

## Data Scientist

Informatique

Notions fondamentales | OLAP /  
Décisionnel | Big Data

Référence : 4-IT-DAS

Durée : 5 jours

Présentiel ou en classe à distance

Mise à jour : 27/11/2023

Tarif Inter : 750 € Prix HT jour / personne

Tarif Intra : 1500 € Prix HT jour / groupe

Durée de validité : du 01/01/2026 au 31/12/2026

### Objectifs

- Savoir mettre en place un DataLake et un DataMart en SQL ou big data
- Savoir mettre en place une stratégie de Machine Learning en Python afin de créer le modèle le plus satisfaisant possible en le mesurant et en affichant les résultats, le tout, en utilisant des algorithmes performants

### Prérequis

- Maîtriser l'algorithmique, avoir une appétence pour les mathématiques
- La connaissance de Python et des statistiques est un plus

### Public concerné

Développeurs, chefs de projets proches du développement, ingénieurs scientifiques sachant coder

### Contenu pédagogique

#### Introduction aux Data Sciences

- Qu'est que la data science ?
- Qu'est-ce que Python ?
- Qu'est que le Machine Learning ?
- Apprentissage supervisé vs non supervisé
- Les statistiques
- La randomisation
- La loi normale

#### Introduction à Python pour les Data Science

- Les bases de Python
- Les listes
- Les tuples
- Les dictionnaires
- Les modules et packages
- L'orienté objet
- Le module math
- Les expressions lambda
- Map, reduce et filter
- Le module CSV
- Les modules DB-API 2 Anaconda

#### Introduction aux DataLake, DataMart et DataWarehouse

- Qu'est-ce qu'un DataLake ?
- Les différents types de DataLake
- Le Big Data

- Qu'est-ce qu'un DataWarehouse ?
- Qu'est ce qu'un DataMart ?
- Mise en place d'un DataMart
- Les fichiers
- Les bases de données SQL
- Les bases de données No-SQL

### Python Package Installer

- Utilisation de PIP
- Installation de package PIP PyPi

### MathPlotLib

- Utilisation de la bibliothèque scientifique de graphes MathPlotLib
- Affichage de données dans un graphique 2D
- Affichages de sous-graphes
- Affichage de polynômes et de sinusoïdales

### Machine Learning

- Mise en place d'une machine learning supervisé
- Qu'est qu'un modèle et un dataset
- Qu'est qu'une régression
- Les différents types de régression
- La régression linéaire
- Gestion du risque et des erreurs
- Quarter d'Ascombe
- Trouver le bon modèle
- La classification
- Loi normale, variance et écart type
- Apprentissage
- Mesure de la performance No Fee Lunch

### La régression linéaire en Python

- Programmer une régression linéaire en Python
- Utilisation des expressions lambda et des listes en intention
- Afficher la régression avec MathPlotLib
- L'erreur quadratique
- La variance
- Le risque

### Le Big Data

- Qu'est-ce que Apache Hadoop ?
- Qu'est-ce que l'informatique distribué ?
- Installation et configuration de Hadoop
- HDFS
- Création d'un datanode
- Création d'un namenode distribué
- Manipulation de HDFS
- Hadoop comme DataLake
- Map Reduce
- Hive

- Hadoop comme DataMart
- Python HDFS

### Les bases de données NoSql

- Les bases de données structurées
- SQL avec SQLite et Postgresql
- Les bases de données non ACID
- JSON
- MongoDB
- Cassandra, Redis, CouchDb
- MongoDB sur HDFS
- MongoDB comme DataMart PyMongo

### Numpy et SciPy

- Les tableaux et les matrices
- L'algèbre linéaire avec Numpy
- La régression linéaire SciPy
- Le produit et la transposée
- L'inversion de matrice
- Les nombres complexes
- L'algèbre complexe
- Les transformées de Fourier Numpy et Mathplotlib

### ScikitLearn

- Régressions polynomiales
- La régression linéaire
- La création du modèle
- L'échantillonnage
- La randomisation
- L'apprentissage avec fit
- La prédiction du modèle
- Les metrics
- Choix du modèle
- PreProcessing et Pipeline
- Régressions non polynomiales

### Nearest Neighbors

- Algorithme des k plus proches voisins (k-NN)
- Modèle de classification
- K-NN avec SciKitLearn
- Choix du meilleur k
- Sérialisation du modèle
- Variance vs Erreurs
- Autres modèles : SVN, Random Forest

