



# Projet Hadoop Big Data (Data Sciences)

Description du projet  
SEPTEMBRE 2024

Par Christophe GERMAIN

## Description

# Projet Big Data

## Table des matières

I. Technologies .....	4
II. Objectifs : .....	5
III. Le projet : .....	6
A. A partir du fichier csv : .....	6
B. Format du fichier : .....	6
IV. Le projet : LOT 1 .....	9
V. Lot2 : HBase et Moteur de Recherche : .....	10
A. Importation HBase : .....	10
1. Contraintes : .....	10
B. Choisir Power BI ou ELK : .....	10
VI. Liens : .....	11

## I. Technologies

- Hadoop + python (HappyBase ...)
- Python Pandas
- Suggestions :
  - `import numpy as np`
  - `import pandas as pd`
  - `import matplotlib.pyplot as plt`

# Projet Big Data

## II. Objectifs :

- Le groupe doit livrer :
- Un ensemble d'applications Big Data et Power BI / ELK
- 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

## III. Le projet :

### A. A partir du fichier csv :

dataw\_fro.csv fourni dans le dossier du projet

### B. 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

	#	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>timbrecde</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

<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

## IV. Le projet : LOT 1

- Contexte :
  1. Une Fromagerie (le client) a un datawarehouse depuis 2004 qui est représenté par le fichier csv fournit dans ce document.
  2. Créer des jobs (job.sh avec l'utilisation hadoop jar ...) 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 2008 et 2012,

Avec uniquement les départements 53, 61,75 et 28
  2. A partir du point 1 : Ressortir dans un tableau les 10 Clients les plus fidèles (Sommes du produit (des points et des quantités) sur l'ensemble des commandes par client) : récupérer les colonnes suivantes : Nom, Prénom, Ville, département du client et, nom de l'objet & la quantité commandée
  3. Exporter le résultat dans un fichier Excel et les 10 graphes (en pdf) par client avec le % de répartitions des objets commandés (produits commandés).

# Projet Big Data

## V. Lot2 : HBase et Moteur de Recherche :

### A. Importation HBase :

Importer toutes les données du fichier csv dans HBase

#### 1. Contraintes :

- Ne pas mettre dans le JSON des données à NULL,
- Date invalide interdite (enregistrement complet refusé),
- Ne pas importer l'année 2004 (enregistrement complet refusé)

### B. Choisir Power BI ou ELK :

Créer un Dashboard afin de manière dynamique de :

- Voir la fidélité des clients sur selon l'intervalle de dates,
- Voir le nombre d'objets commandés par année et les 5 meilleurs (Palmarès),
- Voir les départements les plus représentatifs du programme de fidélité de la fromagerie (Palmarès).

# Projet Big Data

## VI. Liens :

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