

## Examen TP : Big Data

### 1. Partie 1 : Installation

- Installer : *jdk-8, spark-cluster, findspark, pyspark*
- Importez [*findspark, pyspark*] et initialisez *SparkContext*

### 2. Partie 2 : SPARK SQL

- Charger les DataSets *dataset1.csv* et *dataset2.csv* dans deux dataFrames *d1\_df* et *d2\_df*, respectivement.
- Afficher les deux dataFrames.
- Nettoyer les deux dataFrames en gardant uniquement ceux ils n'ont pas d'attributs manquants.
- Afficher des statistiques de l'attribut **BMD** (*bone mineral density*) du deuxième dataFrame *d2\_df*.
- Afficher les différentes valeurs distinctes pour les attributs « *sex, medication* » du dataFrame *d1\_df*.
- Transformer les valeurs des attributs « *sex, medication* » en valeurs numériques.

### 3. Partie 3 : SPARK MLlib

On souhaite créer un modèle de régression qui prédit le «**BMD**» des patients en se basant sur les features : [**age, sex, medication, BMI**], où le **BMI**=  $[\text{weight}/(\text{height})^2]$ , tel que, le **weight** en kg et **height** en mètres

- A partir des deux dataFrames *d1\_df* et *d2\_df* créer le dataset correspondant.
- Splitter le dataset en *train\_df* et *test\_df* [80% , 20%]
- Trainer le modèle
- Afficher les coefficients et les évaluations (*erreurs*)
- Prédire le BMD de certains nouveaux patients