### Pandas

- L'analyse des données avec Pandas
- Les Series
- Les DataFrames
- La théorie ensembliste avec Pandas

- L'importation des données CSV
- L'importation de données SQL
- L'importation de données MongoDB Pandas et SKLearn

## Le Clustering

- Regroupement des données par clusterisation
- Les clusters SKLearn avec k-means
- Autres modèles de clusterisation : AffinityPropagation, MeanShift, ...
- L'apprentissage semi-supervisé

## Jupyter

- Présentation de Jupyter et Ipython
- Installation
- Utilisation de Jupyter avec Mathplotlib et Sklearn

## Python Yield

- La programmation efficace en Python
- Le générateurs et itérateurs
- Le Yield return
- Le Yield avec Db-API 2, Pandas et Sklearn

## Les réseaux neuronaux

- Le perceptron
- Les réseaux neuronaux
- Les réseaux neuronaux supervisés
- Les réseaux neuronaux semi-supervisés
- Les réseaux neuronaux par Hadoop Yarn
- Les heuristiques
- Le deep learning

---

## Moyens pédagogiques

- Réflexion de groupe et apports théoriques du formateur.
- Travail d'échange avec les apprenants sous forme de réunion - discussion.
- Utilisation de cas concrets issus de l'expérience professionnelle.
- Validation des acquis par des questionnaires, des tests d'évaluation, des mises en situation et des jeux pédagogiques.
- Alternance entre apports théoriques et exercices pratiques (en moyenne sur 30 à 50% du temps)

**Modalités pédagogiques :** Présentiel, Distanciel et AFEST

## Moyens techniques

### En formation présentielle

Accueil des apprenants dans une salle dédiée à la formation et équipée avec :

- Ordinateurs
- Vidéo projecteur ou Écran TV interactif
- Tableau blanc ou Paper-Board

### En formation distancielle

A l'aide d'un logiciel comme ® Microsoft Teams ou Zoom, un micro et une caméra pour l'apprenant.

- Suivez une formation en temps réel et entièrement à distance. Lors de la session en ligne, les apprenants interagissent et communiquent entre eux et avec le formateur.
- Les formations en distanciel sont organisées en Inter-Entreprise comme en Intra-Entreprise.
- L'accès à l'environnement d'apprentissage (support de cours, ressources formateur, fichiers d'exercices ...) ainsi qu'aux preuves de suivi et d'assiduité (émergence, évaluation) est assuré.
- Les participants recevront une convocation avec le lien de connexion à la session de formation.
- Pour toute question avant et pendant le parcours, une assistance technique et pédagogique est à disposition par téléphone au 02 35 12 25 55 ou par email à [commercial@xxlformation.com](mailto:commercial@xxlformation.com)

## Modalités d'évaluation

- Positionnement préalable oral ou écrit.
- Feuille de présence signée en demi-journée.
- Evaluation des acquis tout au long de la formation.

- Questionnaire de satisfaction
- Attestation de stage à chaque apprenant
- Evaluation formative tout au long de la formation.
- Evaluation sommative faite par le formateur.

**Profil du formateur**

- Nos formateurs sont des experts dans leurs domaines d'intervention
- Leur expérience de terrain et leurs qualités pédagogiques constituent un gage de qualité

**Adaptation pédagogique et matérielle**

Si vous avez besoin d'adaptation matérielle ou pédagogique, merci de prendre contact avec notre référent Handicap par téléphone au 02 35 12 25 55 ou par email à [handicap@xxlformation.com](mailto:handicap@xxlformation.com)

**Modalités et délais d'accès à la formation**

Les formations sont disponibles selon les modalités proposées sur la page programme. Les inscriptions sont possibles jusqu'à 48 heures ouvrées avant le début de la formation. Dans le cas d'une formation financée par le CPF, ce délai est porté à 11 jours ouvrés.

**Nos sessions INTER 2026**

Sessions de formation à venir :

- Aucune session à venir pour cette formation.

**Nos sessions INTRA 2026**

Pour organiser cette formation en intra-entreprise, veuillez nous contacter par mail à [commercial@xxlformation.com](mailto:commercial@xxlformation.com) ou par téléphone au 02 35 12 25 55