



FACULTÉ DES SCIENCES DE SFAX



## RAPPORT DE PROJET BIG DATA

---

# Pipeline Big Data pour l'Ingestion, le Stockage, et la Visualisation Temps Réel des Données de Transactions Bancaires

---

ELABORÉ PAR :

Ahmed Rami Belghuith

Dhia Elhak Toukebri

Année universitaire : 2024/2025

## BIG DATA



# Table des matières

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Objectif . . . . .	2
1.2	Contexte . . . . .	2
1.3	Analyse des Besoins Fonctionnels . . . . .	2
1.4	Bénéfices attendus . . . . .	2
<b>2</b>	<b>Configuration du Pipeline Big Data</b>	<b>4</b>
2.1	Introduction . . . . .	5
2.2	Configuration du Hadoop . . . . .	5
2.2.1	Configuration de core-site.xml . . . . .	5
2.2.2	Configuration de hdfs-site.xml . . . . .	6
2.2.3	Configuration de yarn-site.xml . . . . .	7
2.2.4	Configuration de zoo.cfg (Zookeeper) . . . . .	9
2.3	Ajout de Hue . . . . .	9
2.3.1	configuration de Hue.ini . . . . .	9
2.4	Configuration d'Apache Kafka . . . . .	11
2.4.1	Configuration de server.properties . . . . .	11
2.5	Configuration d'Apache Hive . . . . .	12
2.5.1	Configuration de hive-site.xml . . . . .	12
2.5.2	L'ajout de connector dans hive-env.sh . . . . .	13
2.6	Configuration de Apache Airflow . . . . .	13
2.6.1	Configuration de l'environnement Airflow . . . . .	13
2.6.2	Implémentation des Composants . . . . .	14
2.6.3	Intégration avec HDFS . . . . .	15
2.6.4	Pipeline Airflow . . . . .	16
2.7	Configuration de Apache superset . . . . .	17
2.7.1	configuration de supersetconfig . . . . .	17
2.8	Configuration de Fichier <b>.bashrc</b> . . . . .	17
2.8.1	Configuration ajouté à <b>.bashrc</b> . . . . .	17

---

<b>3</b>	<b>Réalisation de projet</b>	<b>19</b>
3.1	Architecture . . . . .	20
3.2	Activation des technologies . . . . .	20
3.2.1	Activation de Hadoop et ces services . . . . .	20
3.2.2	Activation de Hive et ces services . . . . .	22
3.2.3	Activation de Superset . . . . .	22
3.2.4	Activation de Hue . . . . .	23
3.2.5	Activation de Kafka . . . . .	24
3.2.6	Activation de Airflow . . . . .	25
3.3	Réalisation finale . . . . .	26
3.3.1	Analyse sur dataset de transactions . . . . .	26
3.3.2	Création de tableau "fraude" avec Hue . . . . .	27
3.3.3	Division et l'envoi vers pipeline d'ingestion . . . . .	28
3.3.4	Execution de pipeline d'ingestion avec airflow . . . . .	30
3.3.5	Mise à jour de tableau et des visualisations en temps réelle . . . . .	31
3.3.6	Dashboard Final avec superset . . . . .	32

# Table des figures

3.1	Architecture BIGDATA . . . . .	20
3.2	activation de HDFS/Yarn/zoo . . . . .	21
3.3	Activation de hive . . . . .	22
3.4	Activation de superset . . . . .	23
3.5	Activation de Hue . . . . .	24
3.6	Activation de Kafka . . . . .	24
3.7	listes des Dags (airflow) . . . . .	25
3.8	Activer scheduler de Airflow . . . . .	26
3.9	Activation Airflow Webserver . . . . .	26
3.10	création de tableau dans hive . . . . .	27
3.11	Division de fichier . . . . .	28
3.12	Envoi les fichiers . . . . .	29
3.13	Execution de Dag . . . . .	30
3.14	Pipeline . . . . .	30
3.15	Visualisation avant l’ingestion totale de données . . . . .	31
3.16	Visualisation après l’ingestion totale de données . . . . .	31
3.17	Dashboard Final avec superset . . . . .	32

## 1.1 Objectif

Ce Projet a pour but de fournir une description détaillée des étapes de configuration et de mise en œuvre d'un pipeline Big Data. Le pipeline vise à automatiser l'ingestion, le stockage, l'analyse, et la visualisation des données en temps réel en utilisant des outils tels qu'Apache Airflow, Kafka, HDFS, Hive, et Superset. Il sert de guide pratique pour déployer la solution et garantir une compréhension claire des processus sous-jacents.

## 1.2 Contexte

Dans un contexte où les entreprises manipulent d'énormes volumes de données, il est crucial de mettre en place des systèmes automatisés pour gérer ces données efficacement. Le projet répond au besoin d'automatiser l'ingestion de fichiers CSV, leur stockage dans un système distribué (HDFS), leur exploitation via une base de données (Hive), et leur visualisation en temps réel à travers un tableau de bord interactif (Superset). L'objectif est de permettre une mise à jour fluide et continue des tableaux de bord sans intervention manuelle, assurant ainsi une prise de décision basée sur des données toujours à jour.

## 1.3 Analyse des Besoins Fonctionnels

Le pipeline couvre les étapes suivantes :

0. **Ingestion de Données** : Surveillance d'un répertoire pour détecter de nouveaux fichiers CSV à traiter.
0. **Transport via Kafka** : Transmission des fichiers détectés vers un cluster HDFS en utilisant Apache Kafka.
0. **Stockage dans HDFS** : Enregistrement des fichiers CSV dans un système de fichiers distribué (HDFS).
0. **Analyse avec Hive** : Intégration des données stockées dans HDFS dans une base de données Hive pour permettre des requêtes SQL.
0. **Visualisation avec Superset** : Création d'un tableau de bord interactif et mise à jour automatique des données visualisées grâce à une connexion JDBC entre Hive et Superset.

## 1.4 Bénéfices attendus

- **Automatisation complète** : Réduction de l'intervention humaine dans le traitement des fichiers.
- **Efficacité** : Traitement rapide et fiable de grandes quantités de données.

- **Visualisation en temps réel** : Mise à jour immédiate des tableaux de bord dès l'ingestion d'un nouveau fichier.
- **Scalabilité** : Extensibilité pour inclure de nouvelles sources de données ou de nouveaux cas d'utilisation.

## 2.1 Introduction

Cette section décrit les configurations nécessaires pour déployer et exécuter le pipeline sur un cluster Hadoop comprenant un nœud maître (master) et deux nœuds esclaves (slave1 et slave2). Elle détaille les configurations des outils utilisés (Airflow, Kafka, HDFS, Hue, Hive, Superset) ainsi que la mise en place des connexions entre les composants.

## 2.2 Configuration du Hadoop

### Nœud maître (master)

- Responsable : Gestion des ressources, coordination des tâches et supervision globale du cluster.
- **Services principaux :**
  - *Namenode* : Gère les métadonnées du système de fichiers HDFS.
  - *ResourceManager* : Alloue les ressources aux applications distribuées.
  - *Zookeeper* : Service de coordination distribué pour la gestion des configurations et le suivi des états.
  - Hue : Interface utilisateur pour interagir avec Hadoop et exécuter diverses tâches (soumission de jobs, gestion des fichiers, etc.).

### Nœuds esclaves (slave1, slave2)

- Responsable : Exécution des tâches distribuées.
- **Services :**
  - *Datanode* : Stocke les données de manière distribuée dans HDFS.
  - *NodeManager* : Gère les conteneurs pour l'exécution des tâches MapReduce ou Spark.
  - *Zookeeper client* : Assure la communication avec le service Zookeeper pour garantir une coordination fiable et une haute disponibilité.

### 2.2.1 Configuration de core-site.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4     <property>
5         <name>fs.defaultFS</name>
```

```
6         <value>hdfs://master:9000</value>
7     </property>
8     <property>
9         <name>dfs.webhdfs.enabled</name>
10        <value>true</value>
11    </property>
12    <property>
13        <name>dfs.webhdfs.address</name>
14        <value>master:50070</value>
15    </property>
16    <property>
17        <name>hadoop.proxyuser.hue.groups</name>
18        <value>*</value>
19    </property>
20    <property>
21        <name>hadoop.proxyuser.hue.hosts</name>
22        <value>*</value>
23    </property>
24 </configuration>
```

Listing 2.1 – core-site.xml

### 2.2.2 Configuration de hdfs-site.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4     <property>
5         <name>dfs.replication</name>
6         <value>2</value>
7     </property>
8     <property>
9         <name>dfs.namenode.name.dir</name>
10        <value>file:/home/master/hdata/hdfs/namenode</value>
11    </property>
12    <property>
```



```
13 <name>dfs.datanode.data.dir</name>
14 <value>file:/home/master/hdata/hdfs/datanode</value>
15 </property>
16 <property>
17 <name>dfs.namenode.checkpoint.dir</name>
18 <value>file:/home/master/hdata/hdfs/namesecondary</value>
19 </property>
20 <property>
21   <name>dfs.webhdfs.enabled</name>
22   <value>true</value>
23 </property>
24 <property>
25   <name>dfs.namenode.http-address</name>
26   <value>master:50070</value>
27 </property>
28 <property>
29   <name>dfs.permissions.enabled</name>
30   <value>false</value>
31 </property>
32 <property>
33   <name>dfs.namenode.acls.enabled</name>
34   <value>false</value>
35 </property>
36 <property>
37   <name>dfs.permissions</name>
38   <value>false</value>
39 </property>
40
41 </configuration>
```

