

# **Diplôme Universitaire Data Scientist**

Programme année 2022-2023

Université Clermont Auvergne

<b>UE 1</b>		<b>Data Analysis with Python</b>
<b>Présentation</b>	Description et objectifs	Cours de programmation dans le langage Python, orienté pour l'analyse de données.
	Compétences visées	<ul style="list-style-type: none"> <li>- programmer à l'aide du langage Python.</li> <li>- analyser des données efficacement avec numpy</li> <li>- représenter des données avec matplotlib</li> </ul>
	Pré-requis	<ul style="list-style-type: none"> <li>- connaissance de bases d'un langage informatique</li> <li>- connaissance de base de mathématiques et statistiques</li> </ul>
	Formation	Lieu : UCA, Campus Cézeaux Durée : 25h Date indicative: Septembre 2022 Formation continue : oui
<b>Contenu</b>	Programme détaillé	1) Introduction to data analysis tools <ul style="list-style-type: none"> <li>- generalities on Python language</li> <li>- basics on Unix and shell</li> <li>- data analysis with Anaconda and Notebooks</li> </ul> 2) Practical introduction <ul style="list-style-type: none"> <li>- object, collections, functions</li> <li>- loops and few pythonics syntax</li> <li>- basic file manipulation</li> </ul> 3) Numpy introduction <ul style="list-style-type: none"> <li>- numpy arrays vs python list</li> <li>- vectorization, (fancy) indexing, broadcasting</li> </ul> 4) Three important tools to know <ul style="list-style-type: none"> <li>- data representation: matplotlib (histo, 1D plot, 2D scatter)</li> <li>- import/manipulate data as numpy array: pandas (loading csv, accessing columns, plotting)</li> <li>- mathematics, physics and engineering: scipy (studied example: model fitting)</li> </ul> 5) Basic of image processing <ul style="list-style-type: none"> <li>- loading/plotting, colors, grey scale</li> <li>- image filters: kernel, blocks, sliding windows</li> </ul>
	Volume horaire	TP : 25h
	Équipe enseignante	Chercheur CNRS du Laboratoire de Physique de Clermont.
	Ressources et moyens	Salle informatique

<b>UE 2</b>		<b>Advanced Statistics</b>
<b>Présentation</b>	Description et objectifs	Cours et exercices pratiques permettant d'appréhender les concepts statistiques essentiels au 'data scientist'.
	Compétences visées	<ul style="list-style-type: none"> <li>- Maîtriser des concepts statistiques avancés.</li> <li>- Analyser des données à l'aide de méthodes numériques.</li> <li>- Appliquer des méthodes d'analyse prédictive.</li> <li>- Construire des modèles statistiques.</li> </ul>
	Pré-requis	Notions de probabilités. Connaissances mathématiques (niveau Licence).
	Formation	Lieu : UCA, Campus Cézeaux Durée : 20h Date indicative: Octobre 2022 Formation continue : oui
<b>Contenu</b>	Programme détaillé	1) Basics on theory of probabilities <ul style="list-style-type: none"> <li>- Properties of distributions</li> <li>- Random variables distributions</li> <li>- Samples and statistical properties of samples</li> </ul> 2) Likelihood <ul style="list-style-type: none"> <li>- Parameter estimation</li> <li>- Hypothesis testing</li> </ul> 3) Confidence level and confidence intervals 4) Bayesian inference
	Volume horaire	CM : 15h TD : 5h
	Équipe enseignante	Enseignant-chercheur UCA du Laboratoire de Physique de Clermont.
	Ressources et moyens	

<b>UE 3</b>		<b>Machine Learning</b>
<b>Présentation</b>	Description et objectifs	Introduction aux techniques multi-variables et aux méthodes d'apprentissage supervisé (Machine Learning).
	Compétences visées	<ul style="list-style-type: none"> <li>- Appliquer des méthodes de classification et de clustering.</li> <li>- Implémentation d'algorithmes d'apprentissage supervisé.</li> <li>- Concevoir des méthodes d'analyse prédictive et d'optimisation.</li> <li>- Utiliser un algorithme de Deep Learning.</li> </ul>
	Pré-requis	Notions de probabilités et de statistiques. Connaissances mathématiques (niveau Licence).
	Formation	Lieu : UCA, Campus Cézeaux Durée : 20h Date indicative : Décembre 2022 Formation continue : oui
<b>Contenu</b>	Programme détaillé	1) Basic concepts <ul style="list-style-type: none"> <li>- Types of learning: supervised, unsupervised</li> <li>- Regression, Classification, Clustering</li> <li>- Training and testing, cross validation</li> <li>- Bias-variance decomposition</li> <li>- Curse of dimensionality</li> </ul> 2) Regression with linear models <ul style="list-style-type: none"> <li>- Simple exemple: polynomial curve fitting</li> <li>- Linear basis function models</li> <li>- Regularization</li> <li>- Likelihood and regression</li> </ul> 3) Classification <ul style="list-style-type: none"> <li>- Linear discriminant analysis</li> <li>- Logistic regression and Gradient descent</li> <li>- Perceptron's algorithm</li> <li>- Towards Neural Networks</li> </ul> 4) Introduction to neural networks <ul style="list-style-type: none"> <li>- building blocks and architecture</li> <li>- Introduction to Pytorch</li> <li>- Deep learning practice with Pytorch</li> </ul>
	Volume horaire	CM : 10h TP : 10h
	Équipe enseignante	Enseignant-chercheur UCA du Laboratoire de Physique de Clermont.
	Ressources et moyens	Salle informatique

