

## FOAD - Inter - Machine learning et méthodes statistiques appliquées aux processus de classification

**Les 6, 7, 8 et 9 novembre 2023**

<b>Objectifs</b>	<ul style="list-style-type: none"> <li>- Maîtriser le vocabulaire spécifique aux méthodes d'apprentissage à finalité de classement</li> <li>- Identifier le contexte et les conditions d'application des méthodes d'apprentissage supervisé et non supervisé</li> <li>- Connaître les objectifs et les différences entre les méthodes de classement</li> <li>- Décrire la méthodologie inhérente à ces méthodes</li> <li>- Mettre en oeuvre et interpréter les résultats des méthodes d'apprentissage supervisé</li> <li>- Connaître les principaux indicateurs de cohérence liés aux méthodes d'apprentissage</li> <li>- Maîtriser les paramètres permettant d'estimer la qualité de ces analyses</li> </ul>
<b>Public</b>	<p>Toute personne souhaitant comprendre le contexte d'utilisation, les concepts, et la mise en oeuvre des méthodes de classements et de prédiction - Ingénieurs et chargés d'études/d'analyses, Chercheurs, Doctorants, Statisticiens, Data scientists</p>
<b>Pré-requis</b>	<ul style="list-style-type: none"> <li>- Connaissances sur les outils statistiques de base : corrélation, écart-type, variance, intervalles de confiance, tests d'hypothèses.</li> <li>- Dans le cas où la formation serait effectuée avec le logiciel R, une connaissance de base de ce logiciel est préconisée</li> </ul> <p>Après votre inscription dans Sirène, compléter le questionnaire suivant :  <a href="https://forms.office.com/Pages/ResponsePage.aspx?id=SdMr0PDW80WLuFdUfKcVNHsI5D6av0tFtLyEn6XmsGhURE5aTzcyTJLSVg4WUdERV0o0QIAxSTdUQSQIQCNjPTEu">https://forms.office.com/Pages/ResponsePage.aspx?id=SdMr0PDW80WLuFdUfKcVNHsI5D6av0tFtLyEn6XmsGhURE5aTzcyTJLSVg4WUdERV0o0QIAxSTdUQSQIQCNjPTEu</a></p> <p><i>Si le lien ne fonctionne pas en cliquant dessus -&gt;copier/coller dans votre barre de recherche Internet</i></p> <p><b>Sans ce questionnaire complété et votre inscription dans Sirène validé par votre supérieur hiérarchique avant la date limite d'inscription, votre demande de formation ne sera pas étudiée</b></p>
<b>Programme</b>	<p><b>VOLET 1 : LES ALGORITHMES</b></p> <p><b>Généralités sur les différentes méthodes d'apprentissage supervisé</b></p> <ul style="list-style-type: none"> <li>- Différences entre méthodes supervisées et non supervisées</li> <li>- Objectifs de l'apprentissage supervisé : Objectifs de description - Objectifs de prédiction</li> <li>- Structure des jeux de données</li> <li>- Présentation générale de l'éventail des méthodes</li> </ul> <p><b>La méthode knn</b></p> <ul style="list-style-type: none"> <li>- Principe de la méthode des plus proches voisins</li> <li>- Algorithme de calcul</li> <li>- Distances entre individus</li> <li>- Choix des proximités</li> <li>- Définition du paramètre k</li> </ul> <p><b>La régression logistique</b></p> <ul style="list-style-type: none"> <li>- Variable explicative et variable expliquée (continue / binaire)</li> <li>- Différences entre la régression linéaire classique et la régression logistique</li> <li>- Variables explicatives qualitatives, variables explicatives quantitatives</li> <li>- Objectifs de la régression logistique</li> <li>- Définition du modèle Logit (courbe sigmoïde)</li> <li>- Conditions d'utilisation à respecter</li> <li>- Estimation et interprétation des coefficients du modèle</li> <li>- Test de significativité du modèle (validation du modèle)</li> <li>- Tests d'apport d'une variable (test de Wald, tests sur les rapports de vraisemblance)</li> <li>- Interprétation du <math>\chi^2</math> de Wald</li> <li>- Odds-ratios</li> <li>- Parallèle odds ratios et risques relatifs</li> <li>- Mise en oeuvre et analyse des résultats d'un modèle de régression logistique multiple</li> <li>- Estimation et interprétation des coefficients du modèle multiple</li> </ul>