Listing 2.2 – hdfs-site.xml

### 2.2.3 Configuration de yarn-site.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
```

```
3 <configuration>
4 <property>
5 <name>yarn.resourcemanager.hostname</name>
6 <value>master</value>
7 </property>
8 <property>
9 <name>yarn.nodemanager.local-dirs</name>
10 <value>file:/home/master/hdata/yarn/local</value>
11 </property>
12 <property>
13 <name>yarn.nodemanager.log-dirs</name>
14 <value>file:/home/master/hdata/yarn/log</value>
15 </property>
16 <property>
17 <name>yarn.log.server.url</name>
18 </property>
19 <property>
20 <name>yarn.nodemanager.aux-services</name>
21 <value>mapreduce_shuffle,spark_shuffle</value>
22 </property>
23 <property>
24 <name>yarn.nodemanager.aux-services.mapreduce_shuffle.class</name>
25 <value>org.apache.hadoop.mapred.ShuffleHandler,
26 org.apache.spark.network.yarn.YarnShuffleService</value>
27 </property>
28 <property>
29 <name>yarn.nodemanager.disk-health-checker
30 .max-disk-utilization-per-disk-percentage</name>
31 <value>100</value>
32 </property>
33 <property>
34 <name>mapred.jobtracker.taskScheduler</name>
35 <value>org.apache.hadoop.mapred.CapacityTaskScheduler</value>
36 </property>
37 <property>
```

```
38 <name>yarn.resourcemanager.proxyuser.hue.groups</name>
39 <value>*</value>
40 </property>
41 <property>
42 <name>yarn.resourcemanager.proxyuser.hue.hosts</name>
43 <value>*</value>
44 </property>
45 </configuration>
```

Listing 2.3 – yarn-site.xml

### 2.2.4 Configuration de zoo.cfg (Zookeeper)

```
1 tickTime=2000
2 initLimit=10
3 syncLimit=5
4 dataDir=/var/lib/zookeeper
5 clientPort=2181
6 server.1=master:2888:3888
7 server.2=slave1:2888:3888
8 server.3=slave2:2888:3888
```

Listing 2.4 – zoo.cfg

## 2.3 Ajout de Hue

### 2.3.1 configuration de Hue.ini

```
1 [desktop]
2 default_user = master
3 default_superuser = true
4 log_dir=/var/log/hue
5 logging_level=DEBUG
6 secret_key=
7 http_host=0.0.0.0
8 http_port=8888
9 time_zone=America/Los_Angeles
```

```
10 django_debug_mode=false
11 http_500_debug_mode=false
12 [[custom]]
13 [[auth]]
14 idle_session_timeout=-1
15 host=localhost
16 port=25
17 user=
18 password=
19 tls=no
20 [hadoop]
21 proxy_user_enable=true
22 [[hdfs_clusters]]
23 [[[default]]]
24 hadoop_username=master
25 fs_defaultfs=hdfs://192.168.1.2:9000
26 hadoop_conf_dir=/home/master/hadoop-3.3.6/etc/hadoop
27 security_enabled=false
28 webhdfs_url=http://192.168.1.2:50070/webhdfs/v1/
29 submit_to=True
30 enable_doas=True
31 [[yarn_clusters]]
32 [[[default]]]
33 resourcemanager_host=192.168.1.2
34 resourcemanager_port=8088
35 submit_to=True
36 resourcemanager_api_url=http://192.168.1.2:8088
37 history_server_api_url=http://192.168.1.2:19888
38 [beeswax]
39 hive_server_host=localhost
40 hive_server_port=10000
41 hive_server_username=hive
42 hive_server_password=password
43 thrift_transport=TBinaryProtocol
44 hive_server2_authentication=NOSASL
```

```
45  
46 [filebrowser]  
47 show_upload_button=true  
48 level = DEBUG  
49 file = /var/log/hue/hue.log  
50 rotation = true  
51 max_size = 10485760  
52 backup_count = 5  
53 console = true
```

Listing 2.5 – Hue.ini

## 2.4 Configuration d'Apache Kafka

### 2.4.1 Configuration de server.properties

Sur chaque nœud :

```
1 broker.id=1 # =2 pour slave1 et =3 pour slave2  
2 listeners=PLAINTEXT://:9092  
3 advertised.listeners=PLAINTEXT://master:9092  
4 num.network.threads=3  
5 num.io.threads=8  
6 socket.send.buffer.bytes=102400  
7 socket.receive.buffer.bytes=102400  
8 socket.request.max.bytes=104857600  
9 log.dirs=/var/lib/kafka  
10 num.partitions=1  
11 offsets.topic.replication.factor=1  
12 transaction.state.log.replication.factor=1  
13 transaction.state.log.min.isr=1  
14 log.retention.hours=168  
15 log.retention.check.interval.ms=300000  
16 zookeeper.connect=master:2181,slave1:2181,slave2:2181  
17 zookeeper.connection.timeout.ms=18000  
18 group.initial.rebalance.delay.ms=0
```

Listing 2.6 – server.properties

## 2.5 Configuration d'Apache Hive

### 2.5.1 Configuration de hive-site.xml

```
1 <property>
2     <name>javax.jdo.option.ConnectionURL</name>
3     <value>jdbc:mysql://localhost/metastore_dbcreateDatabaseIfNotExist=
4     true&amp;allowPublicKeyRetrieval=true</value>
5 </property>
6 <property>
7     <name>javax.jdo.option.ConnectionDriverName</name>
8     <value>com.mysql.cj.jdbc.Driver</value>
9 </property>
10 <property>
11     <name>javax.jdo.option.ConnectionUserName</name>
12     <value>hiveuser</value>
13 </property>
14 <property>
15     <name>javax.jdo.option.ConnectionPassword</name>
16     <value>hivepassword</value>
17 </property>
18 <property>
19     <name>hive.server2.thrift.bind.host</name>
20     <value>localhost</value>
21 </property>
22 <property>
23     <name>hive.server2.thrift.port</name>
24     <value>10000</value>
25 </property>
26 <property>
27     <name>hive.server2.thrift.port</name>
28     <value>10000</value>
29 </property>
```

```
30 <!---autre ----->
```

Listing 2.7 – Hive-site.xml

## 2.5.2 L'ajout de connector dans hive-env.sh

```
1 # export HIVE_AUX_JARS_PATH=
2 export HIVE_CLASSPATH=/opt/hive/lib/mysql-connector-j-8.0.32.jar:$HIVE_CLASSPATH
```

Listing 2.8 – hive-env.sh

## 2.6 Configuration de Apache Airflow

### 2.6.1 Configuration de l'environnement Airflow

```
1 # Configuration de l'environnement
2 AIRFLOW_HOME = '/opt/airflow'
3 DAGS_FOLDER = '/opt/airflow/dags'
4
5 # Configuration des connexions
6 KAFKA_BROKERS = 'master:9092,slave1:9092,slave2:9092'
7 KAFKA_TOPIC = 'csv-topic'
8 HDFS_URL = 'http://master:50070'
9 HDFS_USER = 'hdfs'
10 HDFS_BASE_PATH = '/kafka_ingestion/csv_files'
11
12 MONITORED_FOLDER = os.path.join(os.path.expanduser('~'), 'kafka_ingestion')
13 PROCESSED_FOLDER = os.path.join(MONITORED_FOLDER, 'processed')
14
15 # Configuration du DAG
16 default_args = {
17     'owner': 'airflow',
18     'depends_on_past': False,
19     'email_on_failure': False,
20     'retries': 1,
21     'retry_delay': timedelta(minutes=1),
22 }
```

Listing 2.9 – Configuration de l’environnement Airflow

## 2.6.2 Implémentation des Composants

### 2.6.2.1 Producteur Kafka

```
1 def create_kafka_producer():
2     return Producer({
3         'bootstrap.servers': KAFKA_BROKERS,
4         'client.id': 'airflow_producer',
5         'acks': 'all',
6         'retries': 3,
7         'linger.ms': 1,
8         'batch.size': 16384,
9         'buffer.memory': 33554432
10    })
11
12 def produce_to_kafka(**kwargs):
13     csv_files = kwargs['ti'].xcom_pull(key='csv_files', task_ids='scan_csv_files')
14     producer = create_kafka_producer()
15     try:
16         for csv_file in csv_files:
17             full_path = os.path.join(MONITORED_FOLDER, csv_file)
18             producer.produce(KAFKA_TOPIC, full_path.encode('utf-8'))
19             producer.flush()
20     finally:
21         producer.close()
```

Listing 2.10 – Implémentation du Producteur Kafka

### 2.6.2.2 Consommateur Kafka

```
1 def create_kafka_consumer():
2     return Consumer({
3         'bootstrap.servers': KAFKA_BROKERS,
4         'group.id': 'csv-hdfs-ingestion-airflow',
```



```
5         'auto.offset.reset': 'earliest',
6         'enable.auto.commit': False,
7         'max.poll.interval.ms': 300000
8     })
9
10 def consume_and_ingest_to_hdfs(**kwargs):
11     consumer = create_kafka_consumer()
12     consumer.subscribe([KAFKA_TOPIC])
13     try:
14         while True:
15             msg = consumer.poll(1.0)
16             if msg is None:
17                 break
18             if msg.error():
19                 continue
20             csv_path = msg.value().decode('utf-8').strip()
21             filename = os.path.basename(csv_path)
22             hdfs_destination = os.path.join(HDFS_BASE_PATH, filename)
23
24             with open(csv_path, 'rb') as local_file:
25                 with hdfs_client.write(hdfs_destination, overwrite=True) as hdfs_file:
26                     shutil.copyfileobj(local_file, hdfs_file)
27     finally:
28         consumer.close()
```

Listing 2.11 – Implémentation du Consommateur Kafka

### 2.6.3 Intégration avec HDFS

```
1 def write_to_hdfs(data, hdfs_path):
2     client = InsecureClient(HDFS_URL, user=HDFS_USER)
3     with client.write(hdfs_path, overwrite=True) as writer:
4         writer.write(data)
```

Listing 2.12 – Intégration avec HDFS

### 2.6.4 Pipeline Airflow

Le pipeline Airflow défini ci-dessous orchestre les différentes étapes du processus, de la détection des fichiers CSV à leur ingestion dans HDFS.