<b>UE 4</b>		<b>Data mining</b>
<b>Présentation</b>	Description et objectifs	Introduction aux méthodes de fouille de données et aux techniques d'apprentissage non supervisé.
	Compétences visées	<ul style="list-style-type: none"> <li>- Savoir évaluer les méthodes de fouille de données.</li> <li>- Savoir utiliser les techniques de fouille de données.</li> <li>- Maîtriser une technique de classification non supervisée.</li> <li>- Maîtriser une technique de recherche de règles d'association.</li> <li>- Connaître les différentes approches de fouille de données complexes.</li> </ul>
	Pré-requis	Notions de probabilités et de statistiques. Connaissances mathématiques (niveau Licence).
	Formation	Lieu : UCA, Campus Cézeaux Durée : 20h Date indicative : Novembre 2022 Formation continue : oui
<b>Contenu</b>	Programme détaillé	1) Introduction to big data and data mining <ul style="list-style-type: none"> <li>- data mining vs machine learning</li> <li>- challenges and criticisms</li> <li>- competence needed</li> </ul> 2) Preprocessing data <ul style="list-style-type: none"> <li>- data type and visualisation</li> <li>- data preprocessing scheme</li> <li>- cleaning methods</li> <li>- space transformation</li> <li>- dimensionality reduction</li> <li>- instance reduction</li> </ul> 3) Clustering <ul style="list-style-type: none"> <li>- hierarchical clustering</li> <li>- partitional clustering (k-means, dbscan)</li> <li>- clustering evaluation</li> </ul> 4) Feature reduction and extraction <ul style="list-style-type: none"> <li>- PCA, CA, MCA, MDS</li> <li>- kernel PCA</li> <li>- locally Linear Embedding</li> <li>- t-SNE</li> </ul>
	Volume horaire	TD : 25h
	Équipe enseignante	Eneignant-chercheur UCA du Laboratoire d'Informatique, de Modélisation et d'Optimisation des Systèmes.
	Ressources et moyens	Salle informatique

<b>UE 5</b>		<b>Data Engineering</b>
<b>Présentation</b>	Description et objectifs	Formation aux techniques modernes de calcul scientifique et de traitement des données dispensée par le Centre de Calcul de l'IN2P3 (cc.in2p3.fr).
	Compétences visées	<ul style="list-style-type: none"> <li>- Comprendre les problématiques liées au calcul (architectures modernes, calcul sur la grille et dans le cloud)</li> <li>- Savoir utiliser une grille de calcul.</li> <li>- Savoir choisir le stockage le plus adapté : système de fichiers, stockage dans le cloud et bases de données (SQL et NoSQL).</li> <li>- Programmation avec le langage Python,.</li> <li>- Utilisation d'un notebook de type Jupyter pour mener une analyse de données.</li> </ul>
	Pré-requis	Connaissance minimale d'un système d'exploitation (Linux de préférence) et de la ligne de commande. Connaissance minimale d'un langage de programmation (écriture du code, compilation et exécution)
	Formation	Lieu : Centre de Calcul de l'IN2P3 (Villeurbanne) Durée : 30h Date indicative : Mars 2023 Formation continue : oui
<b>Contenu</b>	Programme détaillé	1) Distributed data processing - introduction to Python language and notebooks - Using a Jupyter notebook  2) Computing - computation architectures - CPU and GPU computing, parallel computing - virtualization and distributed computing: cloud and grid  3) Storage - object storage in the cloud - relational and non-relational databases - HADOOP & SPARK ecosystem  4) Visit of the CCIN2P3 computer centre
	Volume horaire	CM : 15h TP : 15h
	Équipe enseignante	Ingénieurs CNRS du Centre de Calcul de l'IN2P3 (CCIN2P3).
	Ressources et moyens	Ressources informatiques du CCIN2P3.

