## **Machine Learning interprétable**

Prénom: Yoann Semestre: 9 Nom: PULL Année: M2 Nature: CM Volume horaire: 12 H ECTS / Coef: 2 Cours de Big Data Analytics I, II, III et IV Prérequis Cours d'introduction à Python Machine Learning sous Python Ce cours présente les notions de base du Machine Learning interprétable et les principales approches techniques permettant de rendre interprétable des modèles de Machine Learning qui ne le sont pas nativement. Le cours s'articule en 3 parties. La première partie est consacrée aux méthodes de type model-agnostic qui permettent d'expliquer le fonctionnement d'un modèle indépendamment de sa structure et qui peuvent s'appliquer à n'importe quel modèle de classification ou de régression. Parmi ces méthodes, nous présenterons les Partial Dependence Plot (PDP), les Individual Conditional Résumé Expectation (ICE) et les Accumulated Local Exects (ALE). La seconde partie est consacrée aux modèles d'approximation globale et locale (surrogate models) et notamment à la méthode Local interpretable model-agnostic explanations (LIME). La dernière partie sera consacrée aux valeurs de Shapley et aux SHapley Additive exPlanations (SHAP). Les applications seront réalisées sur les logiciels Python et SAS. L'objectif de ce cours double. Il s'agit tout d'abord de familiariser les étudiants aux enjeux de gouvernance de l'Intelligence Artificielle et du Machine Learning, notamment dans le **Objectifs** domaine de la Finance, en faisant un focus sur les notions d'interprétabilité et d'explicabilité des algorithmes. Il s'agit ensuite d'introduire les principaux outils permettant de rendre interprétable des modèles de Machine Learning qui ne le sont pas nativement. ACPR (2018), Artificial intelligence: challenges for the financial sector. Discussion papers publication, December 2018. ACPR (2020). Governance of artificial intelligence in finance. Discussion papers publication, November, 2020. Lipton, Z.C. (2018), The Mythos of Model Interpretability: In machine learning, the concept of interpretability is both important and slippery, Queue, 16(3), 31--57. Bibliographie Miller, T. (2019) Explanation in artificial intelligence: Insights from the social sciences, Artificial Intelligence, 267, 1-38. Molnar C., Casalicchio G. and Bischk B. (2020), Interpretable Machine Learning: A Brief History, State-of-the-Art and Challenge, arXiv:2010.09337v1. Molnar, C (2019), Interpretable machine learning. A Guide for Making Black Box Models Explainable,. <a href="https://christophm.github.io/interpretable-ml-book/">https://christophm.github.io/interpretable-ml-book/</a>.

## **PLAN**

Section 1 : Définitions : Interprétabilité et explicabilité

Section 2 : Principaux enjeux et principales méthodes du machine learning interprétable

Section 3 : PDP et approches similaires

Section 3.1. Partial Dependence Plot (PDP)

Section 3.2. Individual Conditional Expectation (ICE)

Section 3.3. Accumulated Local Exects (ALE)

Section 4 : Local and Global Surrogate Models

Section 4.1. Global surrogate

Section 4.2. Local interpretable model-agnostic explanations (LIME)

Section 5 : Shapley Values

Section 5.1. Shapley Values

Section 5.2. SHapley Additive exPlanations (SHAP)