

**Filière : ITE3**

**Module : Architecture et Infrastructure Big Data**

**Année universitaire : 2023/2024**

**Travaux Pratique 8 :**

**Apache Kafka+ Spark Streaming + Flume**



**Rédigé par :**  
**HMAYDA Abdessamad**

**Encadré par :**  
**Pr. Kalloubi Fouad**

## Table des matières

<b>I. Apache Kafka :</b>	3
1. Mise en place	3
2. Introduction :	4
a. Création d'un topic « ensademo » :	4
b. S'assurer que le topic a été créé :	4
c. Démarrage du producteur :	4
d. Démarrage du consommateur :	4
e. Affichage du contenu de notre topic :	4
f. Travail à faire :	4
<b>II. Spark streaming :</b>	5
1) Lancement de SparkStreamingContext	5
1)- Tout d'abord allez au terminal de votre namenode et tapez : <b>hdfs dfs -chmod -R 777 /</b>	5
2)-Développement d'une application « wordCount » en utilisant Spark Streaming :	6

## I. Apache Kafka :

### 1. Mise en place

-Cloner ce projet, et exécuter tapez « docker compose up » à l'intérieure du dossier du projet

- tapez « docker ps » pour s'assurer que les conteneurs marchent bien comme l'image d'apres docker desktop montre

**Containers** [Give feedback](#)

Container CPU usage ⓘ  
3.21% / 800% (8 cores allocated)

Container memory usage ⓘ  
1.58GB / 11.23GB [Show charts](#) ▾

🔍 Search ☰ ☷ Only show running containers

<input type="checkbox"/>	Name	Image	Status	Port(s)	Last started	CPU (%)	Actions
<input type="checkbox"/>	tpflume7		Running (5/5)		32 minutes ago	3.21%	■ ⓘ ■
<input type="checkbox"/>	flume-1 5000b595b	bde2020/flume:late	Running	44444:44444	32 minutes ago	0.1%	■ ⓘ ■
<input type="checkbox"/>	namenod 21519636b	apache/hadoop:3.3	Running	9870:9870	32 minutes ago	0.43%	■ ⓘ ■
<input type="checkbox"/>	nodeman edf71f04d1	apache/hadoop:3.3	Running		32 minutes ago	0.68%	■ ⓘ ■
<input type="checkbox"/>	datanode 07383797d	apache/hadoop:3.3	Running		32 minutes ago	0.22%	■ ⓘ ■
<input type="checkbox"/>	resourcer 91e5ed22f	apache/hadoop:3.3	Running	8088:8088	32 minutes ago	1.78%	■ ⓘ ■

- la liste des topics



*Ecole Nationale des Sciences  
Appliquées d'El Jadida*

➔ Puisque le nombre de partitions de `ensaDemoPart` est 2, donc au niveau de « `/tmp/kafka-`  
logs » on aura deux partitions des logs de ce topic, comme indiqué dans l'image suivante :



ii. *Créer un topic avec une facteur de réplication qui vaut 2 :*

Le message d'erreur Facteur de réplication : 2 plus grand que les courtiers(brokers) disponibles : 1 indique qu'on essaye de créer un sujet Kafka avec un facteur de réplication supérieur au nombre de courtiers disponibles dans notre cluster Kafka.

## II. Spark streaming :

## 1) Lancement de SparkStreamingContext

1)- Tout d'abord allez au terminal de votre namenode et tapez : **hdfs dfs-chmod-R 777 /**

2)- lancez le mode interactif de spark :

- ✓ Remarques :
  - Spark streaming est utilisé pour traiter les données en mode streaming, il s'appuie sur un service web appelé StreamingContext (on utilisera netcat dans notre cas pour simuler un service web).
  - StreamingContext exécute un traitement perpétuel de données à intervalles réguliers

- Nous ne pouvons pas avoir plusieurs contextes en meme temps, donc on doit arrêter sparkContexte avant de lancer sparkStreaming

3) -arrêtez SparkContext et lancez SparkStreamingContext :

```
scala> :paste
// Entering paste mode (ctrl-D to finish)

import org.apache.spark.streaming._
import org.apache.spark.SparkConf
sc.stop
val conf=new SparkConf().setAppName("Streaming").setMaster("yarn")
val ssc=new StreamingContext(conf,Seconds(10))

// Exiting paste mode, now interpreting.

23/12/09 00:24:29 WARN Client: Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
import org.apache.spark.streaming._
import org.apache.spark.SparkConf
conf: org.apache.spark.SparkConf = org.apache.spark.SparkConf@78ceb8cd
ssc: org.apache.spark.streaming.StreamingContext = org.apache.spark.streaming.StreamingContext@a4ee49e

scala>
```

