

# Homework 3

Giuliano di Giuseppe, Raffaele Russo, Alessandro Vanacore

## Indice

<b>1</b>	<b>Traccia</b>	<b>2</b>
<b>2</b>	<b>Generazione dei dati</b>	<b>2</b>
<b>3</b>	<b>Setup Kafka</b>	<b>2</b>
3.1	Producer . . . . .	2
3.2	Consumer . . . . .	2
3.2.1	Setup Hadoop . . . . .	3
<b>4</b>	<b>Dashboard</b>	<b>4</b>

# 1 Traccia

Tramite Apache Kafka gestire l'acquisizione di uno stream di dati (e.g., messaggi da un web server, misure da una rete di sensori, etc.) da una data sorgente (anche simulata) e la successiva memorizzazione in un file di log su HDFS. Utilizzando poi uno degli strumenti del primo homework (pyspark, pig o hive) effettuare delle query sullo stream di dati mostrandone i risultati su un'apposita dashboard.

# 2 Generazione dei dati

Per la generazione dei dati si utilizza la libreria python *psutils*, con la quale si estraggono statistiche in tempo reale dalla cpu, memoria e disco.

# 3 Setup Kafka

Il primo passo è avviare il server zookeeper e il server kafka con i comandi:

```
$ bin/zookeeper-server-start.sh config/zookeeper.properties
$ bin/kafka-server-start.sh config/server.properties
```

## 3.1 Producer

Si procede alla creazione di un topic e di un producer che può inviare messaggi su questo topic. A tal proposito si utilizza lo script python precedente per automatizzare le operazioni. Il generico messaggio contiene informazioni sul nome del topic, timestamp e statistiche sulla macchina.

Per la creazione del topic e dell'invio da parte del producer su tale topic si utilizza il codice rispettivamente in figura 1 e 2.

```
def connect(name_topic):
    cmds = ["cd /opt/kafka_2.11-0.9.0.0"
            ,"/opt/kafka_2.11-0.9.0.0/bin/kafka-topics.sh --create --zookeeper localhost:2181 --replication-factor 1 --partitions 1 --topic "+name_topic
            ]

    try:
        for cmd in cmds:
            result = subprocess.Popen(cmd, shell=True)
            print('Comando eseguito correttamente.')
            print('Err:\n', result.stderr)

    except subprocess.CalledProcessError as ex:
        print('Esecuzione fallita.')
        print('Err:', ex.stderr)
```

Figura 1: Creazione topic

## 3.2 Consumer

Si utilizza un consumer che legge dal topic precedente e salva tale informazioni su un file di log di Hadoop. A tal proposito, dopo aver configurato correttamente Hadoop, si utilizza lo script bash in figura 3. Quest'ultimo avvia un producer che legge dal topic precedente e salva tale contenuto sul file di log specificato.

```
while(1):
    cmd_producer = "/opt/kafka_2.11-0.9.0.0/bin/kafka-console-producer.sh --broker-list localhost:9092 --topic "+name_topic
    cmd_full = "echo "+ "\n"+name_topic+" "+get_stats()+'\n' | ' + cmd_producer
    cmd_mess = subprocess.run(cmd_full, shell=True, capture_output=True, text=True)
    print(f"Message {i} sent")
    sleep(0.1)
    i+=1
    if i==5000:
        break
```

Figura 2: Invio producer su topic specificato

```
#!/bin/bash

name_topic="launchv5"
hdfs_log_file="/home/parallels/hadoop/logs/hadoop-parallel-namenode-ubuntu-linux-22-04-desktop.log"
command="/opt/kafka_2.11-0.9.0.0/bin/kafka-console-consumer.sh --zookeeper localhost:2181 --topic $name_topic --from-beginning >> $hdfs_log_file"

eval $command
```

Figura 3: Script bash consumer

### 3.2.1 Setup Hadoop

Per la configurazione di Hadoop si usano i comandi in figura 4. Il primo comando esegue la formattazione del NameNode di Hadoop, mentre il secondo comando avvia l'HDFS.

```
subprocess.run("/home/parallels/hadoop/bin/hdfs namenode -format", shell = True)
subprocess.run("/home/parallels/hadoop/sbin/start-dfs.sh", shell = True)
```

