

---

# Projet Involvd : Elicitation interactive de contraintes pour la fouille de données supervisée et semi-supervisée

Christel Vrain\*<sup>†</sup>

<sup>1</sup>LIFO, Université d'Orléans – Université d'Orléans – France

## Résumé

Les recherches récentes en Apprentissage Automatique et Fouille de Données cherchent à automatiser le processus de découverte de connaissances et à réduire les interactions avec l'expert avec de bonnes raisons comme la difficulté à traiter des volumes importants de données (d'autant plus en grandes dimensions) ainsi que les progrès techniques qui ont permis d'alléger les tâches chronophages. Cela a conduit à l'émergence d'offres telles que AutoML service (Google) envisageable dans un cadre supervisé ou les étiquettes des objets peuvent être exploitées pour régler des paramètres ou sélectionner des modèles. En revanche, si les données sont partiellement étiquetées (apprentissage semi-supervisé) ou n'ont pas d'étiquette (non supervisé), la démarche inverse est nécessaire : mettre l'expert dans la boucle d'apprentissage et intégrer ses retours sur les résultats pour améliorer le processus, autrement dit rendre le processus interactif. Cela pose de nouveaux défis comme présenter les résultats pour permettre des retours informés de l'expert, être capable de les expliquer, interagir fréquemment avec l'utilisateur alors qu'AutoML a la possibilité de tourner pendant des heures. Résoudre ces défis non seulement améliore les résultats mais offre un autre avantage : un utilisateur est plus enclin à accepter un résultat si le processus qui a conduit à son émergence est expliqué. Ceci est d'autant plus vrai dans des applications ou les investissements (en argent, temps, vies humaines) reposent sur la justesse des résultats. De plus, les réglementations récentes en Europe et aux États-Unis donnent des droits aux citoyens concernés par des décisions algorithmiques et imposent que les décisions soient expliquées. Ces exigences ont ainsi motivé des recherches sur l'interprétabilité des méthodes de type boîte noire (e.g. apprentissage profond). Pour obtenir des résultats explicables en fouille de données non supervisée ou semi-supervisée, le projet Involvd traite des questions posées par le développement de processus interactif de fouille de données : identification automatique de visualisations faisant sens, explications pour des retours informés, transformation en contraintes opérationnelles et développement de nouveaux systèmes d'apprentissage intégrant ces contraintes. A contrario d'approches de type boîte noire, nous nous fonderons sur le clustering et la recherche de motifs symboliques. Le cas d'usage en chimie-informatique, qui servira de guide tout au long du projet, est un cas typique d'illustration de cette problématique. En conception de médicaments, l'analyse exploratoire de données est capitale : les molécules doivent être comprises en termes de structures et/ou de propriétés chimiques, et les experts ont des connaissances qu'ils ne peuvent expliciter qu'au vu de résultats préliminaires.

---

\*Intervenant

<sup>†</sup>Auteur correspondant: Christel.Vrain@univ-orleans.fr