```
1 with DAG(
2     'kafka_hdfs_ingestion_pipeline',
3     default_args=default_args,
4     description='Pipeline to monitor folder, send to Kafka, and ingest to HDFS',
5     schedule_interval='*/1* * * * *',
6     start_date=datetime(2024, 1, 1),
7     catchup=False,
8 ) as dag:
9
10    scan_task = PythonOperator(
11        task_id='scan_csv_files',
12        python_callable=scan_csv_files,
13        provide_context=True,
14    )
15
16    produce_task = PythonOperator(
17        task_id='produce_to_kafka',
18        python_callable=produce_to_kafka,
19        provide_context=True,
20    )
21
22    ingest_task = PythonOperator(
23        task_id='consume_and_ingest_to_hdfs',
24        python_callable=consume_and_ingest_to_hdfs,
25        provide_context=True,
26    )
27
28    scan_task >> produce_task >> ingest_task
```

Listing 2.13 – Pipeline Airflow

## 2.7 Configuration de Apache superset

### 2.7.1 configuration de supersetconfig

```
1 # Flask application settings
2 FLASK_APP = 'superset.app:create_app()'
3 APP_DIR = '/home/master/spark_env/lib/python3.11/site-packages/superset'
4
5 # Superset specific config
6 ROW_LIMIT = 10000
7 SQLLAB_ASYNC_TIME_LIMIT_SEC = 3600
8
9 SQLLAB_TIMEOUT = 3600
10 SQLLAB_HARD_TIMEOUT = 3660
11 SUPERSET_WEBSERVER_TIMEOUT = 3600 # Timeout serveur web
12 QUERY_TIMEOUT = 3600
13 SQLALCHEMY_QUERY_TIMEOUT = 3600 # Timeout SQLAlchemy
14 CACHE_DEFAULT_TIMEOUT = 60 * 60 * 24 # 24 heures
15 SQLLAB_DOCKER_TIMEOUT = 3600
16 # Flask App Builder configuration
17 SECRET_KEY = 'TdBVVxe3H5dmoJHL2ftGFY1AhH6Nq/7F2oruFOWfFYQ9/UfSGKWLtrT'
18 # Flask-WTF flag for CSRF
19 WTF_CSRF_ENABLED = True
20 WTF_CSRF_EXEMPT_LIST = []
21 WTF_CSRF_TIME_LIMIT = 60 * 60 * 24 * 365 # 1 an
22 # Client-side override
23 OVERRIDE_TIMEOUT = True
24 # Debugging and extra options
25 FLASK_APP_EXTRA_OPTIONS = {'debug': True}
```

Listing 2.14 – superset\_config.py

## 2.8 Configuration de Fichier .bashrc

### 2.8.1 Configuration ajouté à .bashrc

```
1 export JAVA_HOME=/home/master/jdk1.8.0_92
2 export PATH=$JAVA_HOME/bin:$PATH
3 export HADOOP_HOME=$HOME/hadoop-3.3.6
4 export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
5 export HADOOP_MAPRED_HOME=$HADOOP_HOME
6 export HADOOP_COMMON_HOME=$HADOOP_HOME
7 export HADOOP_HDFS_HOME=$HADOOP_HOME
8 export YARN_HOME=$HADOOP_HOME
9 export PATH=$PATH:$HADOOP_HOME/bin
10 export PATH=$PATH:$HADOOP_HOME/sbin
11 export HIVE_HOME=/opt/hive
12 export HIVE_CONF_DIR=$HIVE_HOME/conf
13 export HADOOP_CLASSPATH=$HADOOP_HOME/lib/*:$HADOOP_CONF_DIR
14 export PATH=$PATH:$HIVE_HOME/bin
15 export CLASSPATH=$CLASSPATH:/opt/hive/lib/mysql-connector-j-8.0.32.jar
16 export SPARK_HOME=/opt/spark # Adaptez le chemin selon votre install>
17 export PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin
18 export PYSPARK_PYTHON=/home/master/spark_env/bin/python3
19 export PYSPARK_DRIVER_PYTHON=/home/master/spark_env/bin/python3
20 export PYSPARK_DRIVER_PYTHON_OPTS='notebook'
21 export SUPERSET_CONFIG_PATH="/home/master/superset/superset_config.py"
22 export FLASK_APP=superset
23
24 export PATH="$HOME/.pyenv/bin:$PATH"
25 eval "$(pyenv_init-)"
26 eval "$(pyenv_virtualenv-init-)"
27
28 export NVM_DIR="$HOME/.config/nvm"
29 [ -s "$NVM_DIR/nvm.sh" ] && \. "$NVM_DIR/nvm.sh" # This loads nvm
30 [ -s "$NVM_DIR/bash_completion" ] && \. "$NVM_DIR/bash_completion"
```

Listing 2.15 – .bashrc

## 3.1 Architecture

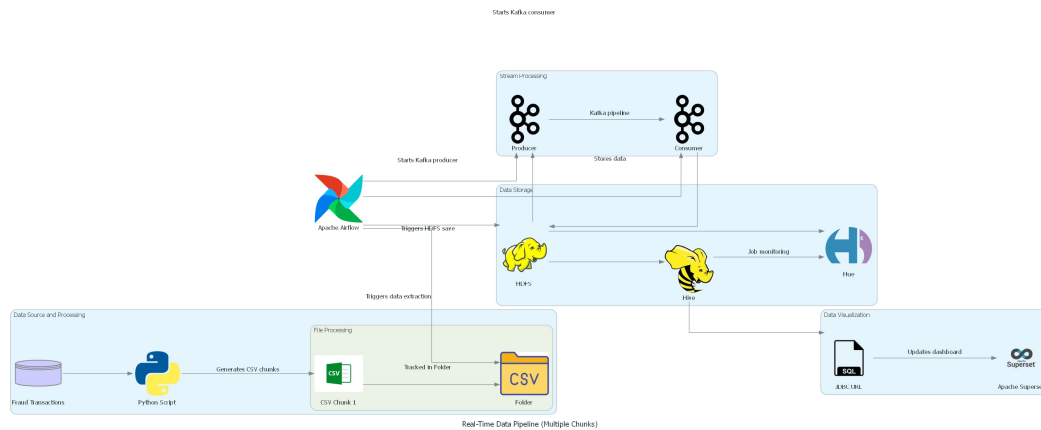


FIGURE 3.1 : Architecture BIGDATA

## 3.2 Activation des technologies

### 3.2.1 Activation de Hadoop et ces services

#### 3.2.1.1 HDFS/YARN/Zookeeper

#### Hadoop Cluster Startup

Pour démarrer les clusters HDFS et YARN :

#### Démarrer les démons HDFS

- Démarrer le daemon **Namenode** sur la machine virtuelle du Namenode :

```
1  hdfs --daemon start namenode
```

- Démarrer les daemons **DataNode** sur les nœuds Datanode :

```
1  hdfs --daemon start datanode
```

#### Démarrer les démons YARN

- Démarrer le daemon **resourcemanager** sur le nœud Namenode :

```
1 yarn --daemon start resourcemanager
```

- Démarrer les daemons **nodemanager** sur les nœuds Datanode :

```
1 yarn --daemon start nodemanager
```

- Démarrer le serveur **JobHistory MapReduce** :

```
1 mapred --daemon start historyserver
```

### Démarrer Zookeeper

#### Configuration

- Modifier les permissions des fichiers Zookeeper :

```
1 sudo chown -R master:master /opt/zookeeper/logs
2 sudo chown -R master:master /var/lib/zookeeper
```

#### Démarrer le serveur Zookeeper

- Démarrer le serveur :