Figura 4: Hadoop

In figura 5 sono mostrate cinque righe salvate dal consumer nel file di log di Hadoop contenenti rispettivamente nome del topic, timestamp e statistiche della macchina.

```
parallels@ubuntu-linux-22-04-desktop: /home/parallels $ cat /home/parallels/namenode-ubuntu-linux-22-04-desktop.log | tail -5
launchv5 ('timestamp': '2023-05-17 16:32:35', 'CPU Physical cores': 2, 'CPU Total cores': 2, 'CPU Core0': 6.1, 'CPU Core1': 8.1, 'RAM Total': '3.82GB', 'RAM Available': '672.25MB', 'RAM Used': '2.72GB', 'RAM Percentage': 82.8, 'RAM SWAPTotal': '2.00GB', 'RAM SWAPFree': '177.89MB', 'RAM SWAPUsed': '1.83GB', 'RAM SWAPPercentage': 91.3, 'DISK Total read': '301.76GB', 'DISK Total write': '4.48GB')
launchv5 ('timestamp': '2023-05-17 16:32:37', 'CPU Physical cores': 2, 'CPU Total cores': 2, 'CPU Core0': 5.1, 'CPU Core1': 4.1, 'RAM Total': '3.82GB', 'RAM Available': '670.95MB', 'RAM Used': '2.72GB', 'RAM Percentage': 82.8, 'RAM SWAPTotal': '2.00GB', 'RAM SWAPFree': '177.89MB', 'RAM SWAPUsed': '1.83GB', 'RAM SWAPPercentage': 91.3, 'DISK Total read': '301.76GB', 'DISK Total write': '4.48GB')
launchv5 ('timestamp': '2023-05-17 16:32:39', 'CPU Physical cores': 2, 'CPU Total cores': 2, 'CPU Core0': 6.1, 'CPU Core1': 5.1, 'RAM Total': '3.82GB', 'RAM Available': '673.60MB', 'RAM Used': '2.72GB', 'RAM Percentage': 82.8, 'RAM SWAPTotal': '2.00GB', 'RAM SWAPFree': '177.89MB', 'RAM SWAPUsed': '1.83GB', 'RAM SWAPPercentage': 91.3, 'DISK Total read': '301.76GB', 'DISK Total write': '4.48GB')
launchv5 ('timestamp': '2023-05-17 16:32:41', 'CPU Physical cores': 2, 'CPU Total cores': 2, 'CPU Core0': 10.2, 'CPU Core1': 8.2, 'RAM Total': '3.82GB', 'RAM Available': '675.23MB', 'RAM Used': '2.71GB', 'RAM Percentage': 82.7, 'RAM SWAPTotal': '2.00GB', 'RAM SWAPFree': '177.91MB', 'RAM SWAPUsed': '1.83GB', 'RAM SWAPPercentage': 91.3, 'DISK Total read': '301.76GB', 'DISK Total write': '4.48GB')
launchv5 ('timestamp': '2023-05-17 16:32:43', 'CPU Physical cores': 2, 'CPU Total cores': 2, 'CPU Core0': 6.1, 'CPU Core1': 5.1, 'RAM Total': '3.82GB', 'RAM Available': '673.10MB', 'RAM Used': '2.72GB', 'RAM Percentage': 82.8, 'RAM SWAPTotal': '2.00GB', 'RAM SWAPFree': '177.91MB', 'RAM SWAPUsed': '1.83GB', 'RAM SWAPPercentage': 91.3, 'DISK Total read': '301.76GB', 'DISK Total write': '4.48GB')
```

Figura 5: 5 righe del file di log

## 4 Dashboard

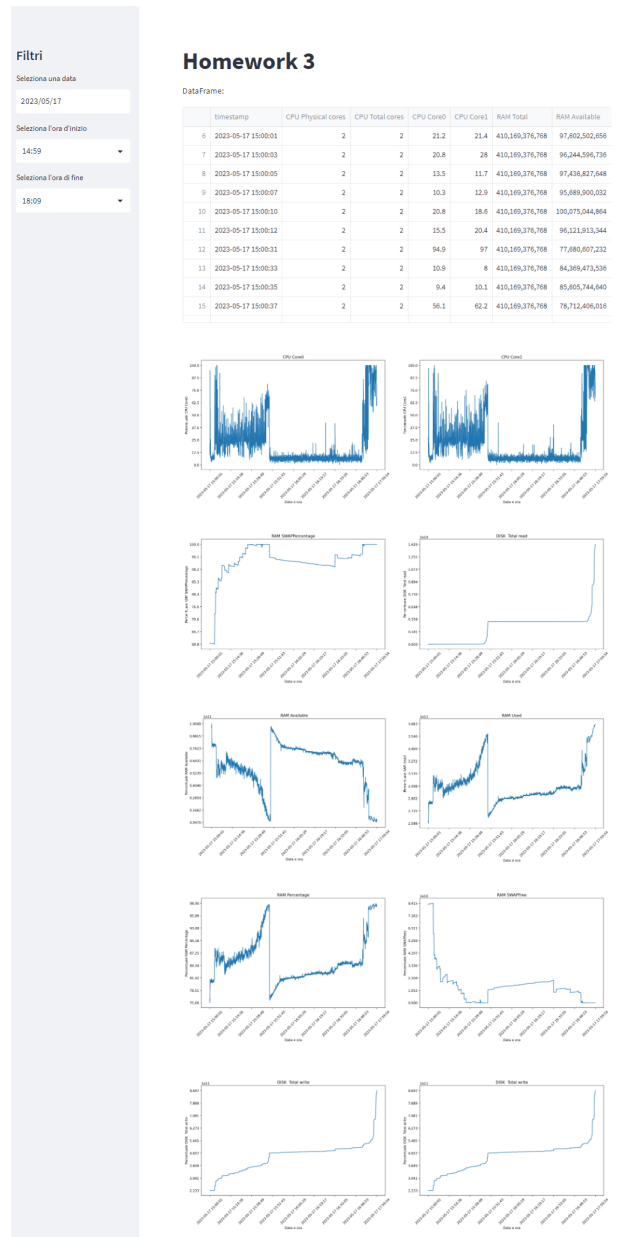


Figura 6: Dashboard