# Twitter Streaming with Spark

Spark Streaming allows us to analyze continuous streams of data. Eg:- Processing log data from websites or servers. This data is aggregated and analyzed at some given interval. We can feed this data to some port like Amazon Kinesis, HDFS, Kafka, Flume, etc.  The 'checkpointing' feature stores data to disk periodically for fault tolerance.

Just like Spark has Dataframes and Datasets, Spark Streaming has 'Dstream' objects, which breaks up a stream of continuous data into distinct RDDs. AN advantage to this is that we can apply a lot of concepts of Dataframes to Dstreams as well.

Streamin can be done on a local port or on any particular app. One thing to keep in mind is that when streaming, our RDDs may contain very less data. For this we can increase the timeframe of our processing through 'windowed operations', which combine results from multiple batches over a sliding time window.

We can also keep track of an event across many batches through time by using 'updateStateByKey()' function.

In this project, we will stream Twitter and keep track of few of the most popular hashtags (#) as we go.

**Creating a Twitter App:**

In order to get data from Twitter, we need to create an app on twitter and accessing its data through our code. It involves the following;

1. Creating the App on Twitter. I have created an App called 'TwitterStreaminWithSpark'. With each app comes a Consumer (or API) Key and Secret. This will be useful for us while establishing connection with Twitter in our code.

# TwitterStreamingWithSpark

| Details | Settings | Keys and Access Tokens | Permissions |

## Application Settings

*Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.*

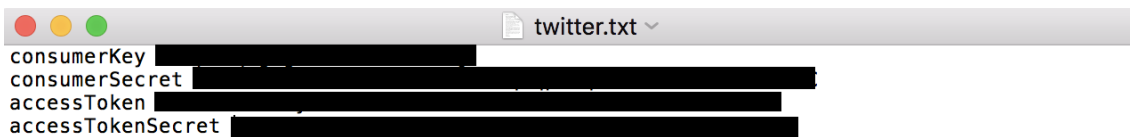| | |
|---|---|
| Consumer Key (API Key) | ██████████ |
| Consumer Secret (API Secret) | █████████████ |
| Access Level | Read and write (modify app permissions) |
| Owner | adi_gogoi |
| Owner ID | ████████ |

2. After creating an app, we need to create an Access Token for our app to function. This token consists of an Access Token and Secret. Keep this recorded for future use.



**Your Access Token**

This access token can be used to make API requests on your own account's behalf. Do not share your access token secret with anyone.

| | |
|---|---|
| Access Token | ███████████████████████ |
| Access Token Secret | ████████████████████████ |
| Access Level | Read and write |
| Owner | adi_gogoi |
| Owner ID | ████████ |

3. Next we need to create the configuration file for Twitter, which will consist of our Consumer and Access Tokens and Keys. This will be reference for our program when it will be accessing Twitter's Streaming data.
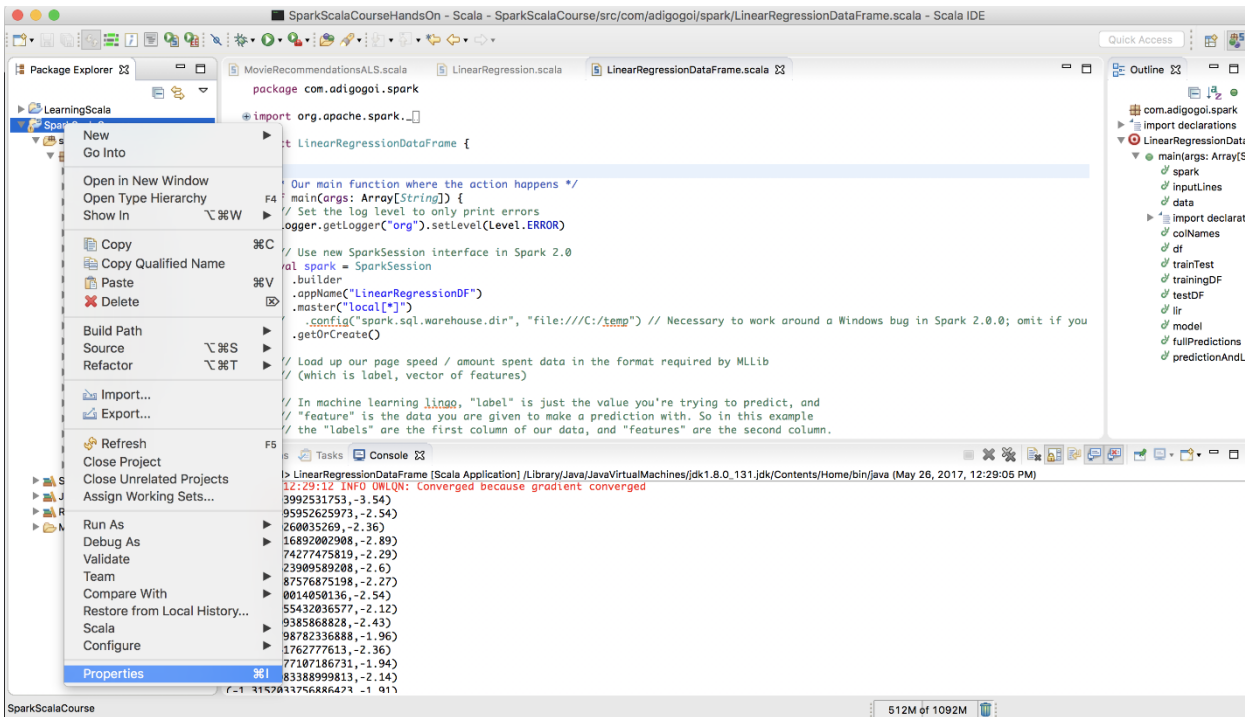


```
consumerKey ███████████████████████████
consumerSecret █████████████████████████████████████████████
accessToken ████████████████████████████████████
accessTokenSecret ██████████████████████████████
```