```
1 /opt/zookeeper/bin/zkServer.sh start
```

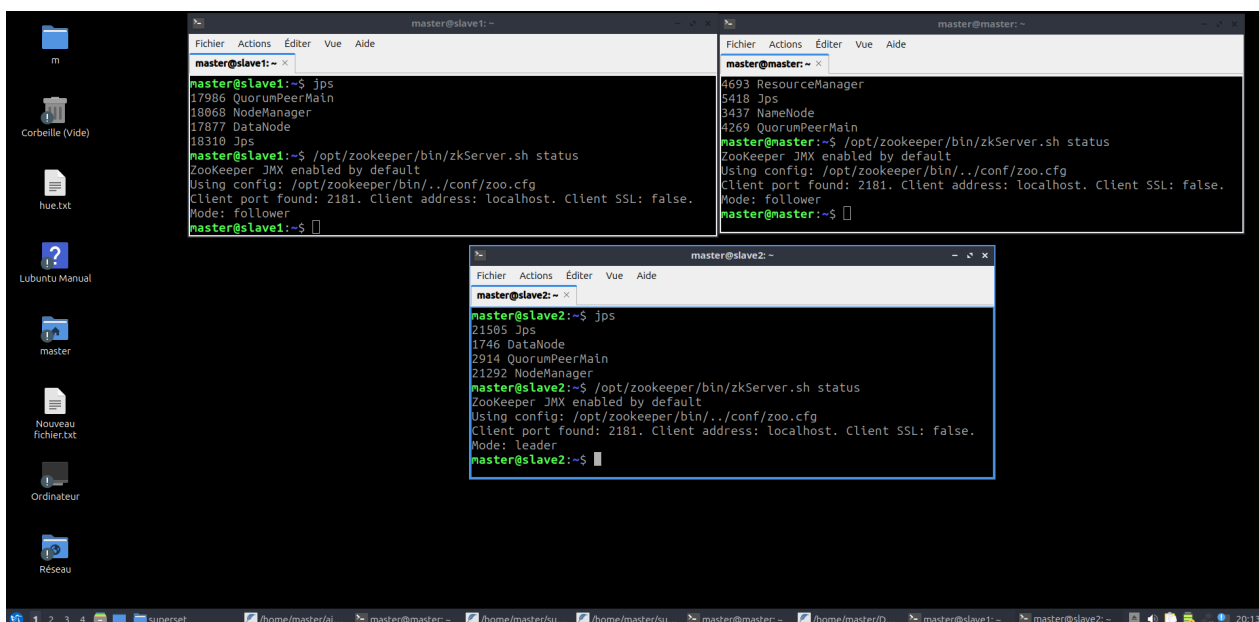


FIGURE 3.2 : activation de HDFS/Yarn/zoo

### 3.2.2 Activation de Hive et ces services

- Démarrer les services Hive :

```
1  hive --service metastore &  
2  hive --service hiveserver2 &
```

- Connexion à Hive via Beeline :

```
1  beeline -u jdbc:hive2://localhost:10000/default -n hive -p password
```

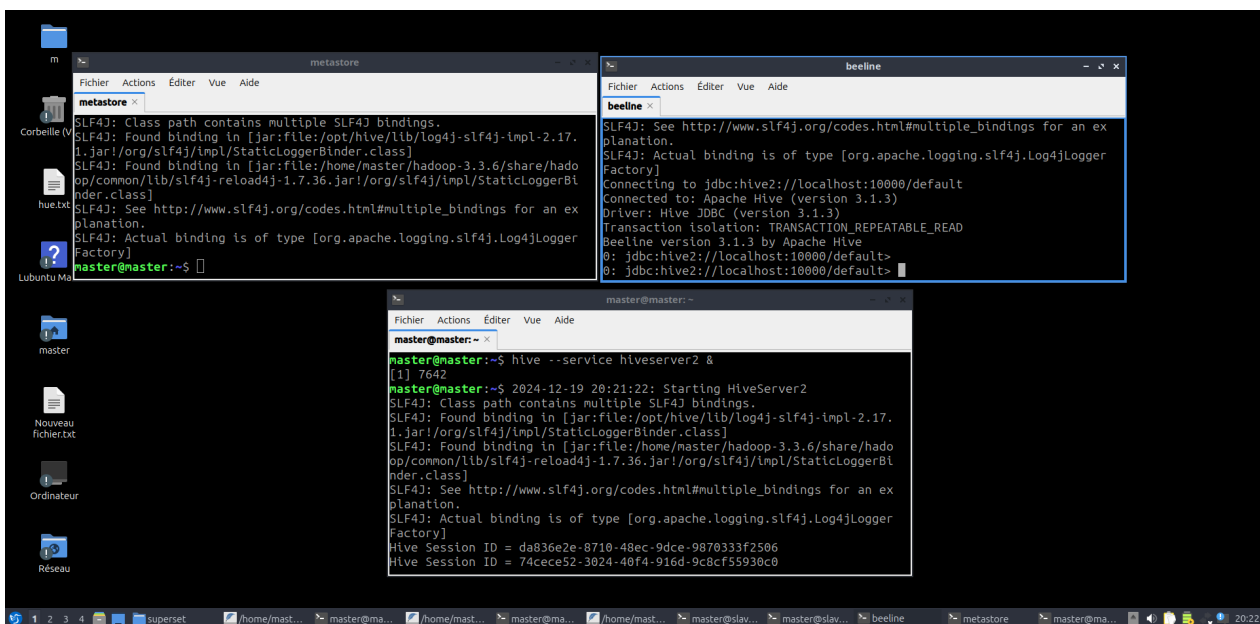


FIGURE 3.3 : Activation de hive

### 3.2.3 Activation de Superset

#### Configuration initiale

- Configurer Superset :

```
1  mkdir -p ~/superset  
2  echo "SECRET_KEY='secret_key'" > ~/superset/superset_config.py  
3  superset db upgrade  
4  superset init
```

#### Démarrer Superset

- Activer l'environnement Superset :

```
1 source superset_env/bin/activate
```

— Démarrer Superset :

```
1 superset run -p 8088 --with-threads --reload --debugger
```

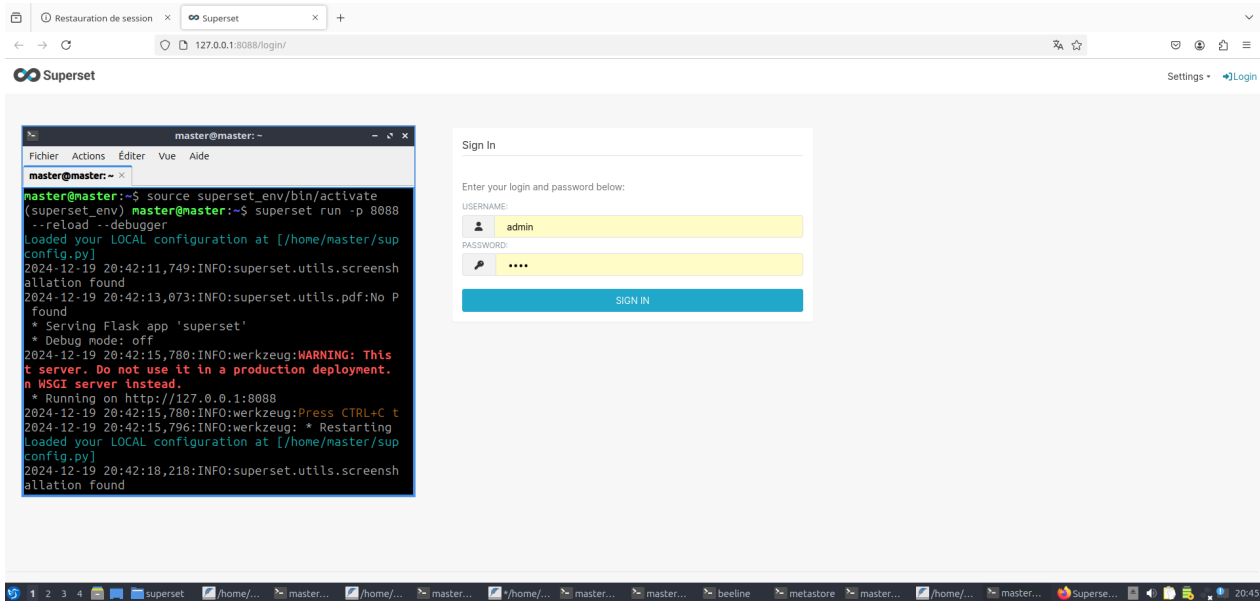


FIGURE 3.4 : Activation de superset

### 3.2.4 Activation de Hue

#### 3.2.4.1 Démarrage de Hue

— Activer l'environnement Hue :

```
1 source /home/master/hue/build/env/bin/activate
```

— Démarrer le serveur Hue :

```
1 sudo -u master /home/master/hue
2 /build/env/bin/python /home/master/hue/build/env/bin/hue runserver 0.0.0.0:8000
```



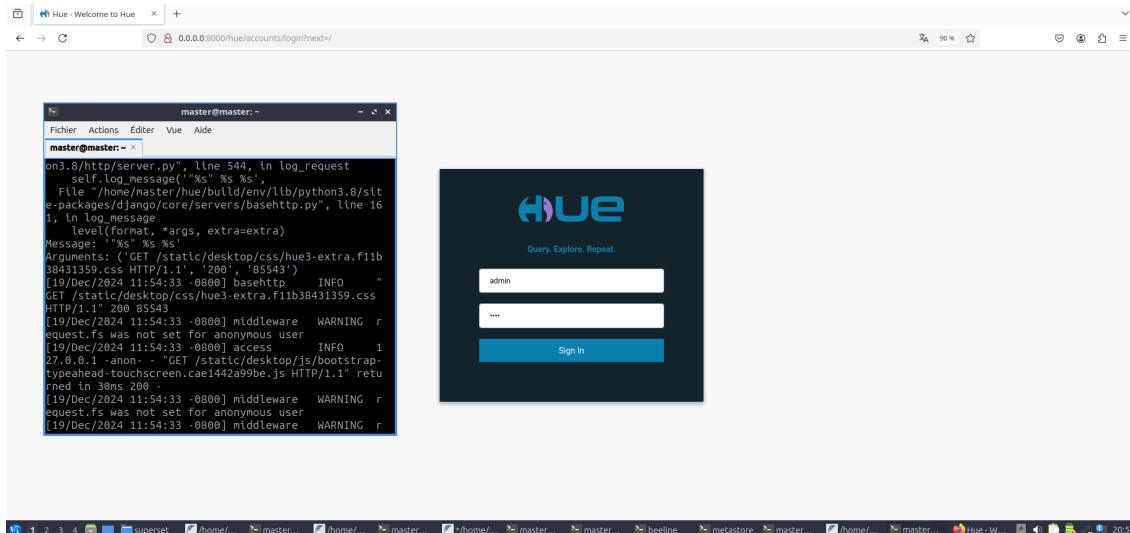


FIGURE 3.5 : Activation de Hue

### 3.2.5 Activation de Kafka

#### 3.2.5.1 Démarrage des services Kafka

— Démarrer le serveur Kafka :

