# Streaming Data Processing TD 3

Minh NGUYEN

November 6th, 2023

**Abstract**

In this programming assignment, we were supposed to get used to the concepts of Apache Kafka and get the first hands-on using a preconfigured Kafka version, named Redpanda.

# 1 Summary

# 2 Introduction

## 2.1 Prerequisites

- Programming language and SDK:

  - Python 3
  - Java 11/13/15/17

- IDE such as IntelliJ (Java) or VS Code / PyCharm with Python.

- Spark 3.5.0

# 3 Install Spark

## 3.1 Other than Windows:

- Download Spark from this link: spark-3.5.0-bin-hadoop3.tgz.

- Extract the downloaded compressed file to a folder as you like.

- Config the PATH environment:

  - If you're using MacOS: do the following steps:
    1. Open ~/.zshrc by an editor (can use Vim)
    2. Add the following lines in the end of the file:
    ```
    export SPARK_HOME="<your-path-to-spark-folder>/spark-3.5.0-bin-hadoop3"
    export PATH="$SPARK_HOME:$SPARK_HOME/bin:$PATH"
    ```
    3. Apply this change by running this command in the terminal:
    ```
    source ~/.zshrc
    ```
    4. Testing by opening another terminal and running this command. If Spark shell appears, then the installation can be said as done.
    ```
    spark-shell
    ```

```
[(base)                          spark-3.5.0-bin-hadoop3 % spark-shell            ]
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLeve
l(newLevel).
23/11/06 10:25:22 WARN NativeCodeLoader: Unable to load native-hadoop library fo
r your platform... using builtin-java classes where applicable
Spark context Web UI available at http://10.10.27.112:4040
Spark context available as 'sc' (master = local[*], app id = local-1699262723042
).
Spark session available as 'spark'.
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /___/ .__/\_,_/_/ /_/\_\   version 3.5.0
      /_/

Using Scala version 2.12.18 (OpenJDK 64-Bit Server VM, Java 17.0.9)
Type in expressions to have them evaluated.
Type :help for more information.

scala>
```

Figure 1: Spark shell interface

# 4   Setup environment for Python

A virtual environment is a directory that contains a specific collection of Python packages that you have installed. For example, you may have one environment with NumPy 1.7 and its dependencies, and another environment with NumPy 1.6 for legacy testing. If you change one environment, your other environments are not affected. You can easily activate or deactivate environments, which is how you switch between them. Working with separate environments help you:

- Profesionalize your workflow, in which Python dependencies can be fulfilled inside an environement
- Separate one environement from another, to avoid the dependencies conflict.

To create virtual environments, you can use either **venv** or **conda**. In this LAB, we will use conda.

## 4.1   Install Conda

- Install Miniconda from this link: https://docs.conda.io/projects/miniconda/en/latest/

## 4.2   Managing environments

Full list of tutorial can be found here :

Creating an environment with required packages (pyspark):

```
1  conda create -n <environment_name> python=3.10 pyspark
```

Activate the environment:

```
1  conda activate <environment_name>
```

# 5   LAB 3

## 5.1   Context

In this LAB we will try to connect Spark Structured Streaming with your previous Kafka topic (from generated data in LAB 2) and do basic task: WordCount.

**EXERCISE 1:** Use the code sample from here to connect with your Kafka cluster and Kafka topic in the previous LAB, and do the wordcount.

**Notes:**

```
1  """
2  Consumes messages from one or more topics in Kafka and does wordcount.
3   Usage:
4
5  python structured_kafka_wordcount.py <bootstrap-servers> <subscribe-type> <topics>
6     <bootstrap-servers> The Kafka "bootstrap.servers" configuration. A
7     comma-separated list of host:port.
8     <subscribe-type> There are three kinds of type, i.e. 'assign', 'subscribe',
9     'subscribePattern'.
10    |- <assign> Specific TopicPartitions to consume. Json string
11    |   {"topicA":[0,1],"topicB":[2,4]}.
12    |- <subscribe> The topic list to subscribe. A comma-separated list of
13    |   topics.
14    |- <subscribePattern> The pattern used to subscribe to topic(s).
15    |   Java regex string.
16    |- Only one of "assign, "subscribe" or "subscribePattern" options can be
17    |   specified for Kafka source.
18    <topics> Different value format depends on the value of 'subscribe-type'.
19  """
```