

Rapport – Data Mining



Enseignant : SOUIDI Mohamed

FERREIRA Alicia – KHALIFA Marina – NEJMI Manal

ITS 2

Année 2021-2022

Introduction

Ce projet consiste à développer une application Python qui consiste à un tableau de bord interactifs en temps réel sur les flux de données et celles-ci seront visualisées sur Superset.

Host est notre machine Windows physique qui héberge différentes machines virtuelles. Chaque logiciel possède une machine virtuelle, donc en tout nous aurons 4 machines virtuelles.

Dans ce projet, nous allons réaliser différentes étapes :

- Nous allons dans un premier temps récupérer les données que nous souhaitons traiter
- Nous allons écrire un programme python qui récupère les données et qui les envoie sur Kafka
- Spark récupère les données sur Kafka pour les traiter puis les sauvegarde sur MySQL
- MySQL
- Superset se connecte à MySQL pour visualiser les données

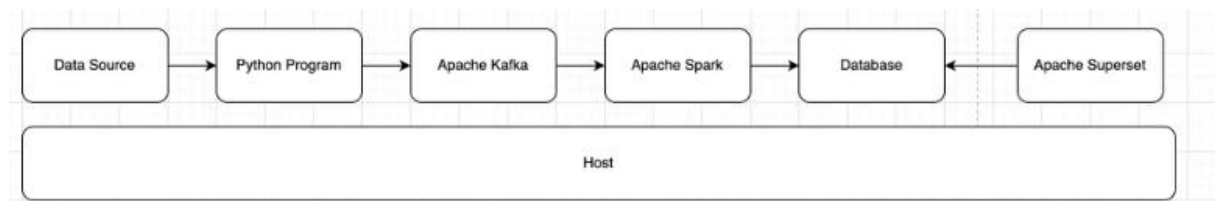


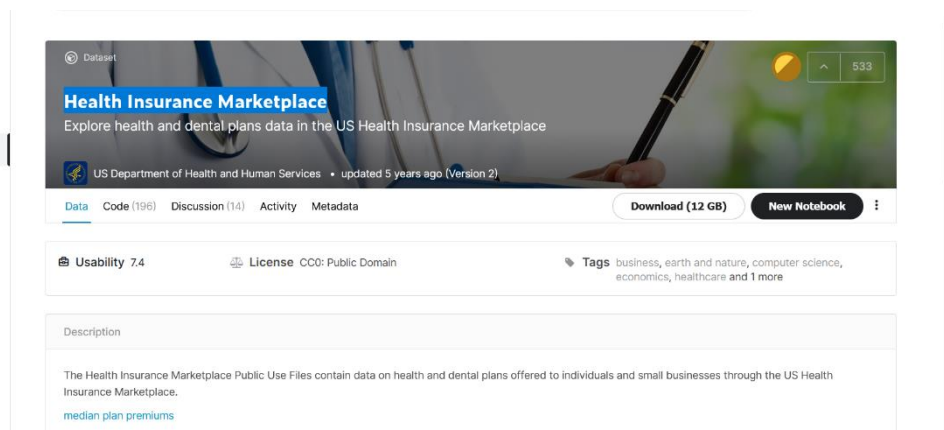
Figure 1 : Use Case du projet

Nous allons attribuer une adresse IP à chaque machine virtuelle afin que celles-ci puissent communiquer entre elles pour visualiser les données.

Data source

Nous avons prévu de traiter des données avec un fichier csv. Nous allons donc récupérer des données sur le site Kaggle sur Health Insurance Marketplace.

Nous avons commencé à créer notre code python qui permet de récupérer les données et les streamers. Ce code sera similaire à celui de wordcount.py vu en classe.



Configuration des machines

Nous avons configuré 4 machines virtuelles pour chaque logiciel : Spark, Kafka, MySQL et Superset.

En effet, nous avons configuré une adresse IP différente pour chaque VM via un vagrantfile.

Kafka

Kafka est une plateforme de streaming d'évènements qui capture des données en temps réel tel que des bases de données.