```
1 bin/kafka-server-start.sh config/server.properties
2 (executer dans tout les noeuds)
```

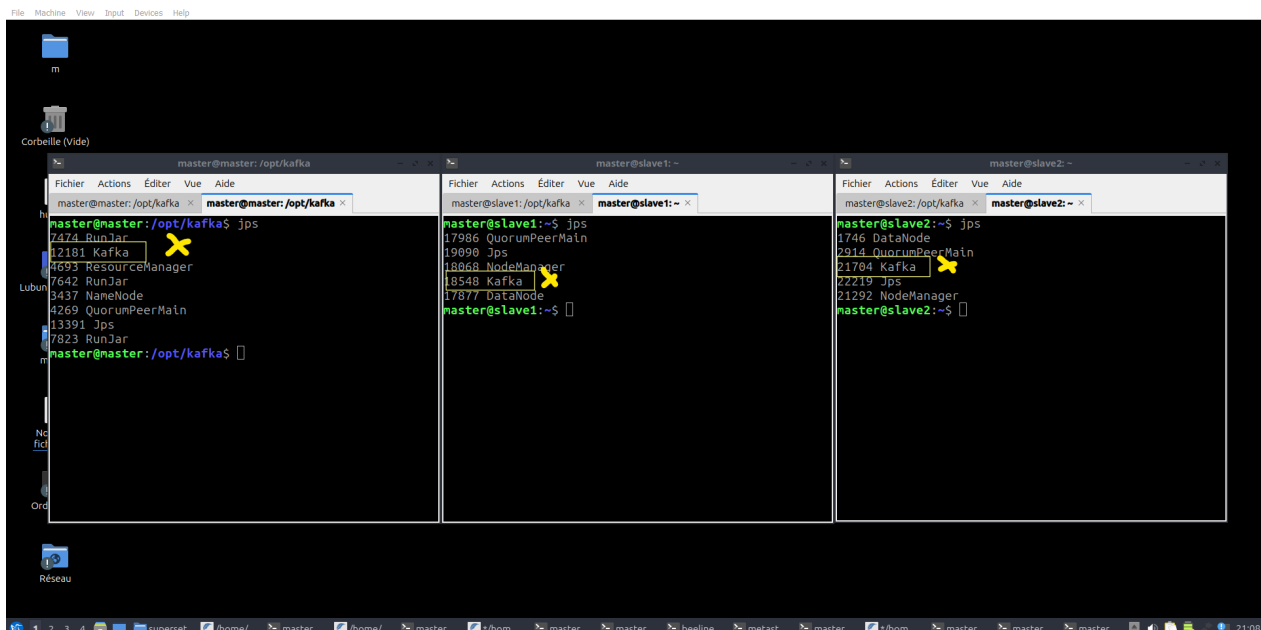


FIGURE 3.6 : Activation de Kafka

### 3.2.6 Activation de Airflow

- Afficher la liste des DAGs :

```
1 airflow dags list
```

- Démarrer le scheduler :

```
1 airflow scheduler
```

- Démarrer le serveur web Airflow :

```
1 airflow webserver --port 8081
```

```

master@master:~$ source /home/master/spark_env/bin/activate
(master_env) master@master:~$ airflow dags list
/home/master/spark_env/lib/python3.11/site-packages/airflow/cli/commands/dag_command.py:48 UserWarning: Could not import graphviz. Rendering graph to the graphical format will not be possible.
dag_id                                     fileloc                                     owners  is_paused
-----
conditional_dataset_and_time_based_timetable /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
consume_1_and_2_with_dataset_expressions    /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
consume_1_or_2_with_dataset_expressions     /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
consume_1_or_both_2_and_3_with_dataset_expressions /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
dataset_alias_example_alias_consumer        /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_alias_example_alias_consumer_with_no_taskflow /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_alias_example_alias_producer        /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_alias_example_alias_producer_with_no_taskflow /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_consumes_1                         /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
dataset_consumes_1_and_2                  /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
dataset_consumes_1_never_scheduled        /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
dataset_produces_1                        /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
dataset_produces_2                        /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_datasets.py airflow True
dataset_s3_bucket_consumer                /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_s3_bucket_consumer_with_no_taskflow /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_s3_bucket_producer                /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_s3_bucket_producer_with_no_taskflow /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_dataset_alias.py airflow True
dataset_with_extra_by_context              /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_outlet_event.py airflow True
dataset_with_extra_by_yield                /home/master/spark_env/lib/python3.11/site-packages/airflow/example_dags/example_outlet_event.py airflow True

```

FIGURE 3.7 : listes des Dags (airflow)

```
master@master:~$ source spark_env/bin/activate
(spark_env) master@master:~$ airflow scheduler

