

**IA School**  
M2 IA et Management

# Projet Big Data

Mehdi KARECH

---

Creation d'une pipeline de données pour la prévision de la météo

---

Superviseur :

Karim KOUKI

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Architecture de la solution</b>	<b>2</b>
2.1	Ingestion de données . . . . .	2
2.2	Traitement de données . . . . .	3
2.3	Visualisation . . . . .	3
2.4	Partie ML . . . . .	3
2.5	Orchestration . . . . .	3
<b>3</b>	<b>Conclusion</b>	<b>5</b>

## CHAPTER 1

# Introduction

---

Dans ce projet nous allons nous intéresser à faire la prévision de la météo pour la ville de Paris (savoir s'il va pleuvoir par exemple) en implémentant une pipeline de données qui récupère ces dernières, en utilisant l'api du site (<https://openweathermap.org/>), l'api permet de récupérer des données tel que la température maximale, minimale et moyenne, la pression, la vitesse du vent, l'humidité et enfin la météo du jour, la récupération est faite sur la dernière année (on peut pas récupérer plus sans faire un abonnement) sur des intervalles de points d'une heure.

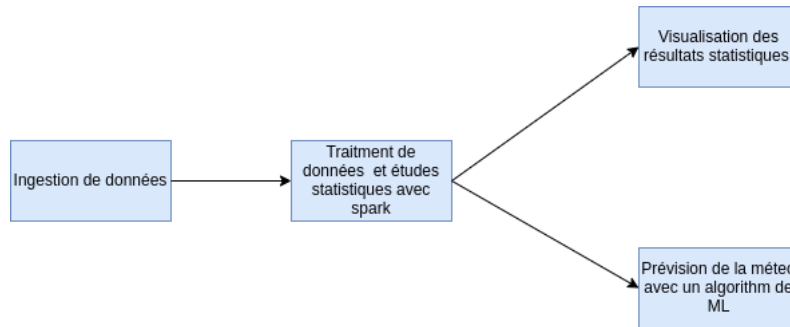


Figure 1.1: Pipeline général.

### **Ingestion de données :**

Pour l'ingestion de données nous allons utiliser un producer kafka qui va récupérer les données en utilisant l'api (openweather map), puis un kafka consumer va construire un fichier csv à partir des données injectées.

### **Traitement de données et études statistique avec Spark :**

Pour l'étude statistique du data set construits nous allons utiliser un script python avec pyspark

**Visualisation des résultats statistiques** Un autre script s'occupera de la visualisation des données qui sera dynamique selon les données injectées.

**Partie ML :** Une fois le dataset fini et qu'on aura récupéré toutes les données un modèle de Machine learning (une régression logistique) sera lancé sur le dataset finale, pour permettre la prévision de la météo.

## CHAPTER 2

# Architecture de la solution

---

Dans ce chapitre nous allons présenter en détails chaque composants de notre pipeline

## 2.1 Ingestion de données

Pour l'ingestion de données nous avons utilisé l'api d'open weather pour récupérer les données, l'api permet seulement de récupérer les données de la derniere année (les douzes derniers mois), la récupération se fait avec des intervalles d'une heure.

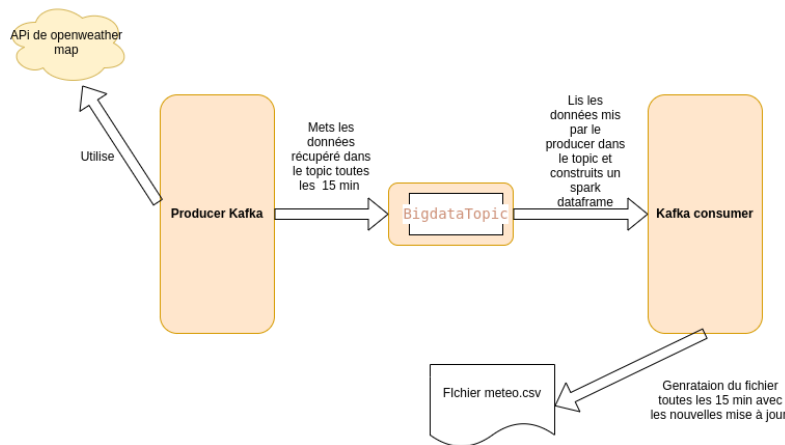


Figure 2.1: Ingestion de données

**configuration de kafka :** Pour utiliser notre architecture nous devons d'abord configurer kafka, pour commencer nous lançons le serveur zookeeper

```
$ sudo systemctl start zookeeper
```

Figure 2.2: Lancement de zookeeper

Puis Kafka :

```
mehdi@mehdi-G5-5587:~$ sudo systemctl start kafka
```

Figure 2.3: Lancement de kafka

une fois le serveur kafka lancé on crée le topic **BigDataTopic** qui servira de canal de communication entre le producteur kafka et le consumer kafka.

