

Référence	4-IT-DAS
Durée	5 jours (35 heures)
Éligible CPF	NON
Mise à jour	27/11/2023

Data Scientist



OBJECTIFS PÉDAGOGIQUES

- 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



PUBLIC CONCERNÉ

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



PRÉREQUIS

Maîtriser l'algorithmique, avoir une appétence pour les mathématiques

La connaissance de Python et des statistiques est un plus



MOYENS PÉDAGOGIQUES

- Réflexion de groupe et apports théoriques du formateur
- Travail d'échange avec les participants 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.
- Remise d'un support de cours.



MODALITÉS D'ÉVALUATION

- 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,
- Positionnement préalable oral ou écrit,
- Evaluation formative tout au long de la formation,
- Evaluation sommative faite par le formateur ou à l'aide des certifications disponibles



MOYENS TECHNIQUES EN PRÉSENTIEL

Accueil des stagiaires dans une salle dédiée à la formation équipée à minima d'un vidéo projecteur et d'un tableau blanc et/ou paperboard.

Pour les formations nécessitant un ordinateur, un PC est mis à disposition de chaque participant.



MOYENS TECHNIQUES EN DISTANCIEL

A l'aide d'un logiciel (Teams, Zoom...), d'un micro et éventuellement d'une caméra les apprenants interagissent et communiquent entre eux et avec le formateur.

Sessions organisées en inter comme en intra entreprise.

L'accès à l'environnement d'apprentissage ainsi qu'aux preuves de suivi et d'assiduité (émargement, évaluation) est assuré.

Pour toute question avant et pendant le parcours, assistance technique à disposition au 04 67 13 45 45.



ORGANISATION

Délai d'accès : 5 jours ouvrés
 (délai variable en fonction du financeur)

Les cours ont lieu de 9h à 12h30 et de 13h30 à 17h



ACCESSIBILITÉ

Les personnes en situation d'handicap sont invitées à nous contacter directement, afin d'étudier ensemble les possibilités de suivre la formation.

Pour tout renseignement, notre référent handicap reste à votre disposition : mteyssedou@ait.fr



PROFIL FORMATEUR

Nos formateurs sont des experts dans leurs domaines d'intervention

Leur expérience de terrain et leurs qualités pédagogiques constituent un gage de qualité.



CERTIFICATION POSSIBLE

Aucune

Data Scientist

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 DATA LAKE, DATA MART ET

DATA WAREHOUSE

- Qu'est-ce qu'un Data Lake ?
- Les différents types de Data Lake
- Le Big Data
- Qu'est-ce qu'un Data Warehouse ?
- Qu'est qu'un Data Mart ?
- Mise en place d'un Data Mart
- 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 sinusoidales

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 Data Lake
- Map Reduce
- Hive
- Hadoop comme Data Mart
- 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 Data Mart 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

SCIKIT LEARN

- 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