Kafka est un bus qui permet :

- Publier et s'abonner à des flux d'évènements en comprenant l'importation et l'exportation continue de nos données à partir d'autres systèmes
- Stocker les flux d'évènements de manière durable et fiable
- Traiter les flux d'évènements au fur et à mesure qu'ils se produisent

Kafka est reconnu comme Producer c'est-à-dire que c'est un client qui écrit des évènements.

Nous avons attribué à Kafka l'adresse IP suivante : 192.168.33.13, puis nous avons installé et lancé les deux serveurs zookeeper et Kafka, finalement on a créé un topic avec les commandes vu durant le cours.

Au niveau de cette partie, nous avons envoyé des évènements (messages) via un terminal qui est le Producer et nous avons pu les visualiser via à un autre terminal qui agit en tant que consumer.

Producer sur Kafka

```
PS C:\kafka> vagrant ssh
Welcome to Ubuntu 18.04.5 LTS (GNU/Linux 4.15.0-151-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:        https://ubuntu.com/advantage

System information as of Sun Jan  9 16:51:45 UTC 2022

System load:  0.5          Processes:    106
Usage of /:   3.4% of 61.80GB Users logged in: 1
Memory usage: 79%         IP address for eth0: 10.0.2.15
Swap usage:   2%          IP address for eth1: 192.168.33.13

This system is built by the Bento project by Chef Software
More information can be found at https://github.com/chef/bento
Last login: Sun Jan  9 14:01:50 2022 from 10.0.2.2
vagrant@vagrant:~$ bin/kafka-console-producer.sh --bootstrap-server 192.168.33.13:9092
-bash: bin/kafka-console-producer.sh: No such file or directory
vagrant@vagrant:~$ cd kafka
vagrant@vagrant:~/kafka$ bin/kafka-console-producer.sh --topic quickstart-events --bootstrap-server 192.168.33.13:9092
>Bonjour
>Manal
>Alicia
>Marina
>
```

Consumer sur Kafka

```
* Support:        https://ubuntu.com/advantage

System information as of Sun Jan  9 17:04:13 UTC 2022

System load:  0.22          Processes:    106
Usage of /:   3.4% of 61.80GB Users logged in: 1
Memory usage: 61%         IP address for eth0: 10.0.2.15
Swap usage:   2%          IP address for eth1: 192.168.33.13

This system is built by the Bento project by Chef Software
More information can be found at https://github.com/chef/bento
Last login: Sun Jan  9 16:52:57 2022 from 10.0.2.2
Welcome to Ubuntu 18.04.5 LTS (GNU/Linux 4.15.0-151-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:        https://ubuntu.com/advantage

System information as of Sun Jan  9 17:04:13 UTC 2022

System load:  0.22          Processes:    106
Usage of /:   3.4% of 61.80GB Users logged in: 1
Memory usage: 61%         IP address for eth0: 10.0.2.15
Swap usage:   2%          IP address for eth1: 192.168.33.13

This system is built by the Bento project by Chef Software
More information can be found at https://github.com/chef/bento
Last login: Sun Jan  9 16:52:57 2022 from 10.0.2.2
vagrant@vagrant:~$ cd kafka
vagrant@vagrant:~/kafka$ bin/kafka-console-consumer.sh --topic quickstart-events --from-beginning --bootstrap-server @192.168.33.13:9092

Bonjour
Manal
Alicia
Marina
```

Spark

Spark est un Framework open source qui permet de traiter de grande quantité de données puis il récupère les informations dont on a besoin. La particularité de Spark est que les données s'exécutent en mémoire. Il a une autre fonction : c'est un consumer.

Spark streaming permet de décorréler un producer à un consumer. Le consumer a pour fonction de lire et de traiter ces événements écrits par Kafka.