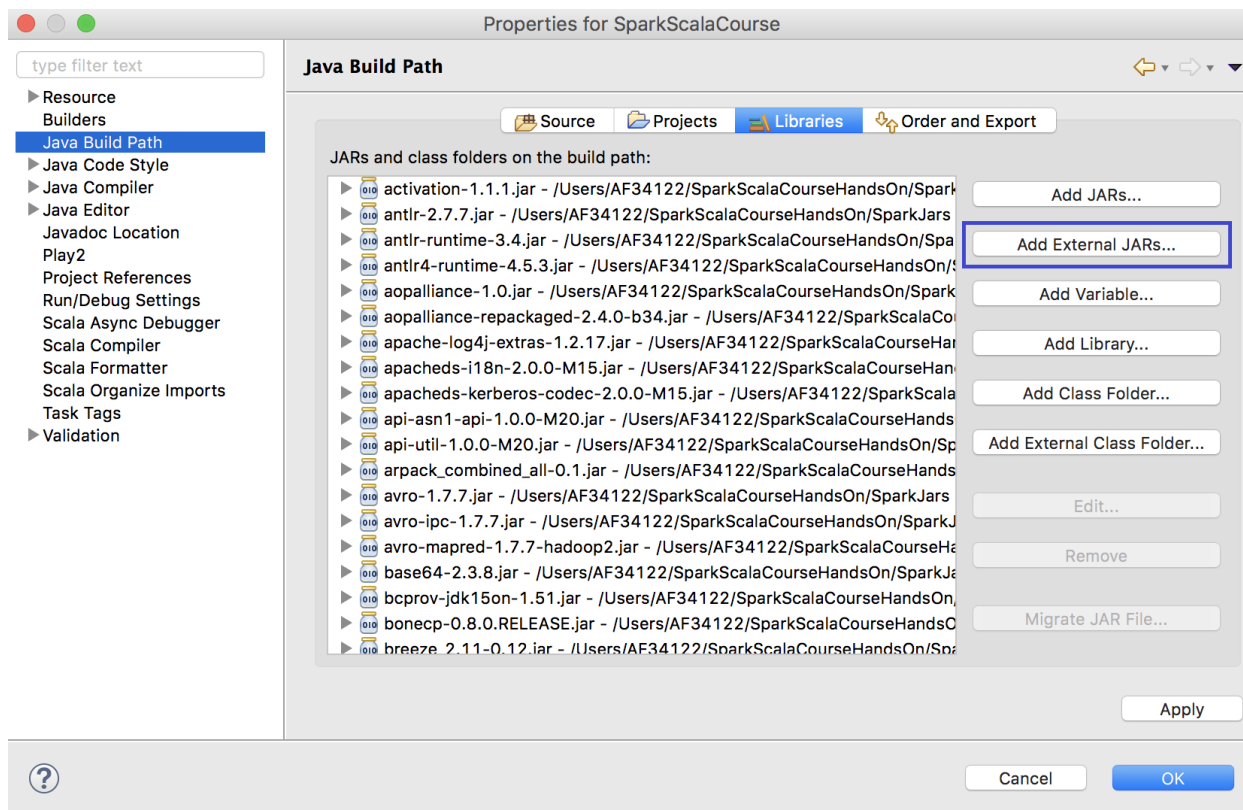**Including Twitter Jars**

Working with twitter is not possible until we include libraries that allow Spark to integrate with the Twitter API. This can be done by including external jars in our project.
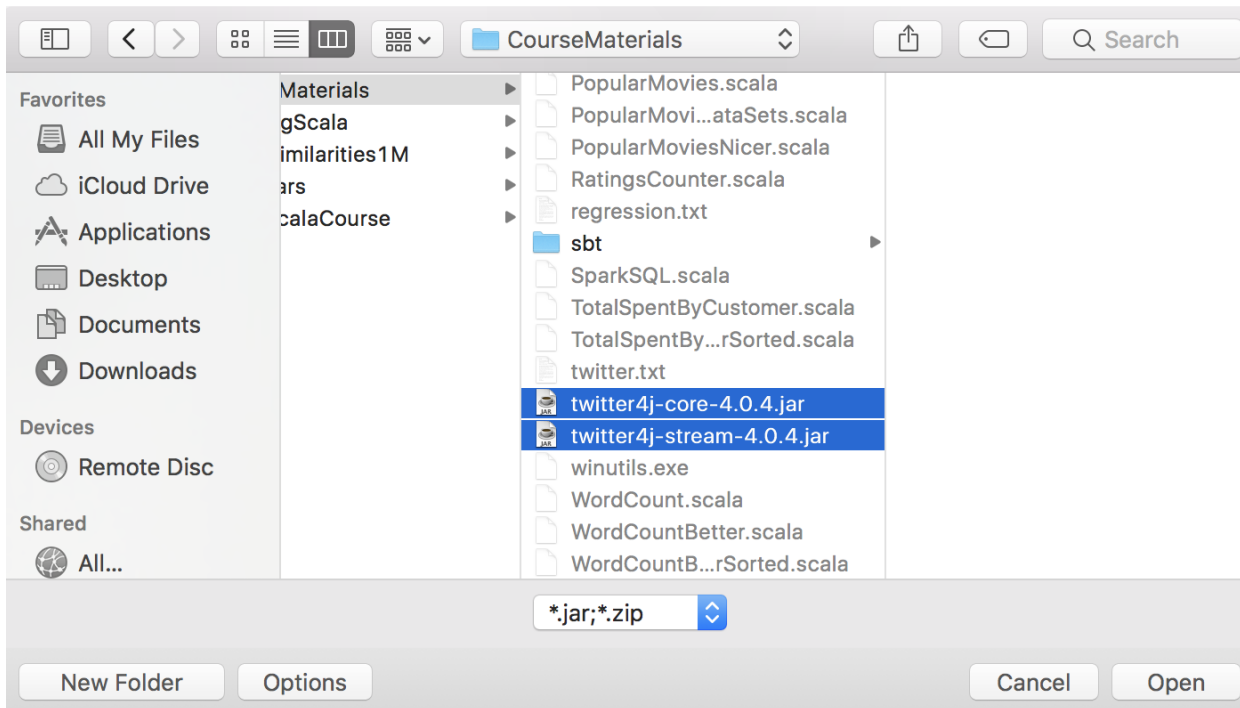
1. Go to the Properties of the Project folder.

2. Got to 'Java Build Path' and then 'Add External Jars'.



3. Select the three JARs necessary for Twitter Streaming i.e. dstream-twitter_2.11-0.1.0-SNAPSHOT, twitter4j-core-4.0.4, and twitter4j-stream-4.0.4 and click Open.

**Twitter Streaming**

After including the necessary jars, we can Stream Twitter data to find the most popular Hashtags. But before that, we should know how our code works:

1. We use the createStream function of TwitterUtils to create the Dstream and store it in 'tweets' variable.

2. We then map the status of the tweets in this Dstream and get them stored in the variable 'statuses'.

3. Using flatMap, we split the tweets into words and store it in 'tweetwords'.

4. We then use filter to get just those words that begin with '#' and store it in the variable 'hashtags'.

5. Then we can convert these hashtags into key-value pairs, with the hashtags being the key and their count being the value. This involves first mapping each hashtag with 1 and then reduce with count.

6. Next we sort the key-value pairs in descending order and display them.

Another major part of the code is the 'setuptwitter' function, which reads the configuration file for twitter, gets the keys and sets up the System for Twitter Authentication.

We also need to note that unlike the usual SparkContext, we need to use the StreamingContext for this project. We also need to setup a checkpoint folder before starting the context for Spark Streaming.

Let's check the first few screenshots of the most popular hashtags of our 5 second window.

