

Spark Lab 2 Report

Building a Streaming Data Pipeline with Apache Spark

Contents


Starting the Cluster	2
Task 1	3
Task 2	4
Task 3	6

André Plancha
andre.plancha@hotmail.com

December 06th, 2025


Starting the Cluster

Dockerfile on the server

 Docker


```
1 FROM apache/spark:4.0.1
2 # switch to root to install packages
3 USER root
4 RUN pip install --no-cache-dir "pandas==2.3.2" "pyarrow==21.0.0"
5 # switch back to spark user
6 USER spark
```

compose.yaml on the server

 YAML

```
1 services:
2   spark:
3     build: .
4     hostname: apache-spark
5     ports:
6       - "7077:7077"      # Spark master port
7       - "8080:8080"      # Spark master web UI
8       - "8081:8081"      # Spark worker web UI
9       - "15002:15002"     # Spark Connect server port
10      - "4040:4040"       # Spark Connect web UI
11     command: >
12       bash -c "/opt/spark/sbin/start-master.sh;
13               /opt/spark/sbin/start-connect-server.sh;
14               /opt/spark/sbin/start-worker.sh spark://192.168.1.7:7077;
15               sleep infinity"
```

```
1 $ ls
2 compose.yaml  Dockerfile
3 $ docker compose up -d
4 [+] Running 1/1
5   ✓ Container spark-docker-spark-1  Started
   0.5s
```

 Shell

```
In [1]: 1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import udf
3
4 spark = SparkSession.builder \
5     .remote("sc://192.168.1.7:15002") \
6     .appName("UDFTransformation") \
7     .config("spark.sql.ansi.enabled", "false") \
8     .config("spark.sql.execution.pythonUDF.arrow.enabled", "true") \
9     .getOrCreate()
10
11 # limit() shows a nice HTML table in Jupyter, while show() prints plain text
12 spark.conf.set('spark.sql.repl.eagerEval.enabled', True)
13
14 spark
```

Python

```
Out[1]: 1 <pyspark.sql.connect.session.SparkSession at 0x24595594cd0>
```

```
In [2]: 1 DataStreamReader = spark.readStream \
2     .format("socket") \
3     .option("host", "192.168.1.7") \
4     .option("port", 9999)
5
6 linesDF = DataStreamReader.load()
7 lines_writer = linesDF.writeStream \
8     .outputMode("append") \
9     .format("memory") \
10    .queryName("lines_table")
11
12 DataStreamReader
```

Python

```
Out[2]: 1 <pyspark.sql.connect.streaming.readwriter.DataStreamReader at 0x245955953f0>
```

```
In [3]: 1 def stop_queries():
2     for query in spark.streams.active:
3         query.stop()
```

Python

Task 1

When you receive an sentences from socket, put your name before each word in sentences , and put last name after that(use UDF)

```
In [ ]: 1 @udf
```

Python

```

2 def name_transform(sentence, first="André", last="Plancha"):
3     words = sentence.split()
4     return ','.join([f"{first}{word}{last}" for word in words])
5
6 transformedDF = linesDF.withColumn("value", name_transform(linesDF.value))
7
8
9 stop_queries()
10
11 task1_writer =
12 transformedDF.writeStream.outputMode("append").format("memory").queryName("task
13 task1_query = task1_writer.start()
14 task1_writer

```

Out[]: 1 <pyspark.sql.connect.streaming.readwriter.DataStreamWriter at 0x1932daf3ee0>

```

1 ssh plancha@192.168.1.7 -t /usr/bin/nc -lk 9999
2 plancha@192.168.1.7's password:
3 Sweden is a good country
4 Finland is a better country
5 Denmark is an awesome country
6 I love programming in Spark
7 Data is the new oil

```

console

In []:

```

1 task1_query.stop()
2
3 spark.sql("SELECT * FROM task1_table").show(20, truncate=False)

