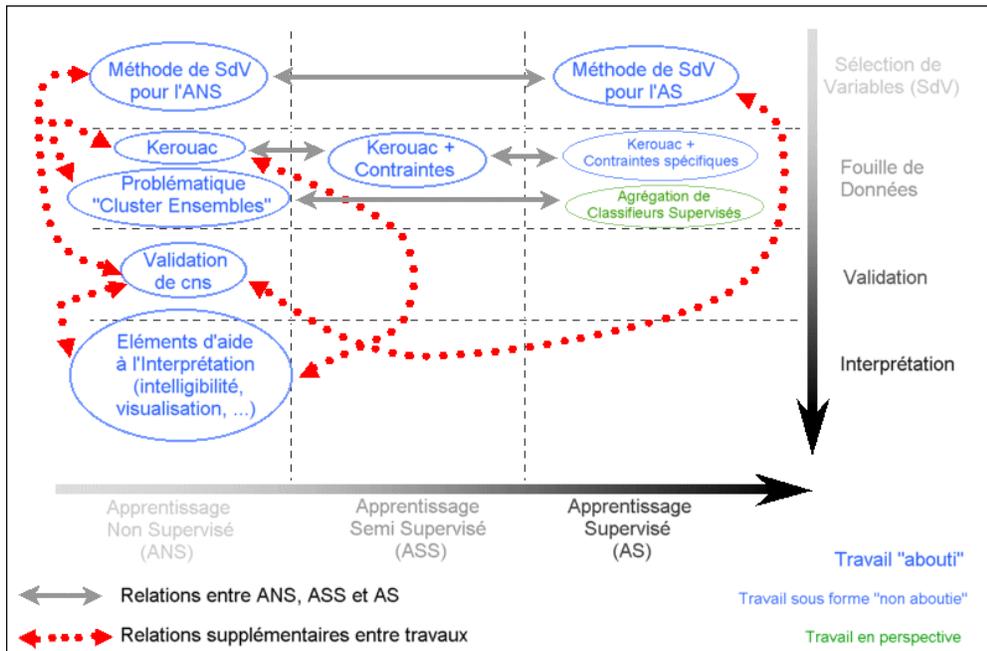


FIG. 7.1 –: Synthèse des contributions, relations entre contributions et relations unissant entre "Supervisé", "Semi-Supervisé" et "Non Supervisé"

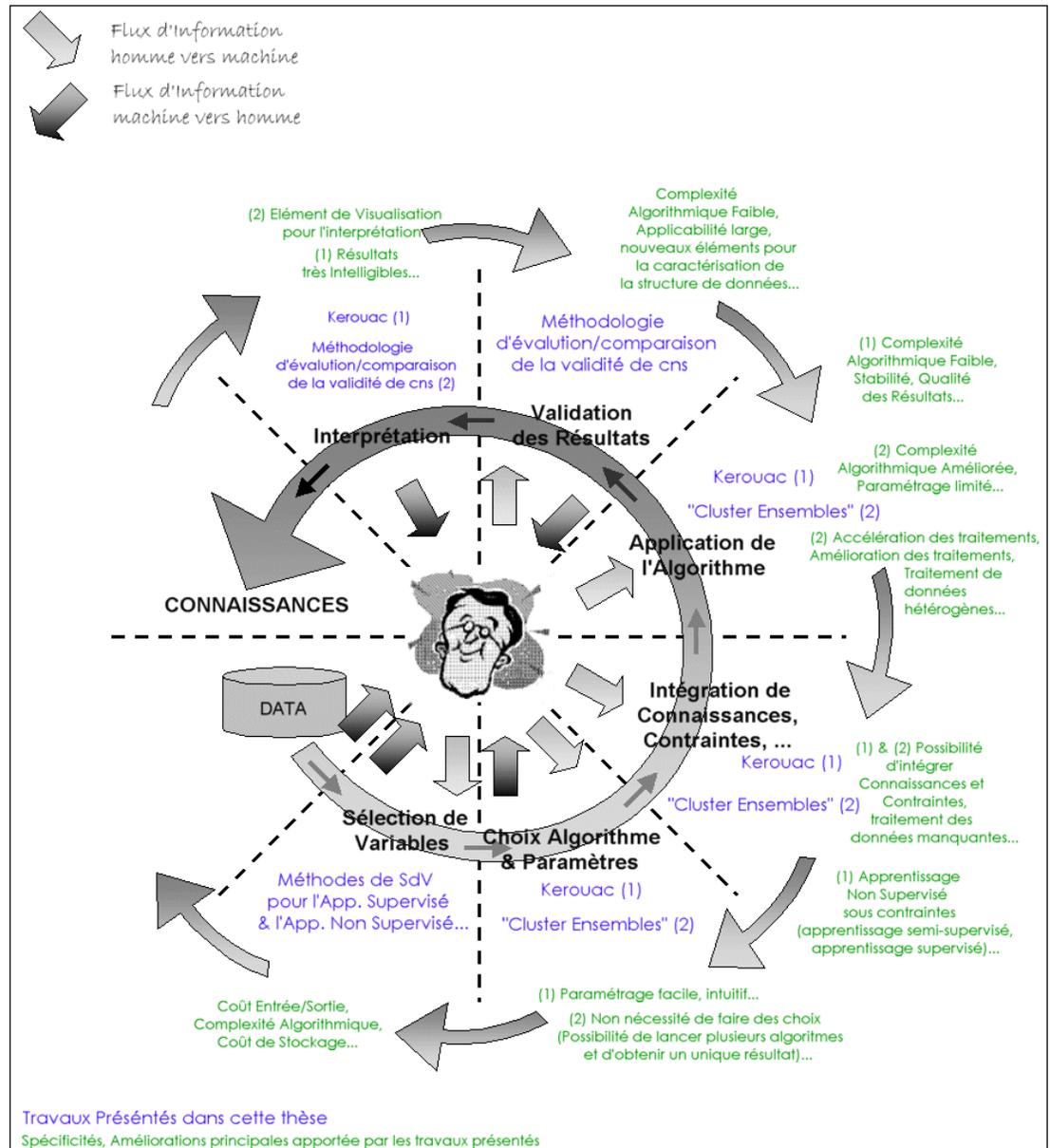


FIG. 7.2 –: Synthèse des contributions