Une fois le topic créé le **Producer kafka est lancé**, ce dernier est implémenter dans le fichier **producer.py**. Dans ce script les données sont récupérés en utilisant l'api cité précédement, le script est implémenter de sorte qu'on récupère les données d'un mois toutes les 15 minutes et on les mets sous format **.json** dans le topic **BigDataTopic** les données acquies sont (time stamp, température maximale, température minimale, température moyenne, pression, humidité, vitesse de vent, et météo)

Après le lancement du script **producer.py**, le script **consumer.py** est exécuté, ce dernier récupère les objets **.json** à partir de **BigDataTopic**, construit un dataframe et crée un fichier **meteo.csv** toute les 15 min (l'ancien fichier est écrasé et remplacé par le nouveau fichier et le dataframe est mis à jour à chaque fois)

## 2.2 Traitement de données

Une fois le premier fichier **meteo.csv** créé, le script **traitement.py** est lancé, dans ce dernier des statistiques sont effectuées sur chaque colonne du dataset, et le script est réexécuté toute les 15 minutes pour prendre en compte la nouvelle mise à jour des données, le script génère un fichier **stats.csv** qui servira à faire la visualisation dans l'étape suivante

## 2.3 Visualisation

Notre solution a été développée sous linux, et on voulait connecter **Tableau** avec **spark**, On a eu des erreurs d'intégration qu'on a pas pu malheureusement régler à temps, on a décidé donc de faire nos visualisations en utilisant **matplotlib** et **seaborn** en intégrant dans un script **visualisation.py**, le script de visualisation est également exécuté toutes les 15 min une fois le premier fichier **stats.csv** généré, on peut trouver quelques visualisations ci-dessous

## 2.4 Partie ML

Une fois toutes les données obtenues (un dataset final est obtenu avec toutes les données sur une année), le script **ML.py** est lancé manuellement, le script applique une régression logistique sur le dataset final obtenu, pour donner la météo à une heure donnée (4 classes possibles, tous les attributs sont numériques donc pas de one hot encoding à effectuer), nous avons utilisé **mlib** pour entraîner notre modèle sur 80% des données et le tester sur les 20% restantes, et on a obtenu une accuracy de 67%

## 2.5 Orchestration

Pour l'orchestration des données, et vu qu'on a utilisé linux comme environnement de développement, on a utilisé **cron** pour lancer les scripts toutes les 15 minutes, seuls les

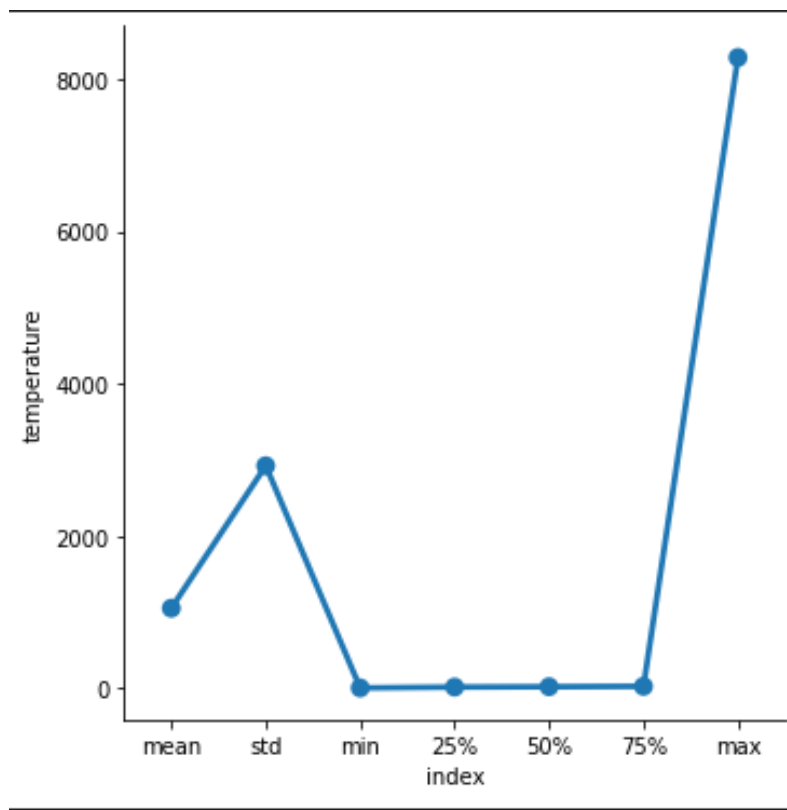


Figure 2.4: Distribution de la population température

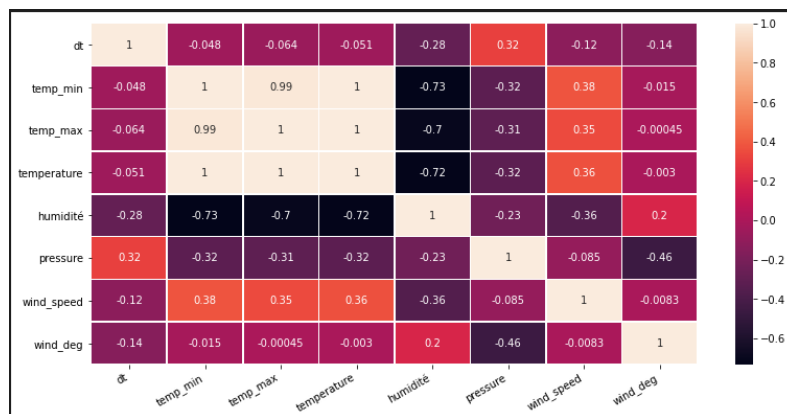


Figure 2.5: Coorelation des données

scripts **producer.py**, **consumer.py** et **ML.py** son exécuter manuellement.

Une meilleure manière de faire serait d'intégrer Airflow ou Luigi pour orchestrer tout le processsus.( j'ai pas eu le temps de chercher plus par rapport à cet aspect vu que j'ai travaillé tout seul)

## CHAPTER 3

# Conclusion

---

Ce projet a été très benifique pour moi et j'ai vraiment appris beaucoup de choses, j'aurais aimé pouvoir intégrer tableau et Airflow mais malheureusement j'ai pas eu le temps pour le faire, je vous remercie pour vos effortx tout au long du cours et je vous souhaite une très bonne continuation pour la suite.