<b>UE 6</b>		<b>Statistical tools : R language</b>
<b>Présentation</b>	Description et objectifs	Ce cours vise à se familiariser avec les fonctionnalités de programmation matricielles et objet du logiciel R.
	Compétences visées	<ul style="list-style-type: none"> <li>- Se familiariser avec le logiciel R</li> <li>- Savoir programmer sous R(fonctions et objets).</li> <li>- Outils de visualisation graphique</li> <li>- Connaître les interfaçages Python, Shiny, ggplots.</li> <li>- Utiliser des routines C ou Fortran.</li> </ul>
	Pré-requis	Connaissance en statistique
	Formation	Lieu : UCA, Campus Cézeaux Durée : 20h Date indicative : Avril 2023 Formation continue : oui
<b>Contenu</b>	Programme détaillé	1) R software and packages for data analysis - objects, structure, packages, functions and graphics  2) Statistics with R - basics, tests, regression, classification  3) Use case studies  4) Links R/python and R/C
	Volume horaire	CM:10h TP: 10h
	Équipe enseignante	Enseignant-chercheurs UCA du laboratoire de mathématiques Blaise Pascal.
	Ressources et moyens	Salle informatique

<b>UE 7</b>		<b>Time series</b>
<b>Présentation</b>	Description et objectifs	Il s'agit d'une introduction à la modélisation et la prévision de séries temporelles.
	Compétences visées	-Connaître les concepts de base des séries temporelles. -Connaître les principaux outils statistiques d'analyse d'une série temporelle. -Savoir modéliser des séries temporelles par les modèles classiques fondés sur des processus à temps discrets usuels et savoir en tirer des prévisions.
	Pré-requis	Connaissance des probabilités, de la statistique avancée, de régression linéaire.
	Formation	Lieu : UCA, Campus Cézeaux Durée : 20h Date indicative : Mai 2023 Formation continue : oui
<b>Contenu</b>	Programme détaillé	1) Trend modelling and seasonality, seasonal adjustment. 2) Main properties of Fourier analysis, periodogram 3) Modeling and forecasting by ARMA process. 4) Elements on the ARIMA and SARIMA models 5) Applications with the R software.
	Volume horaire	CM:8h TD: 4h TP : 8h
	Équipe enseignante	Enseignant-chercheurs UCA du laboratoire de mathématiques Blaise Pascal.
	Ressources et moyens	Salle informatique



UE 8		Deep Learning
Présentation	Description et objectifs	Introduction aux méthodes d'apprentissage profond pour la classification, la régression, la réduction de dimensionalité et l'apprentissage de métriques. Introduction aux réseaux de neurones récurrents
	Compétences visées	- Concevoir et entraîner un réseau de neurones profond supervisé. Perceptron multi-couches, DCNN, autoencodeur - Savoir implémenter un réseau de neurones sur une architecture GPUs. - Savoir analyser et évaluer les performances d'un entraînement.
	Pré-requis	Connaissances mathématiques (niveau Licence). Notions de probabilités, de statistiques et d'optimisation. Notions de base en Machine Learning et réseaux de neurones. Programmation Python.
	Formation	Lieu : UCA, Campus Cézeaux Durée : 20h Date indicative : Janvier/Février 2023 Formation continue : oui
Contenu	Programme détaillé	1) Introduction aux réseaux de neurones profonds - architecture générale, fonctions d'activation (Relu, ...), softmax. - réseaux convolutifs (DCNN)  2) Entraînement d'un réseau - fonctions de coûts pour la classification et la régression - techniques d'optimisation et de régularisation, - stratégies d'entraînement pour limiter le surapprentissage  3) Architectures de réseaux pour l'analyse d'images : classification, détection, segmentation, réseaux siamois  4) Autoencodeurs - Principe et exemples d'utilisations : réduction de dimension, débruitage, génération  5) Réseaux de neurones récurrents - Structure générale et entraînement - Architectures GRU, LSTM  <u>Deep learning practice with Pytorch</u> 1) Réseaux de neurones pour la classification, influence des hyperparamètres et régularisation 2) Auto-encodeurs et application au débruitage d'image 3) Utilisation d'un réseau existant et transfert d'apprentissage 4) Réseau de neurones récurrent et prédiction de séquence
	Volume horaire	CM : 10h TP : 10h

