



**Department of Computer Science and Engineering
University of Puerto Rico
Mayagüez Campus**

CIIC 8995/5995 5016 – Big Data Analytics Spring 2017

Project 2: Streaming for Twitter Analysis Due Date: April 27, 2017, 11:59 PM

Objectives

1. Use Spark Streaming, SparkSQL, and Kafka to analyze trends contained in a collection of tweets.
2. Become familiar with streaming concepts

Overview

You will design, implement and test a series of programs that will analyze live tweets. Your solution will:

1. Capture the tweets from the Tweet Stream API
<https://dev.twitter.com/streaming/overview>

For this purpose, you can use:

- python twitter (<https://pypi.python.org/pypi/twitter>)
 - pip install twitter
- tweetpy (<http://www.tweepy.org/>)

2. Put the tweets into Kafka
3. Read the tweets from Kafka with Spark Streaming
4. Use Spark, Spark Streaming, Hive, and SparkSQL to implement the following operations:
 - a. Capture the top 10 trending hashtags (most viewed) in the last 60 minutes.
Refresh every 10 minutes
 - b. Capture the top 10 trending keywords (most viewed) in the last 60 minutes
Refresh every 10 minutes. (No stop words here)
 - c. Capture the top 10 participants (most tweets posted) in the last 12 hours
Refresh every hour
5. Count the number of occurrences for these keyword, in intervals of 1 hours, on each day,
 - a. Trump
 - b. MAGA

- c. Dictator
- d. Impeach
- e. Drain
- f. Swamp

You Must accumulate statistic for at least 3 days

Your solution will consist of a collection of Python programs, and SQL queries that perform tasks 1-5.

Visualization

Provide a means to visualize the results of the tasks 1-5, using the D3.js library. You are free to use the charts that you think best fits the visualization.

Deliverables

- **GitHub repo with all the code**

Grading

- **Project will be graded via demonstration of working code, running in cluster mode, forked from GitHub repo.**

PROJECT DUE DATE: 11:59 PM – April 27, 2017.