Nous avons attribué à Spark l'adresse IP suivante : 192.168.33.12

Pipeline spark avec kafka : Streaming en temps réel

Au niveau de cette partie, nous avons relié kafka avec spark. Le consumer (kafka) et le producer (spark) avec la configuration de l'adresse ip 193.168.33.13.

```

root@kafka:vagrant ssh
Welcome to Ubuntu 18.04.5 LTS (GNU/Linux 4.15.0-151-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/advantage

System information as of Sun Jan 9 16:51:45 UTC 2022

System load: 0.5          Processes:           106
Usage of /:   3.4% of 61.0GB      Users logged in:     1
Memory usage: 70%            IP address for eth0: 10.0.2.15
Swap usage:   2%              IP address for eth1: 192.168.33.15

This system is built by the Bento project by Chef Software
More information can be found at https://github.com/chef/bento
Last login: Sun Jan 9 14:01:50 2022 from 10.0.2.2
vagrant@vagrant:~$ bin/kafka-console-producer.sh --topic quickstart-events --bootstrap-server 192.168.33.1:9092
^C
-bash: bin/kafka-console-producer.sh: No such file or directory
vagrant@vagrant:~$ cd kafka
vagrant@vagrant:~/kafka$ bin/kafka-console-producer.sh --topic quickstart-events --bootstrap-server 192.168.33.1:9092
^C
(Nothing sent)
vagrant@vagrant:~/kafka$ bin/kafka-console-producer.sh --topic quickstart-events --bootstrap-server 192.168.33.1:9092
^C
(Nothing, 'ITS 2')
vagrant@vagrant:~/kafka$ bin/kafka-console-producer.sh --topic quickstart-events --bootstrap-server 192.168.33.1:9092
^C
(Nothing, 'Big data')

```

```
vagrant@vagrant:~$ sudo mysql
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 2
Server version: 5.7.36-0ubuntu0.18.04.1 (Ubuntu)

Copyright (c) 2000, 2021, Oracle and/or its affiliates.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql>
mysql>
mysql>
mysql>
mysql> ls
->
-> SELECT user,authentication_string,plugin,host FROM mysql.user;
SELECT user,authentication_string,plugin,host FROM mysql.user;
^C
mysql> SELECT user,authentication_string,plugin,host FROM mysql.user;
+-----+-----+-----+-----+
| user                | authentication_string | plugin                | host                |
+-----+-----+-----+-----+
| root                | *THISISNOTAVALIDPASSWORDTHATCANBEUSEDHERE | auth_socket          | localhost           |
| mysql.session       | *THISISNOTAVALIDPASSWORDTHATCANBEUSEDHERE | mysql_native_password | localhost           |
| mysql.sys           | *THISISNOTAVALIDPASSWORDTHATCANBEUSEDHERE | mysql_native_password | localhost           |
| debian-sys-maint    | *3F7CD2E03C9E68782E5D59A643D5FDB6DC72BE4F | mysql_native_password | localhost           |
+-----+-----+-----+-----+
4 rows in set (0.00 sec)
```

Superset

Superset est un logiciel open source de datavisualisation et d'exploration des données. Ce logiciel permet de traiter des données massives.

Nous avons attribué à Spark l'adresse IP suivante : 192.168.33.10

```
PS C:\superset> vagrant ssh
Welcome to Ubuntu 18.04.5 LTS (GNU/Linux 4.15.0-151-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/advantage

System information as of Sun Jan  9 17:13:38 UTC 2022

System load:  0.0               Processes:    100
Usage of /:   4.9% of 61.80GB   Users logged in:  0
Memory usage: 16%              IP address for eth0: 10.0.2.15
Swap usage:   0%               IP address for eth1: 192.168.33.10

This system is built by the Bento project by Chef Software
More information can be found at https://github.com/chef/bento
Last login: Sun Jan  9 16:42:54 2022 from 10.0.2.2
vagrant@vagrant:~$
```

Problèmes rencontrés

Nous avons rencontré un problème lors de l'installation de Kafka. Cependant, grâce au dernier cours nous avons réussi à le résoudre. Par la suite, nous avons réalisé l'installation de MySQL et superset. L'installation de superset n'est pas encore tout à fait finie car nous rencontrons un problème, mais nous pensons que cette erreur est dû au fait que nous n'avons pas encore téléchargé notre base de données.

Afin de finir ce projet, il nous reste quelques étapes à réaliser : le traitement de notre base de données sur MySQL, modification du fichier python qui permet d'envoyer les données sur Kafka et puis la partie visualisation sur superset.