## 2)-Développement d'une application « wordCount » en utilisant Spark Streaming :

Tout d'abord nous allons créer une application Spark en utilisant Scala(2.11.0) et sbt (1.3.3) et apres importer spark et spark-streaming

```
SparkStreaming.scala  build.sbt x

ThisBuild / version := "v1"

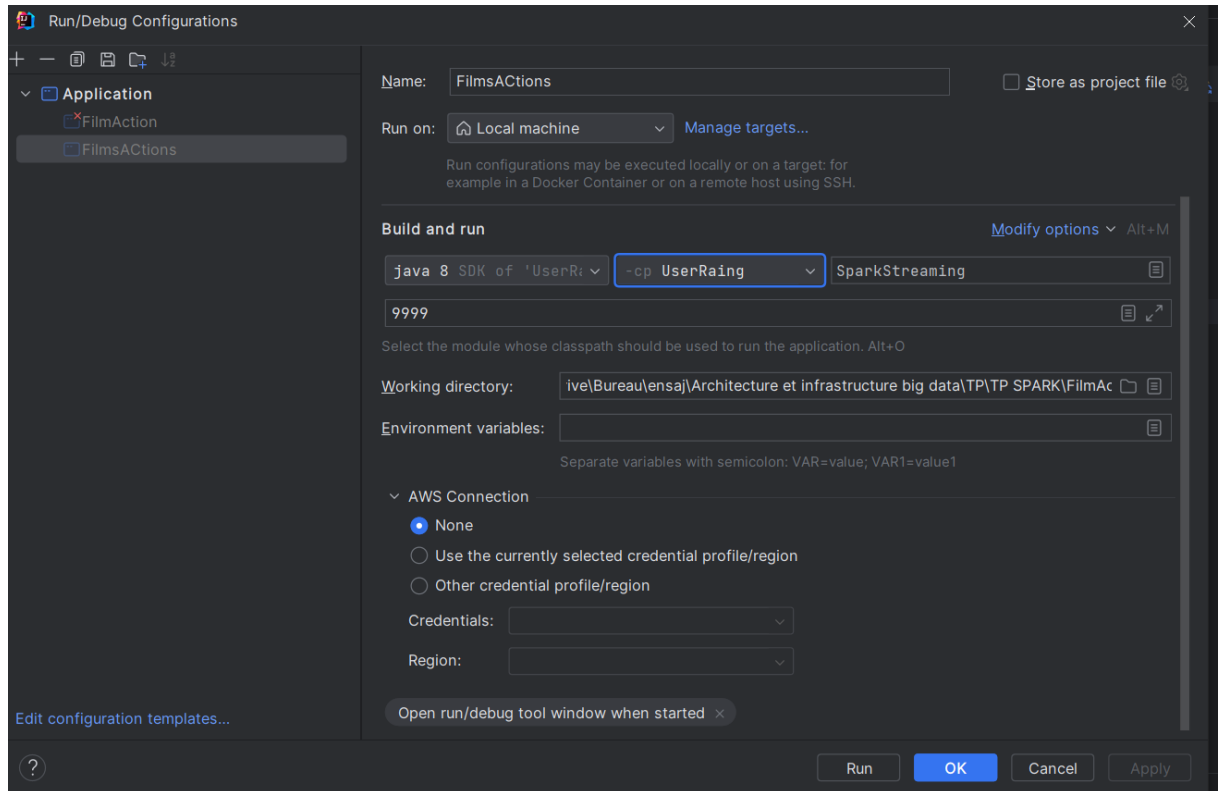
2
3 ThisBuild / scalaVersion := "2.11.0"
4
5
6 libraryDependencies += "org.apache.spark" %% "spark-core" % "1.6.0"
7 libraryDependencies += "org.apache.spark" %% "spark-streaming" % "1.6.0"
8 libraryDependencies += "org.apache.logging.log4j" % "log4j-api" % "2.17.1"
9 libraryDependencies += "org.apache.logging.log4j" % "log4j-core" % "2.17.1"
10
11 lazy val root = (project in file("."))
12   .settings(
13     name := "SparkStreaming"
14   )
15
```

