



Université Claude Bernard



DIPLÔME NATIONAL DE DOCTORAT

(Arrêté du 25 mai 2016)

Date de la soutenance : **31 mars 2023**

Nom de famille et prénom de l'auteur : **Monsieur VILLIE Antoine**

Titre de la thèse : « *Quantifier l'incertitude de l'explicabilité en apprentissage automatique : Inférence post-sélection sur des caractéristiques biologiques interprétables* »



Résumé

Les réseaux de neurones artificiels ont récemment été utilisés avec succès pour faire des prédictions sur des séquences biologiques. Ces réseaux sont principalement évalués pour leurs capacités prédictives, et sont souvent critiqués pour leur manque d'interprétabilité. D'un autre côté, plusieurs méthodes issues de la littérature bioinformatique cherchent à aider à comprendre les mécanismes biologiques sous-jacents, en sélectionnant des variants génomiques significativement associés avec le trait biologique d'intérêt. Bien qu'elles aient mené à de nombreuses découvertes durant les dernières années, ces méthodes restent soumises à certaines limitations. Récemment, des outils cherchent à expliquer les prédictions des réseaux de neurones, en extrayant des caractéristiques biologiques interprétables de ces réseaux entraînés. Les réseaux de neurones peuvent alors être compris comme des méthodes permettant de sélectionner des variants génomiques, et peuvent permettre de dépasser certaines des limitations préalablement mentionnées. Mais à notre connaissance, la quantification de la significativité de l'association entre les caractéristiques biologiques extraites de ces réseaux et les traits biologiques d'intérêt n'a reçu que peu d'attention. Nous proposons donc de dépasser la notion d'explicabilité pour l'apprentissage automatique, en cherchant à quantifier statistiquement l'association entre les variants issus de réseaux de neurones et le phénotype, afin de participer à créer un lien entre les méthodes d'apprentissage automatiques et celles provenant de la biologie computationnelle. En particulier, nous formalisons le lien entre réseaux de neurones et sélection de variants biologiques, et nous proposons différentes modifications à ces réseaux, afin d'améliorer leurs performances en tant que méthodes de sélection. Nous proposons également une procédure de test valide pour les variants ainsi sélectionnés, issue des avancées récentes en inférence post-sélection.