



# Projet Hadoop Big Data (Data Sciences)

Description du projet  
Décembre 2023

Par Christophe GERMAIN

# Description

# Projet Big Data

## Technologies

- Hadoop + python (HappyBase ...)
- Python Pandas
- Suggestions :

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

# Projet Big Data

Objectifs :

- Le groupe doit livrer :
  - Un ensemble d'applications Big Data et Power BI
  - Un dossier comprenant :
    - L'analyse de la compréhension de la problématique
    - Des données qualifiées
    - Des procédures d'import des données
    - Des procédures de structuration
    - Des algorithmes d'analyse des données
    - Vos recommandations par rapport au déroulement du projet

# Projet Big Data

Le projet :

- A partir du fichier csv : **dataaw\_fro.csv**
- Format du fichier :

Options spécifiques au format :

Colonnes séparées par :

;

Colonnes entourées par :

"

Colonnes échappées avec :

"

Lignes terminées par :

AUTO

Remplacer NULL par :

NULL

☐ Retirer les caractères de fin de ligne à l'intérieur des colonnes

☒ Afficher les noms de colonnes en première ligne

# Projet Big Data

Le projet (suite) :

- Entête du fichier :

	#	Nom	Type	Interclassement	Attributs	Null	Valeur par défaut
<input type="checkbox"/>	1	<b>codcli</b>	int(11)			Non	<i>Aucun(e)</i>
<input type="checkbox"/>	2	<b>genrecli</b>	varchar(8)	utf8mb4_general_ci		Oui	<i>NULL</i>
<input type="checkbox"/>	3	<b>nomcli</b>	varchar(40)	utf8mb4_general_ci		Oui	<i>NULL</i>
<input type="checkbox"/>	4	<b>prenomcli</b>	varchar(30)	utf8mb4_general_ci		Oui	<i>NULL</i>
<input type="checkbox"/>	5	<b>cpcli</b>	varchar(5)	utf8mb4_general_ci		Oui	<i>NULL</i>
<input type="checkbox"/>	6	<b>villecli</b>	varchar(50)	utf8mb4_general_ci		Oui	<i>NULL</i>
<input type="checkbox"/>	7	<b>codcde</b>	int(11)			Non	<i>Aucun(e)</i>
<input type="checkbox"/>	8	<b>datcde</b>	datetime			Oui	<i>NULL</i>
<input type="checkbox"/>	9	<b>timbrecli</b>	float			Oui	<i>NULL</i>
<input type="checkbox"/>	10	<b>timbrecede</b>	float			Oui	<i>NULL</i>
<input type="checkbox"/>	11	<b>Nbcolis</b>	tinyint(4)			Oui	<i>NULL</i>
<input type="checkbox"/>	12	<b>cheqcli</b>	float			Oui	<i>NULL</i>
<input type="checkbox"/>	13	<b>barchive</b>	bit(1)			Oui	<i>NULL</i>
<input type="checkbox"/>	14	<b>bstock</b>	bit(1)			Oui	<i>NULL</i>
<input type="checkbox"/>	15	<b>codobj</b>	int(11)			Oui	<i>NULL</i>
<input type="checkbox"/>	16	<b>qte</b>	smallint(6)			Oui	<i>NULL</i>

# Projet Big Data

Le projet (suite) :

- Entête du fichier (suite) :

<input type="checkbox"/>	17	<b>Colis</b>	int(11)	Oui	<i>NULL</i>
<input type="checkbox"/>	18	<b>libobj</b>	varchar(50) utf8mb4_general_ci	Oui	<i>NULL</i>
<input type="checkbox"/>	19	<b>Tailleobj</b>	varchar(50) utf8mb4_general_ci	Oui	<i>NULL</i>
<input type="checkbox"/>	20	<b>Poidsobj</b>	double	Oui	<i>NULL</i>
<input type="checkbox"/>	21	<b>points</b>	int(11)	Oui	<i>NULL</i>
<input type="checkbox"/>	22	<b>indispobj</b>	bit(1)	Oui	<i>NULL</i>
<input type="checkbox"/>	23	<b>libcondit</b>	varchar(50) utf8mb4_general_ci	Oui	<i>NULL</i>
<input type="checkbox"/>	24	<b>prixcond</b>	double	Oui	<i>NULL</i>
<input type="checkbox"/>	25	<b>puobj</b>	double	Oui	<i>NULL</i>

# Projet Big Data

## Le projet (suite) : LOT 1

- Contexte :
  - Une Fromagerie (le client) a un datawarehouse depuis 2004 qui est représenté par le fichier csv fournit dans ce document.
  - Créer des jobs pour limiter le flux d'information (Mapper-Reducer) pour obtenir uniquement les informations voulues pour répondre au besoin du client décrit ci-dessous :
  - Le client désire les statistiques suivantes :
    1. Filtrer les données selon les critères suivants :  
Entre 2006 et 2010,  
Avec uniquement les départements 53, 61 et 28
    2. A partir du point 1 : Ressortir dans un tableau des 100 meilleures commandes avec la ville, la somme des quantités des articles et la valeur de « timbrece » (la notion de meilleures commandes : la somme des quantités la plus grande ainsi que le plus grand nombre de « timbrece » )
    3. Exporter le résultat dans un fichier Excel.



# Projet Big Data

Le projet (suite) : LOT 2

- Contexte :
  - (Comme le LOT 1)
  - Le client désire les statistiques suivantes :
    1. Filtrer les données selon les critères suivants :  
Entre 2011 et 2016,  
Avec uniquement les départements 22, 49 et 53
    2. A partir du point 1 : Ressortir de façon aléatoire de 5% des 100 meilleures commandes avec la ville, la somme des quantités des articles sans « timbrecli » (le timbrecli non renseigné ou à 0) avec la moyenne des quantités de chaque commande)  
Avoir un PDF avec un graphe (PIE) (par Ville)
    3. Exporter le résultat dans un fichier Excel.

# Projet Big Data

Le projet (suite ) : LOT 3

- Mettre en place une base NoSQL HBASE pour stocker le contenu du fichier CSV et de mettre en œuvre un moteur de recherche avec Power BI pour interroger ce Data Warehouse.
  - Pour répondre au Lot 1 et Lot 2 au niveau des résultats avec les graphes,
  - Mise en place d'un Dashboard interactif

# Projet Big Data

Liens :

- [Python\\_Complet](#)
- [https://pandas.pydata.org/docs/getting\\_started/index.html](https://pandas.pydata.org/docs/getting_started/index.html)