

# Assignment 3

---

Spring 2020

Due Date: Sunday Apr 19, 2020

## Instructions

- This assignment will involve writing code in the form of Scala classes. The code for the 2 parts should be in different classes in a single Scala Spark project. The project should be compiled into a jar file that can run on AWS cluster. You should include build.sbt for dependency management.
- All instructions for compiling and running your code must be placed in the README file.
- You should use a cover sheet, which can be downloaded from [here](#)
- You are allowed to work in pairs i.e. a group of two students is allowed. Please write the names of the group members on the cover page.
- **You have a total of 4 free late days for the entire semester. You can use at most 2 days for any one assignment. After four days have been used up, there will be a penalty of 10% for each late day. The submission for this assignment will be closed 2 days after the due date.**
- Please ask all questions on Piazza, not via email.

# 1 Spark Streaming with Twitter and Kafka

In this part, you will create a Spark Streaming application that will continuously read data from Twitter about a topic. These Twitter feeds will be analyzed for their sentiment, and then analyzed using ElasticSearch. To exchange data between these two, you will use Kafka as a broker. Everything will be done locally on your computer. Below are the steps to get started:

## Setting up your development environment

You will need to perform the following steps to get your environment set up.

- You will need to create a Twitter app and get credentials. It can be done at <https://apps.twitter.com/>
- Download Apache Kafka and go through the quickstart steps: <https://kafka.apache.org/quickstart>
- Windows users might want to consider WSL, which gives a Unix-like environment on Windows: <https://docs.microsoft.com/en-us/windows/wsl/install-win10>
- I have provided a Spark Scala project that gives an example of how to connect Spark Structured Streaming to a Kafka broker. The project can be downloaded from <http://www.utdallas.edu/~axn112530/cs6350/kafka.zip>.  
The next section gives you the instructions on how to compile and run the project.
- Later, you will also need to set up Elasticsearch and Kibana environment to visualize data. You will need to download Elasticsearch, Kibana, and Logstash from <https://www.elastic.co/downloads/>

## Running the Example

The example project can be downloaded from:

<http://www.utdallas.edu/~axn112530/cs6350/kafka.zip>

It illustrates a simple way to communicate between Structured Streaming and Kafka. Below are the steps:

- Make sure that you first start Zookeeper and then Apache Kafka. After ensuring that these two services are running, create two topics, which we will call topicA and topicB. After this, use the producer utility to generate some messages on topicA, and start a consumer on topicB.
- For building this project, you need to have Scala Build Tool (SBT) installed. It can be downloaded from <https://www.scala-sbt.org>  
After downloading, you should be able to go to root of your project and build by the following commands.

```
sbt
> assembly
```

This will create a fat jar (i.e. a jar containing all its dependencies) under **target/scala-2.11** directory.

- Run the Spark project by using the following command:

```
spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.0
--class ClassName PathToJarFile topicA
```

The above will accept data from topicA, perform word count, and display the results on the console. This is accomplished by the Spark Structured Streaming application.

## Starting the Project

For this project, you will need to perform the following steps:

- Create a Twitter application using Spark and Scala. An example of how to do this is available here:  
<http://bahir.apache.org/docs/spark/current/spark-streaming-twitter/>
- The application should perform search about a topic and gather Tweets related to it. The next step is sentiment evaluation of the Tweets. Various libraries are available for this task:
  - [CoreNLP Scala Examples](#)
  - [Databricks example](#) - Note: you can't use Databricks for this assignment.
  - [Sentiment Analysis in Scala](#)
  - [Sentiment Analysis of Social Media Posts using Spark](#)

The sentiment evaluation should happen continuously using a stream approach. At the end of every window, a message containing the **sentiment** should be sent to Kafka through a topic.

- In the next step, you will configure Logstash, Elasticsearch, and Kibana to read from the topic and set up visualization of sentiment.

## Visualizing the data using Elasticsearch and Kibana

You were able to read the data using a console consumer in a previous step. Let's now move forward and visualize the data using Elasticsearch. To send data from Kafka to Elasticsearch, you need a processing pipeline that is provided by [Logstash](#). Download Elasticsearch, Kibana, and Logstash from <https://www.elastic.co/downloads>.

After downloading the data, you need to go to the appropriate directories, and start the services in the following order:

1. **Elasticsearch** using the following command in the elasticsearch-5.6.8/bin directory:  
`./elasticsearch`
2. **Kibana** using the following command in the kibana-5.6.8-darwin-x86\_64/bin directory:  
`./kibana`

3. **Logstash:** Go to the logstash-5.6.8 directory and create a file logstash-simple.conf with following content:

```
input {
  kafka {
    bootstrap_servers => "localhost:9092"
    topics => ["YourTopic"]
  }
}

output {
  elasticsearch {
    hosts => ["localhost:9200"]
    index => "YourTopic-index"
  }
}
```

Then run the following command

```
bin/logstash -f logstash-simple.conf
```

This sets up the right pipeline between Kafka and Elasticsearch.

If everything is set up properly, you should be able to go to <http://localhost:5601> and use Kibana to visualize your data in real-time. You will have to search for the appropriate topic index, which is **YourTopic-index** in the example shown above. You can read more about Kibana [here](#).

Note: If you have trouble setting up all of this on your local computer, please use online sources for help with troubleshooting.

## What to Submit

You are required to submit the following:

- Your project in a compressed zip format. Please delete the .jar file from the target folder as it will be quite large.
- A graphical plot showing how the sentiment for your topic varied over a time interval. You are free to choose the time interval - but it should be at least a few hours.
- A brief summary of how the sentiment varied over the time frame, and what insights does it give you.
- A README file indicating how to run your code. Please do not use any hard coded paths - all paths should be specified as arguments.

## 2 Analyzing Social Networks using GraphX

In this part, you will use Spark GraphX to analyze social network data. You are free to choose any one of the Social network datasets available from the [SNAP repository](#).

You will use this dataset to construct a GraphX graph and run some queries and algorithms on the graph. You will write a class in Scala that will read in the data using an argument, and perform following steps:

### 2.1 Loading Data

Load the data into a dataframe or RDD using Spark. Define a parser so that you can identify and extract relevant fields. Note that GraphX edges are directed, so if your dataset has undirected relationships, you might need to convert those into 2 directed relationships. That is, if your dataset contains an undirected friendship relationship between X and Y, then you might need to create 2 edges one from X to Y and the other from Y to X.

### 2.2 Create Graphs

Define edge and vertex structure and create property graphs.

### 2.3 Running Queries

Run the following queries using the GraphX API and write your output to a file specified by the output parameter.

- a. Find the top 5 nodes with the highest outdegree and find the count of the number of outgoing edges in each
- b. Find the top 5 nodes with the highest indegree and find the count of the number of incoming edges in each
- c. Calculate PageRank for each of the nodes and output the top 5 nodes with the highest PageRank values. You are free to define the threshold parameter.
- d. Run the connected components algorithm on it and find the top 5 components with the largest number of nodes.
- e. Run the triangle counts algorithm on each of the vertices and output the top 5 vertices with the largest triangle count. In case of ties, you can randomly select the top 5 vertices.

## 2.4 What to submit

You are required to submit the following:

- Your project in a compressed zip format. There should be two arguments: input dataset path, and output dataset path
- Output file with the results of the queries
- A brief summary indicating whether your results make sense and what insights you get by this analysis.
- A README file indicating how to run your code. Please do not use any hard coded paths - all paths should be specified as arguments.