```

Python

```

1 +-----+
2 |value|
3 +-----+
4 |AndréSwedenPlancha,AndréisPlancha,AndréaPlancha,AndrégoodPlancha,AndrécountryPlancha|
5 |AndréFinlandPlancha,AndréisPlancha,AndréaPlancha,AndrébetterPlancha,AndrécountryPlancha|
6 |AndréDenmarkPlancha,AndréisPlancha,AndréanPlancha,AndréawesomePlancha,AndrécountryPlancha|
7 |AndréIPlancha,AndrélovePlancha,AndréprogrammingPlancha,AndréinPlancha,AndréSparkPlancha|
8 |AndréDataPlancha,AndréisPlancha,AndréthePlancha,AndrénewPlancha,AndréoilPlancha|
9 +-----+
10

```

Task 2


Average word length from all the words received from socket. Use SQL.

In [41]:

```

1  stop_queries()
2
3  linesDF.createOrReplaceTempView("lines_table")
4
5  task2_DF = spark.sql("""
6      from lines_table |>
7      select value, split(value, ' ') as phrase |>
8      select explode(phrase) as word |>
9      aggregate avg(length(word)) as average_word_length
10  """)
11  task2_writer =
12  task2_DF.writeStream.outputMode("complete").format("memory").queryName("task2_t
13  task2_query = task2_writer.start()
14  task2_writer

```


 Python

Out[41]: 1 <pyspark.sql.connect.streaming.readwriter.DataStreamWriter at 0x245a7fc76d0>

```

1  ssh plancha@192.168.1.7 -t /usr/bin/nc -lk 9999
2  plancha@192.168.1.7's password:
3  Sweden is a good country
4  Finland is a better country
5  Denmark is an awesome country
6  I love programming in Spark
7  Data is the new oil


```

 console

```

In [39]: 1  lines = [
2      "Sweden is a good country",
3      "Finland is a better country",
4      "Denmark is an awesome country",
5      "I love programming in Spark",
6      "Data is the new oil"
7  ]
8  words = " ".join(lines).split(" ")
9  lengths = [len(word) for word in words]
10  average_length = sum(lengths) / len(lengths)
11  average_length

```

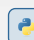
 Python

Out[39]: 1 4.24

```

In [42]: 1  task2_query.stop()
2  spark.sql("SELECT * FROM task2_table").show(20, truncate=False)

```

 Python

```

1  +-----+
2  |average_word_length|

```

```
3 +-----+
4 |4.24      |
5 +-----+
6
```

Task 3

Use SQL to filter only the even numbers from the streaming data.

In [59]:

```
1 stop_queries()
2
3 linesDF.createOrReplaceTempView("lines_table")
4
5 task3_DF = spark.sql("""
6     from lines_table |>
7     select cast(value as int) as value |>
8     where value % 2 == 0
9 """)
10 task3_writer =
11 task3_DF.writeStream.outputMode("append").format("memory").queryName("task3_tab
12 task3_query = task3_writer.start()
13 task3_writer
```

Out[59]:

```
1 <pyspark.sql.connect.streaming.readwriter.DataStreamWriter at 0x245a7fc62c0>
```

```
1 $ ssh plancha@192.168.1.7 -t /usr/bin/nc -lk 9999
2 plancha@192.168.1.7's password:
3 3413
4 12315
5 1231
6 121516
7 574
8 2325
9 252
10 5785
11 23123
12 54756712
13 312312
14 31
15 543646856
16 231425
```

In [51]:

```
1 numbers = [
```

```

2     3413,
3     12315,
4     1231,
5     121516,
6     574,
7     2325,
8     252,
9     5785,
10    23123,
11    54756712,
12    312312,
13    31,
14    543646856,
15    231425
16 ]
17 [num for num in numbers if num % 2 == 0]

```


Out[51]: 1 [121516, 574, 252, 54756712, 312312, 543646856]

In [60]:

```

1 task3_query.stop()
2 spark.sql("SELECT * FROM task3_table").show(20, truncate=False)

```

 Python

```

1 +-----+
2 |value   |
3 +-----+
4 |574     |
5 |543646856|
6 |54756712 |
7 |252     |
8 |312312  |
9 |121516  |
10 +-----+
11

```