[2024-12-19 21:25:37.808+0100] [utils.py:161] INFO - NumExpr defaulting to 4 threads.
[2024-12-19 21:25:37.817+0100] [executor_loader.py:254] INFO - Loaded executor: SequentialExecutor
[2024-12-19 21:25:37.817+0100] [15279] [INFO] Starting gunicorn 21.2.0
[2024-12-19 21:25:37.817+0100] [15279] [INFO] Listening at: http://[::]:8793 (15279)
[2024-12-19 21:25:37.817+0100] [15279] [INFO] Using worker: sync
[2024-12-19 21:25:37.817+0100] [15280] [INFO] Booting worker with pid: 15280
[2024-12-19 21:25:37.817+0100] [scheduler_job_runner.py:950] INFO - Starting the scheduler
[2024-12-19 21:25:37.817+0100] [scheduler_job_runner.py:957] INFO - Processing each file at most -1 times
[2024-12-19 21:25:37.817+0100] [manager.py:174] INFO - Launched DagFileProcessorManager with pid: 15281
[2024-12-19 21:25:37.817+0100] [scheduler_job_runner.py:1949] INFO - Adopting or resetting orphaned tasks for active dag runs
[2024-12-19 21:25:37.817+0100] [settings.py:63] INFO - Configured default timezone UTC
[2024-12-19 21:25:37.817+0100] [manager.py:406] WARNING - Because we cannot use more than 1 thread (parsing_processes = 2) when using sqlite. So we set parallelism to 1.
[2024-12-19 21:25:37.817+0100] [15282] [INFO] Booting worker with pid: 15282
[2024-12-19 21:25:37.817+0100] [scheduler_job_runner.py:1972] INFO - Marked 1 SchedulerJob instances as failed
[2024-12-19 21:25:37.817+0100] [dag.py:4180] INFO - Setting next_dagrun for kafka_hdfs_ingestion_pipeline to 2024-12-19 20:25:00+00:00, run_after=2024-12-19 20:26:00+00:00
Dag run in running state
Dag information Queued at: 2024-12-19 20:25:38.576724+00:00 hash info: 3d195943ef207075d6569b7a62ffe3d7
[2024-12-19 21:25:38.918+0100] [scheduler_job_runner.py:435] INFO - 1 tasks up for execution:
<TaskInstance: kafka_hdfs_ingestion_pipeline.scan_csv_files scheduled_2024-12-19T20:24:00+00:00 [scheduled]>
[2024-12-19 21:25:38.918+0100] [scheduler_job_runner.py:507] INFO - DAG kafka_hdfs_ingestion_pipeline has 0/16 running and queued tasks
[2024-12-19 21:25:38.918+0100] [scheduler_job_runner.py:646] INFO - Setting the following tasks to queued state:
<TaskInstance: kafka_hdfs_ingestion_pipeline.scan_csv_files scheduled_2024-12-19T20:24:00+00:00 [scheduled]>
[2024-12-19 21:25:38.918+0100] [15282] [INFO] Trying to enqueue tasks: [<TaskInstance: kafka_hdfs_ingestion_pipeline.scan_csv_files scheduled_2024-12-19T20:24:00+00:00 [scheduled]>] for executor: SequentialExecutor(parallelism=32)
[2024-12-19 21:25:38.923+0100] [scheduler_job_runner.py:692] INFO - Sending TaskInstanceKey(dag_id='kafka_hdfs_ingestion_pipeline', task_id='scan_csv_files', run_id='scheduled_2024-12-19T20:24:00+00:00', try_number=1, map_index=-1) to SequentialExecutor with priority 3 and queue default
[2024-12-19 21:25:38.923+0100] [base_executor.py:169] INFO - Adding to queue: ['airflow', 'tasks', 'run', 'kafka_hdfs_ingestion_pipeline', 'scan_csv_files', 'scheduled_2024-12-19T20:24:00+00:00', '--local', '--subdir', 'DAGS_FOLDER/kafka_hdfs_ingestion_listen_dag.py']
[2024-12-19 21:25:38.954+0100] [sequential_executor.py:84] INFO - Executing command: ['airflow', 'tasks', 'run', 'kafka_hdfs_ingestion_pipeline', 'scan_csv_files', 'scheduled_2024-12-19T20:24:00+00:00', '--local', '--subdir', 'DAGS_FOLDER/kafka_hdfs_ingestion_listen_dag.py']
```

FIGURE 3.8 : Activer scheduler de Airflow

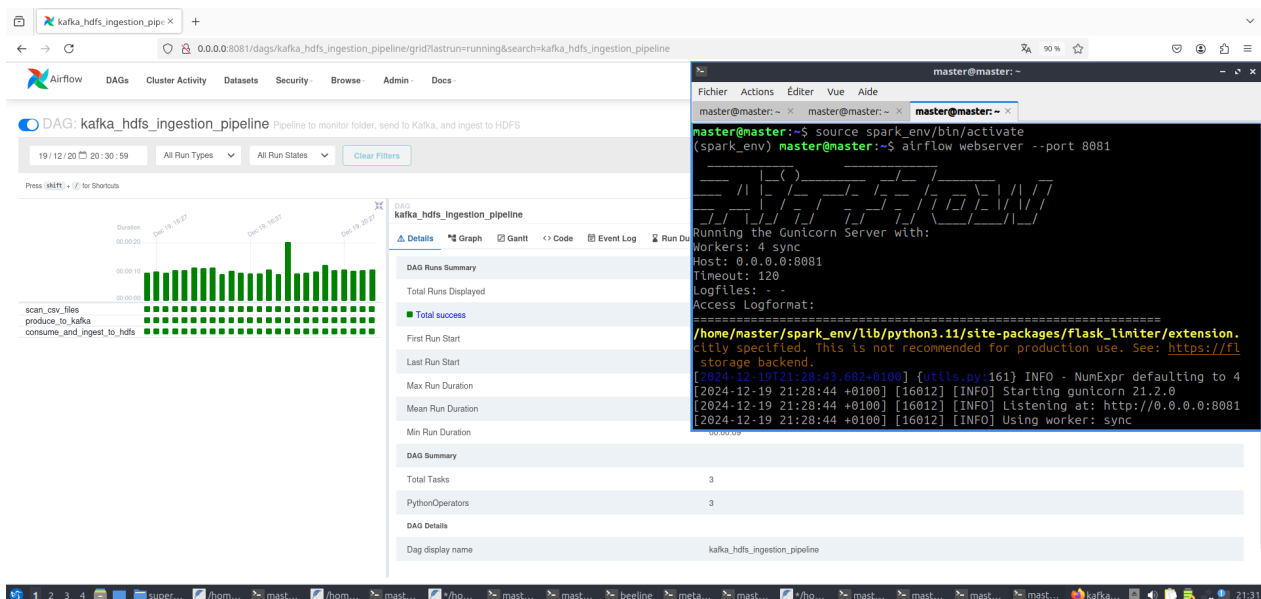


FIGURE 3.9 : Activation Airflow Webserver

## 3.3 Réalisation finale

### 3.3.1 Analyse sur dataset de transactions

La base de données transactions.csv contient un ensemble de **1 852 394** transactions, avec **23** colonnes variées décrivant chaque transaction. Cette vaste base de données semble concerner des transactions financières, incluant des informations telles que les **dates et heures des transactions**, les **numéros de carte de crédit** (cc\_num), les **commerçants impliqués**, les **catégories de dépenses**, ainsi que les **montants des transactions**

(amt). En outre, elle contient des données personnelles sur les utilisateurs, notamment leurs **prénoms, noms de famille, sexe, adresses, villes, états, codes postaux**, et **dates de naissance**.

Certaines informations **géographiques** telles que les **coordonnées GPS** (latitude et longitude) des commerçants et des utilisateurs sont également présentes, ce qui pourrait permettre de croiser des informations géographiques avec des habitudes de consommation. Les transactions sont accompagnées de données telles que le **numéro de transaction**, un identifiant **UNIX** pour la transaction (`unix_time`), et des informations sur l'occupation des utilisateurs, incluant leur **métier**. Enfin, la colonne `is_fraud` est utilisée pour indiquer si une transaction a été classifiée comme **frauduleuse**, avec une valeur de **1** pour frauduleuse et **0** pour non frauduleuse.

La base de données semble être principalement utilisée pour l'analyse de **fraudes** dans les transactions financières, avec un fort potentiel d'application dans les modèles d'**apprentissage automatique** et de détection des fraudes, notamment en croisant des informations **transactionnelles, géographiques et personnelles** pour prédire les comportements suspects.

### 3.3.2 Création de tableau "fraude" avec Hue

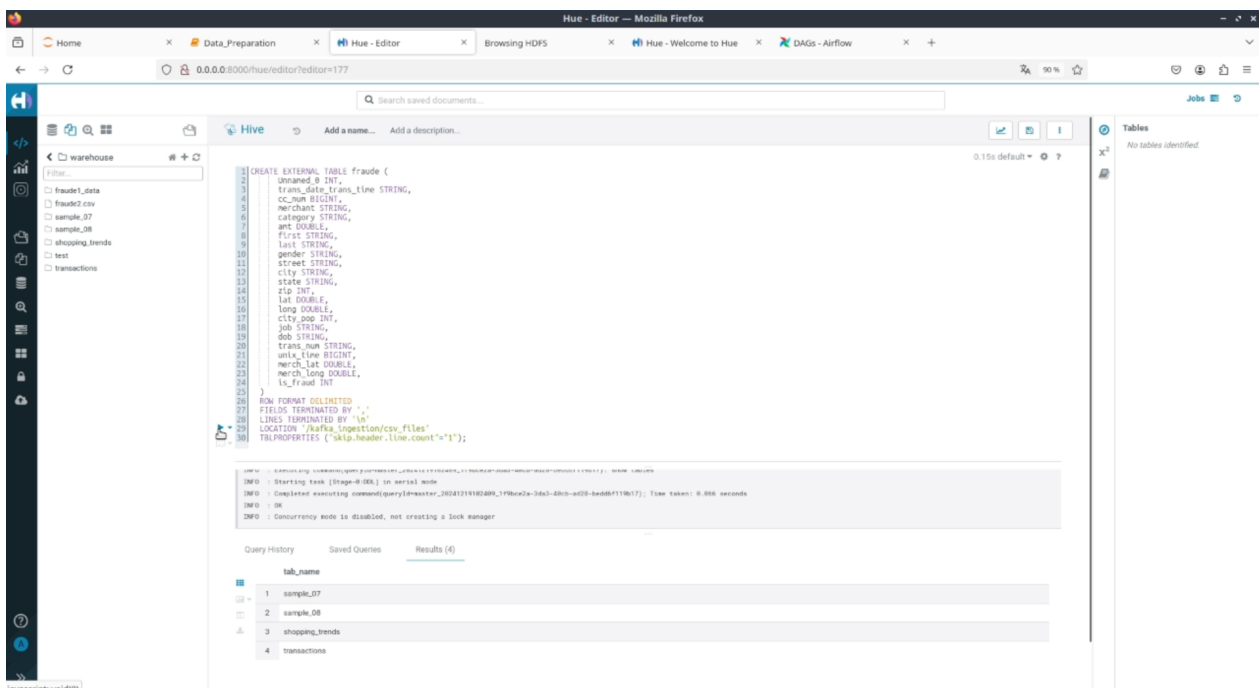


FIGURE 3.10 : création de tableau dans hive

Ce code SQL crée une **table externe** appelée **fraude** dans un système de gestion de bases de données (**Apache Hive**) s fichiers stockés à l'extérieur du système de base de données, comme dans **HDFS** (Hadoop Distributed File System) ou un autre emplacement de stockage.

La structure de la table fraude est définie avec **23 colonnes**, comprenant des informations telles que l'**ID de la transaction**, la **date et l'heure de la transaction**, le **numéro de carte de crédit**, le **commerçant**, la **catégorie**, le **montant de la transaction**, des informations personnelles sur les utilisateurs (**prénom, nom, sexe, adresse**, etc.), ainsi que des informations géographiques et de localisation (**latitude et longitude du commerçant et de l'utilisateur**). La colonne `is_fraud` est utilisée pour indiquer si la transaction est frauduleuse.

Le format des données est spécifié comme étant **délimité**, avec une virgule (,) comme séparateur de champs et un saut de ligne (`\n`) comme délimiteur de lignes. Le chemin d'accès aux fichiers CSV source est défini comme étant `/kafka_ingestion/csv_files`, ce qui indique où les fichiers de données sont stockés. Une propriété de table supplémentaire, `skip.header.line.count`, est définie à **1**, ce qui permet d'ignorer la première ligne du fichier (souvent utilisée pour les en-têtes de colonnes).

### 3.3.3 Division et l'envoi vers pipeline d'ingestion

Une fonction `split_csv_by_parts` qui permet de diviser le fichier CSV en plusieurs morceaux. La fonction prend en entrée le fichier `transactions.csv` (`input_file`), un répertoire de sortie (`output_dir`) et un nombre spécifié de morceaux (`num_parts`). Elle commence par créer le répertoire de sortie si nécessaire, puis elle compte le nombre total de lignes dans le fichier CSV (en excluant l'en-tête). Ensuite, elle calcule la taille de chaque morceau, lit le fichier par morceaux et écrit chaque morceau dans un nouveau fichier CSV. Un message est affiché pour chaque morceau sauvegardé. Dans l'exemple d'utilisation, le fichier `transactions_1.csv` est divisé en 10 morceaux et sauvegardé dans le répertoire `csv_batches`.