Voici le code de notre application

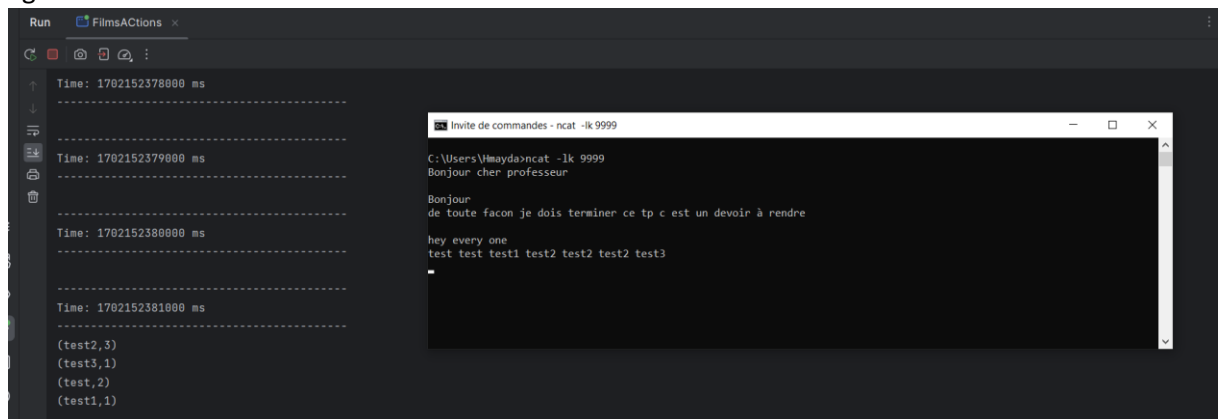
```
1 import org.apache.log4j.{Level, Logger}
2 import org.apache.spark._
3 import org.apache.spark.streaming._
4
5 object SparkStreaming{
6   def main (args : Array[String]) : Unit = {
7     if (args.length != 1) {
8       println("Usage: NetcatReceiver <port>")
9       System.exit( status = 1)
10    }
11
12    Logger.getLogger ( name = "org").setLevel(Level.ERROR)
13
14    val port=args(0).toInt
15    val conf = new SparkConf().setAppName("ActionMovies").setMaster("local[*]")
16    val ssc=new StreamingContext(conf,Seconds(1))
17
18    val lines = ssc.socketTextStream( hostname = "localhost", port)
19
20    val words=lines.flatMap(_.split( regex = " "))
21    val pairs=words.map(word=>(word,1))
22    val wordCounts=pairs.reduceByKey(_+_ )
23    wordCounts.print()
24    ssc.start()
25    ssc.awaitTermination()
26  }
27 }
```

SparkStreaming > main(args: Array[String])

Exécuter cette application en passant 9999 comme argument

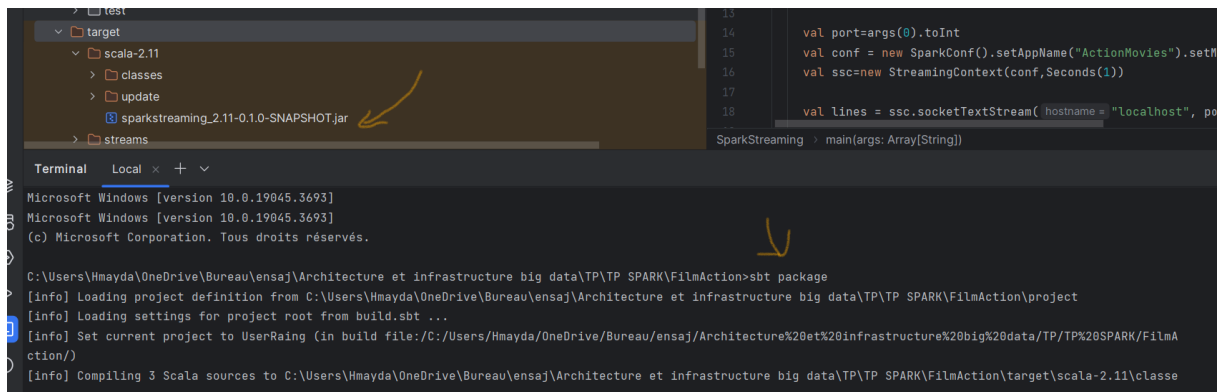


Également nous allons simuler un service web en utilisant netcat :



-donc notre application repond à nos besoin , on doit maintenant contsruire le jar de cette application :

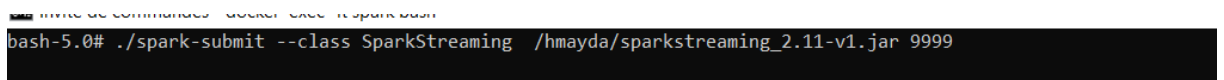




- Ensuite, on doit mettre le jar dans HDFS, pour le faire déplacer vers le dossier du volume hmayda et exécuter cette commande de hdfs :

```
sh-4.2$ hdfs dfs -put "/hmayda/sparkstreaming_2.11-v1.jar" /hmayda/
WARNING: log4j.properties is not found. HADOOP_CONF_DIR may be incomplete.
sh-4.2$
```

- Lancez le jar avec spark-submit



- utiliser netcat et envoyez des messages sur le port 9999 et voici le résultat du traitement effectué par l'application spark