### **L'analyse factorielle discriminante**

- Structure du jeu de données et contexte d'application
- Objectifs détaillés de l'AFD
- Notions de classement et de discrimination
- Méthodologie de l'AFD
- Comparaison avec l'ACP
- Interprétation des sorties logiciel : cercle factoriels, corrélations variables x axes
- Qualité de l'AFD (de la discrimination obtenue) : Tests univariés et multivariés (lambda de Wilks) - Graphique des individus - Matrice de confusion (et éventuellement courbe ROC)
- Les confusions et erreurs à ne pas commettre

### **Les supports vecteurs machines (SVM)**

- Démarche des svm
- Notions de marge
- Séparation linéaire
- Séparation non linéaire
- Fonction noyau

### **Les arbres de décision**

- Structure du jeu de données
- Principes, vocabulaire et objectifs
- Notion d'échantillon d'apprentissage, de validation et de test
- Comparaison de méthodes de type régression linéaire / logistique aux arbres de décision
- Principe de la segmentation selon le type de variable : Arbre de régression ou arbre de classification
- Définir les conditions d'arrêt de construction d'un arbre : Notion de pré-élagage
- Définition des groupes après construction de l'arbre
- Indicateurs de qualité
- Comparaison d'arbre de décision selon un certain type d'algorithme : CHAID vs CART
- Avantages et inconvénients : limites des arbres de décision
- Mise en oeuvre et interprétation des résultats obtenus après application d'une analyse par arbre de décision

### **De l'arbre à la forêt - Random Forest**

- Pourquoi avoir recourt aux forêts aléatoires ?
- Principes et objectifs o Instabilité de l'arbre : Notion de Bagging - Les erreurs liées à l'échantillonnage (Out-Of-Bag) - Prédiction avec un algorithme de Forêt aléatoire : Les paramètres
- Evaluer l'importance des variables : Notion d'importance - Comportement de l'importance - Lien entre diversité des arbres et l'importance - Influence des paramètres
- Sélection de variables : Généralités et principes de la sélection - Procédure de sélection - Les paramètres de sélection - Validation

### **VOLET 2 : VALIDATIONS DES METHODES, MESURE DES PERFORMANCES**

- Partitionnement des données disponibles : Jeu d'entraînement - Jeu de validation - Tests sur le jeu d'entraînement - Tests sur le jeu de validation
- La validation croisée : Leave one out : - K fold - Leave v out
- Compromis biais / variance
- Mesures des performances des classifications : Matrices de confusions - Courbe Roc - Aire sous la courbe (AUC) - Sensibilité & spécificité

#### **Dates**

Date : **Les 6, 7, 8 et 9 novembre 2023**  
Inscriptions avant **le 20 septembre 2023** sur <https://www.sirene.inserm.fr/>

#### **FOAD : formation organisée à distance**

Formation à distance sur votre lieu de travail ou à domicile.  
Au sein d'un groupe restreint, vous suivez une formation en direct avec un formateur dédié qui dispose d'une solution de visio-conférence  
Conditions pour pouvoir suivre la formation : Ordinateur avec micro (caméra non obligatoire), accès à internet, deuxième écran facilitateur mais non obligatoire.

#### **Contact**

Catherine Marcilhac  
INSERM – DR Paris IdF Sud  
@ : [formation.paris11@inserm.fr](mailto:formation.paris11@inserm.fr)