```
Chunk 1 sauvegardé : csv_batches/part_1.csv
Chunk 2 sauvegardé : csv_batches/part_2.csv
Chunk 3 sauvegardé : csv_batches/part_3.csv
Chunk 4 sauvegardé : csv_batches/part_4.csv
Chunk 5 sauvegardé : csv_batches/part_5.csv
Chunk 6 sauvegardé : csv_batches/part_6.csv
Chunk 7 sauvegardé : csv_batches/part_7.csv
Chunk 8 sauvegardé : csv_batches/part_8.csv
Chunk 9 sauvegardé : csv_batches/part_9.csv
Chunk 10 sauvegardé : csv_batches/part_10.csv
Chunk 11 sauvegardé : csv_batches/part_11.csv
```

**FIGURE 3.11 :** Division de fichier

Les fichiers CSV d'un répertoire source vers un répertoire de destination (report où kafka lire produce et consume et l'envoie vers hdfs ), un fichier à la fois. Il commence par lister tous les fichiers CSV dans le répertoire source `source_directory` et les trie par ordre alphabétique. Si le répertoire de destination `destination_directory` n'existe pas, il est créé. Ensuite, pour chaque fichier CSV, le programme attend que l'utilisateur appuie sur Entrée avant de déplacer le fichier du répertoire source vers le répertoire de destination. Un message est affiché après chaque transfert pour indiquer que le fichier a été envoyé avec succès. Enfin, une fois tous les fichiers envoyés, un message final confirme que tous les fichiers ont été transférés.

