

# DATAComp2 - Project Spark & ELASTICSEARCH

Yacine Mokhtari  
Lilia Izri  
Alexandre Combeau

May 4, 2022

## 1 Introduction and Motivations

The goal of our project is to make a short program to process and analyze real-time stream data (in this case, tweets) using Sparklab Streaming[3], index some of this data with ElasticSearch[1] and evaluate some queries with it. Our project will be divided in three main parts that we will detail in this report.



Figure 1: Schema

## 2 Processing Twitter Data

### 2.1 Settings: Defining Twitter developer credentials and Authenticating the connection

The first step was to create a [Twitter Application](#)[2] so as we can get both the Consumer Keys and the Authentication Tokens. Using the `tweepy` library[4], we can retrieve Tweets from the [Twitter API](#) and store them into sockets before making any analysis on them as illustrated in figure1. For this project, we focused on tweets containing the keyword `netflix`.

### 2.2 Stream Listener Class

In order to get the tweets in real-time, we have to implement a new class that inherits from the `Stream` module of `Tweepy`. The constructor takes a `Socket` object and calls the super constructor, pretty straightforward. Then, we overwrite the `on_data`[5] method to adapt it to our needs.

```
# We combine data and metadata to send them
# We add a key ###:field:### so we can split the fields easily
# We remove '\n' from a tweet and put one '\n' between tweets
tweet_info = ("###:field:### user: " + user + " ###:field:### tweet: " + text + "\n"
              "###:field:### date: " + date + " ###:field:### location: " + location + "\n"
              "###:field:### hashtags: " + hashtags)
              .replace('\n', ' ') + '\n'

print(tweet['user']['location'])
print(tweet_info)

# Send to socket : We convert this tweet into a bite code (since spark takes easily this kind of data)
self.client_socket.send(tweet_info.encode('utf-8'))
```

Figure 2:

The `on_data` method will format the received JSON file (a tweet) into a string. We only keep some features such as the name, the location, hashtags (if any) and the date. These different fields are delimited by the tag `###field### field_name:`.

Note that the `\n` is replaced by a simple space `" "`. So we make sure that each tweet corresponds to just one line and it's correctly read when creating the DStream in spark with the function `socketTextStream()`.

### 2.3 Establish connexion with client

---

## 3 Spark Analysis

### 3.1 Process data

We process the data we receive by first, splitting the incoming string so each row in our datastream `'tweets'` be a tuple where each element will correspond to a field of the tweet received. And using `textblob`, we perform sentiment analysis on the text of the tweet and we add this sentiment to the previous tuple. So each tweet will be represented by a tuple of the form (user, text, date, location, hashtags, sentiment).

Where sentiment equals -1 if the tweet is more likely negative. It will equal 1 if it's more positive, and 0 otherwise.

**Remark:** We chosed to take into account the polarity returned by the `TextBlob` `text processing library` and to ignore the subjectivity.

## 4 Elasticsearch Indexing

## References

- [1] E. Elasticsearch. Official documentation. <https://www.elastic.co/guide/en/elasticsearch/reference/current/index.html>.
- [2] T. Inc. Twitter developer portal. <https://dev.twitter.com/apps/new>.
- [3] S. Streaming. Official documentation. <https://spark.apache.org/docs/latest/streaming-programming-guide.html>.
- [4] Tweepy. Index - official documentation. <https://docs.tweepy.org/en/stable/index.html>.
- [5] Tweepy. Streaming. <https://docs.tweepy.org/en/stable/streaming.html>.