```
Appuyez sur Entrée pour envoyer chaque fichier un par un.
Prêt à envoyer part_1.csv? Appuyez sur Entrée pour continuer...
Fichier part_1.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_10.csv? Appuyez sur Entrée pour continuer...
Fichier part_10.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_11.csv? Appuyez sur Entrée pour continuer...
Fichier part_11.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_2.csv? Appuyez sur Entrée pour continuer...
Fichier part_2.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_3.csv? Appuyez sur Entrée pour continuer...
Fichier part_3.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_4.csv? Appuyez sur Entrée pour continuer...
Fichier part_4.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_5.csv? Appuyez sur Entrée pour continuer...
Fichier part_5.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_6.csv? Appuyez sur Entrée pour continuer...
Fichier part_6.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_7.csv? Appuyez sur Entrée pour continuer...
Fichier part_7.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_8.csv? Appuyez sur Entrée pour continuer...
Fichier part_8.csv envoyé vers /home/master/kafka_ingestion/.
Prêt à envoyer part_9.csv? Appuyez sur Entrée pour continuer...
Fichier part_9.csv envoyé vers /home/master/kafka_ingestion/.
Tous les fichiers ont été envoyés.
```

**FIGURE 3.12 :** Envoi les fichiers

le processus d'intégration et de traitement des fichiers transactions à l'aide d'Apache Kafka, ainsi que leur ingestion dans Hadoop HDFS, permet de gérer efficacement de grandes quantités de données distribuées. Le fichier source est découpé en segments plus petits et chaque segment est envoyé sous forme de message Kafka, garantissant une gestion distribuée et parallèle des données. Kafka agit comme un médiateur entre le producteur et le consommateur, permettant une ingestion continue et fiable des transactions.

Une fois les messages consommés, les données sont transformées et stockées dans HDFS, assurant une gestion et un traitement optimisés des données pour des analyses futures. L'automatisation de ce flux de travail avec des outils comme Apache Airflow renforce l'efficacité du système, garantissant une mise à jour en temps réel des données dans l'écosystème Hadoop.

Ce mécanisme permet non seulement de maintenir l'ordre et la cohérence des données, mais aussi de tirer parti des capacités de scalabilité et de fiabilité de Kafka et Hadoop pour des applications de grande envergure.

### 3.3.4 Execution de pipeline d'ingestion avec airflow

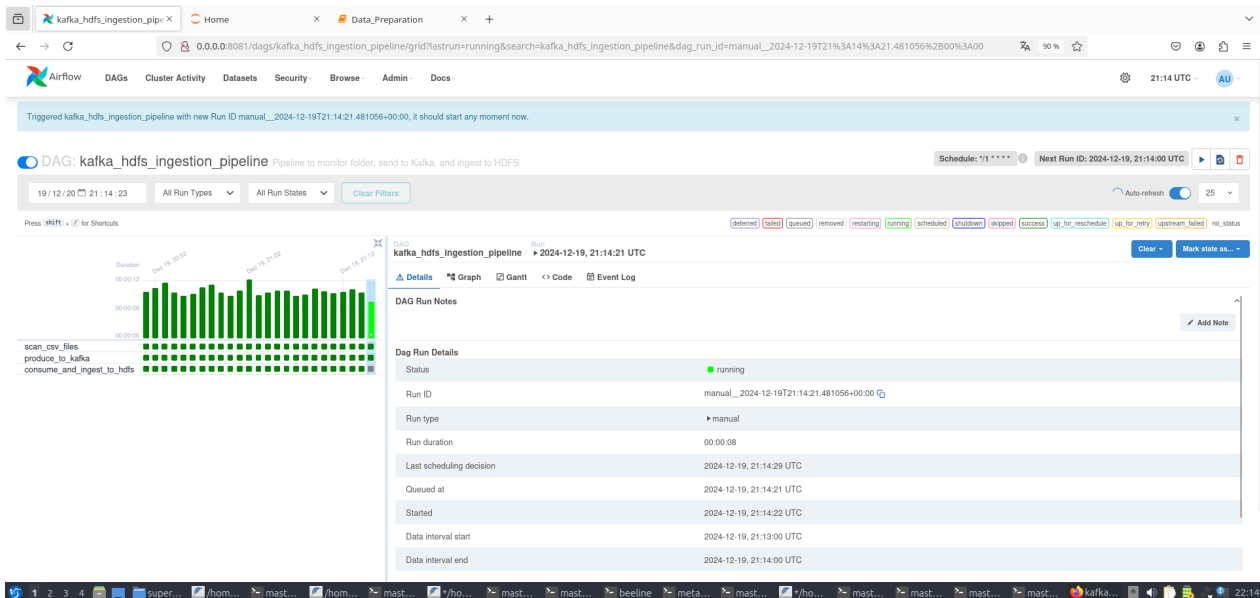


FIGURE 3.13 : Execution de Dag

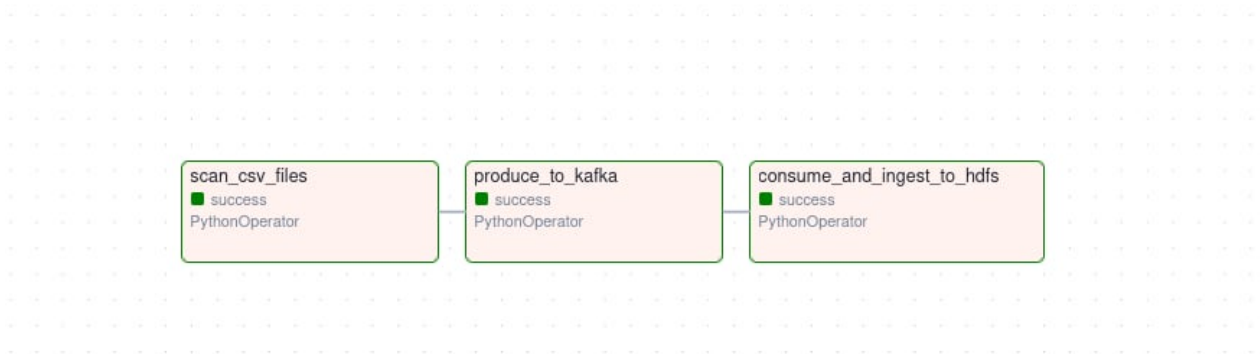


FIGURE 3.14 : Pipeline

L'intégration d'Airflow dans ce processus permet de coordonner toutes les étapes de l'ingestion, du traitement et de la visualisation des données de manière fluide et automatisée. Cela garantit une mise à jour continue des tableaux de bord, permettant aux utilisateurs finaux d'avoir accès aux données les plus récentes en temps réel. Grâce à Kafka, HDFS, Hive, et Apache Superset, ce processus devient hautement évolutif et fiable, capable de gérer de grandes quantités de données tout en minimisant les erreurs humaines.

### 3.3.5 Mise à jour de tableau et des visualisations en temps réelle

#### Avant L'ingestion total des transactions

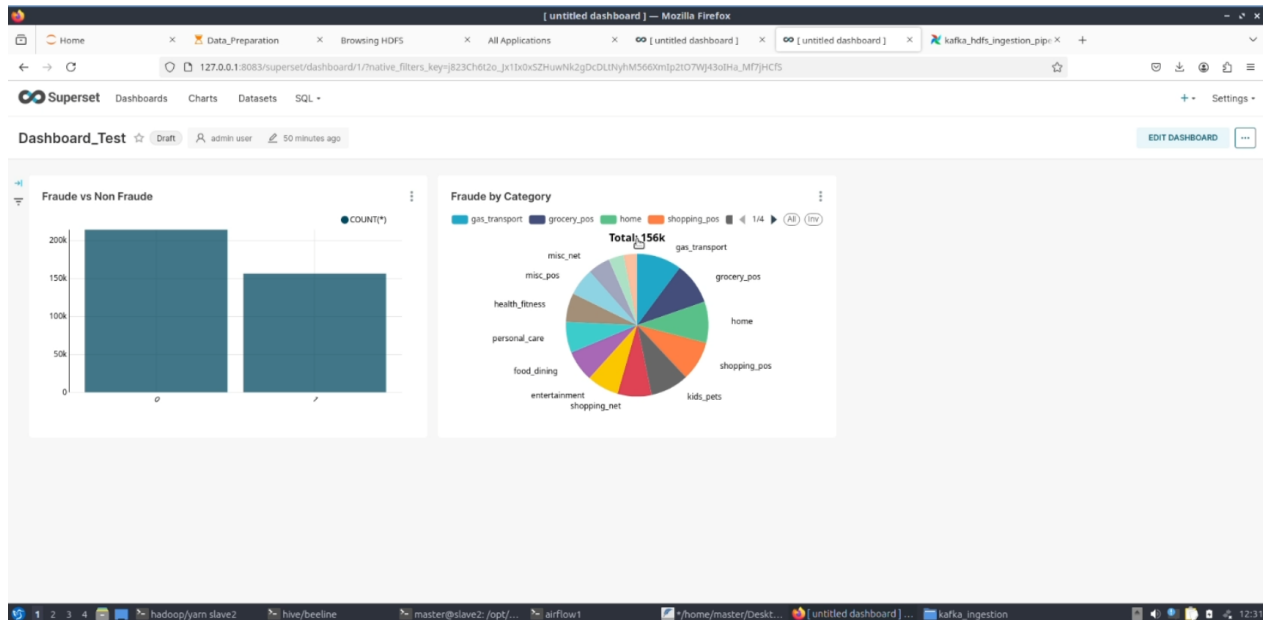


FIGURE 3.15 : Visualisation avant l'ingestion totale de données

#### Après L'ingestion total des transactions

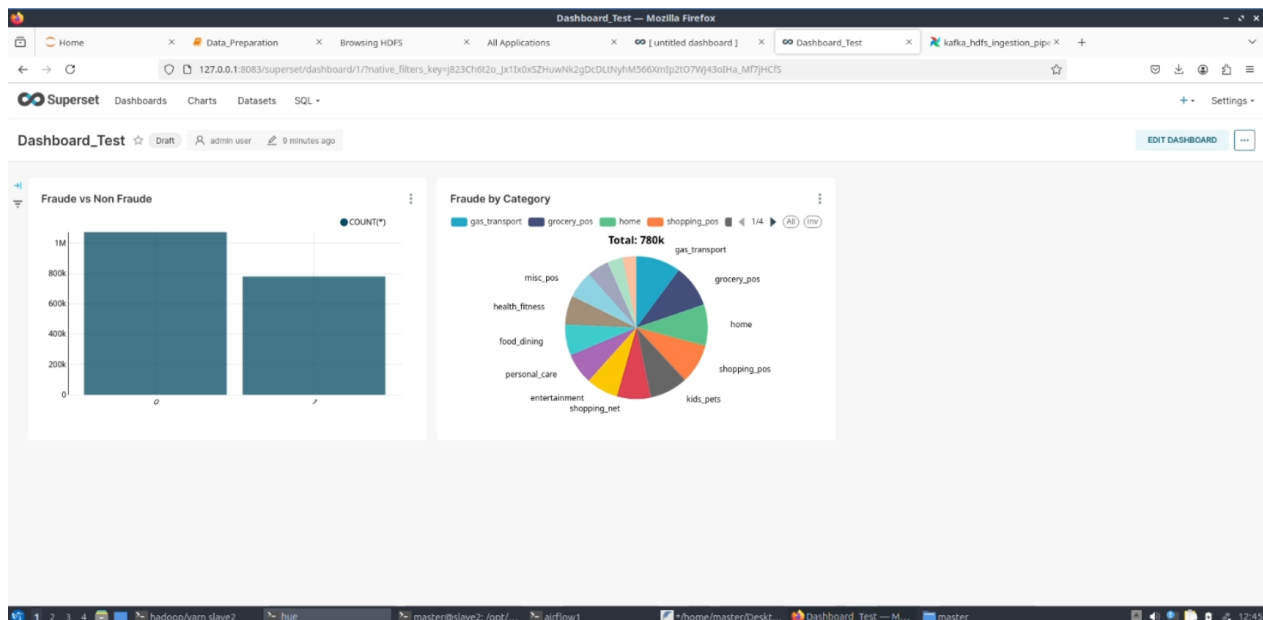


FIGURE 3.16 : Visualisation après l'ingestion totale de données



### 3.3.6 Dashboard Final avec superset

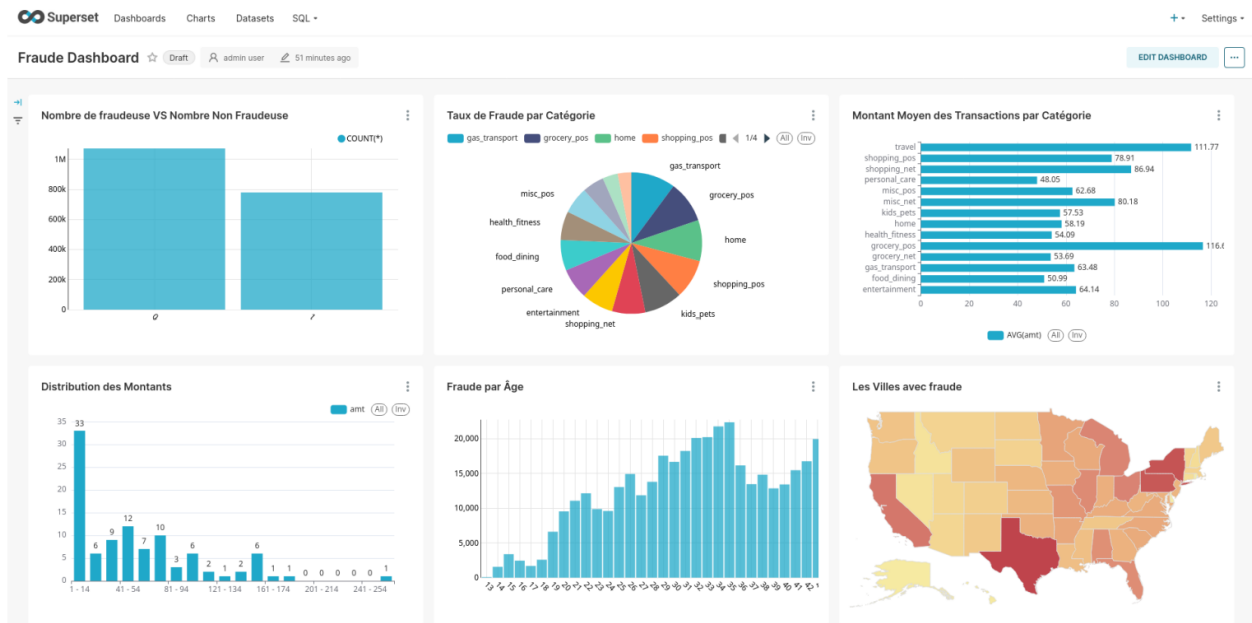


FIGURE 3.17 : Dashboard Final avec superset

## Conclusion Globale

Ce projet a permis de mettre en œuvre une architecture **Big Data** robuste et évolutive pour le traitement, la gestion, et la visualisation des données en temps réel. En intégrant des technologies modernes telles que **Apache Kafka**, **Hadoop** (HDFS, Hive), **Apache Airflow**, et **Apache Superset**, nous avons conçu un pipeline complet répondant aux exigences de traitement des volumes de données massifs tout en assurant une automatisation et une optimisation des performances.

### Bilan du Projet

- **Ingestion et stockage des données** : Grâce à Kafka, les fichiers de transactions sont efficacement collectés et traités, avant d’être stockés dans HDFS, garantissant ainsi une persistance fiable et scalable.
- **Traitement et mise à jour automatisée** : L’utilisation de Hive pour structurer les données a permis d’automatiser les mises à jour des tables et des partitions, facilitant ainsi leur exploitation pour l’analyse.
- **Visualisation en temps réel** : Avec Superset, les données sont rendues accessibles via des tableaux de bord interactifs et dynamiques, offrant une compréhension immédiate et approfondie des tendances et métriques.
- **Orchestration des processus** : Airflow a joué un rôle crucial en orchestrant l’ensemble du pipeline, garantissant la coordination des tâches, la surveillance des erreurs, et la planification des exécutions, réduisant ainsi la complexité opérationnelle.

### Perspectives d’Amélioration

Bien que ce projet ait démontré son efficacité, des améliorations futures pourraient être envisagées :

- **Intégration d’un système de monitoring avancé** (comme Prometheus ou Grafana) pour surveiller la performance des services.
- **Optimisation des performances** des requêtes Hive en utilisant des formats de données comme Parquet ou ORC.
- **Mise en œuvre de mécanismes de sécurité avancés**, comme le chiffrement des données et l’authentification renforcée, pour protéger les données sensibles.

### **Conclusion Finale**

En conclusion, ce projet illustre l'efficacité des technologies **Big Data** pour relever les défis liés au traitement de grandes quantités de données en temps réel. Il offre une solution flexible, scalable et automatisée, tout en posant les bases pour des extensions futures. Cette expérience a également renforcé notre compréhension des systèmes distribués et des outils de traitement des données modernes, offrant une valeur ajoutée significative pour des projets similaires à venir.