**Package Explorer** ⌷

- ▸ LearningScala
- ▾ SparkScalaCourse
  - ▾ src
    - ▾ com.adigogoi.spark
      - ▸ DataFrames.scala
      - ▸ DegreesOfSeparation.scala
      - ▸ FriendsByAge.scala
      - GraphX.scala
      - ▸ LinearRegression.scala
      - ▸ LinearRegressionDataFrame.s
      - ▸ MaxTemperatures.scala
      - ▸ MinTemperatures.scala
      - ▸ MostPopularSuperhero.scala
      - ▸ MovieRecommendationsALS.
      - ▸ MovieSimilarities.scala
      - ▸ MovieSimilarities1M.scala
      - ▸ MyCustomerSpending.scala
      - ▸ PopularHashtags.scala
      - ▸ PopularMovies.scala
      - ▸ PopularMoviesDataSets.scala
      - ▸ PopularMoviesNicer.scala
      - ▸ RatingsCounter.scala
      - ▸ SparkSQL.scala
      - ▸ WordCount.scala
      - ▸ WordCountBetter.scala
      - ▸ WordCountBetterSorted.scal
  - ▸ Scala Library container [ 2.11.8 ]
  - ▸ JRE System Library [JavaSE-1.8]
  - ▸ Referenced Libraries
  - ▸ MovieData

Tabs: MovieRecommendationsALS.scala | LinearRegression.scala | LinearRegressionDataFrame.scala | PopularHashtags.scala | GraphX.scala

```scala
package com.adigogoi.spark

import org.apache.spark._

/** Listens to a stream of Tweets and keeps track of the most popular
object PopularHashtags {

  /** Makes sure only ERROR messages get logged to avoid log spam. */
  def setupLogging() = {
    import org.apache.log4j.{Level, Logger}
    val rootLogger = Logger.getRootLogger()
    rootLogger.setLevel(Level.ERROR)
  }

  /** Configures Twitter service credentials using twiter.txt in the main workspace directory */
  def setupTwitter() = {
    import scala.io.Source

    for (line <- Source.fromFile("../SparkScalaCourse/Datasets/twitter.txt").getLines) {
      val fields = line.split(" ")
      if (fields.length == 2) {
        System.setProperty("twitter4j.oauth." + fields(0), fields(1))
      }
    }
  }

  /** Our main function where the action happens */
  def main(args: Array[String]) {
```

**Outline** ⌷

- ⬢ com.adigogoi.spark
- ▸ import declarations
- ▾ PopularHashtags
  - ▾ setupLogging()
    - ▸ import declaratio
    - ⬥ rootLogger
  - ▾ setupTwitter()
    - ▸ import declaratio
  - ▾ main(args: Array[Stri
    - ⬥ ssc
    - ⬥ tweets
    - ⬥ statuses
    - ⬥ tweetwords
    - ⬥ hashtags
    - ⬥ hashtagKeyValue
    - ⬥ hashtagCounts
    - ⬥ sortedResults

**Problems** | **Tasks** | **Console** ⌷

PopularHashtags [Scala Application] /Library/Java/JavaVirtualMachines/jdk1.8.0_131.jdk/Contents/Home/bin/java (May 27, 2017, 6:59:42 PM)
```
(#VzlaVenceLaCensura,2)
(#Debbie4Congress,1)
(#태풍,1)
...

-----------------------------------------
Time: 1495925991000 ms
-----------------------------------------
(#MTVFRESCAKIKA,2)
(#MTVACTAPPCAELI,2)
(#MTVCRUSHREYES,2)
(#PREMIOSMTVMIAW,2)
(#MTVSUPERCAEYOSS,2)
(#MTVEXPOPREYES,2)
(#MTVINSTAMXCANTU,2)
(#VzlaVenceLaCensura,2)
(#Debbie4Congress,1)
(#태풍,1)
...
```

778M of 1061M 🗑

**Problems** | **Tasks** | **Console** ⌷

PopularHashtags [Scala Application] /Library/Java/JavaVirtualMachines/jdk1.8.0_131.jdk/Contents/Home/bin/java (May 27, 2017, 6:59:42 PM)
```
(#MTVINSTAMXCANTU,4)
(#MTVFRESCAKIKA,3)
(#MTVCRUSHREYES,3)
...

-----------------------------------------
Time: 1495926007000 ms
-----------------------------------------
(#PREMIOSMTVMIAW,10)
(#YoungAwardsMX17,9)
(#DiscoPreferido,9)
(#CantantePreferidaDelAño,9)
(#MTVACTAPPCAELI,8)
(#MTVSUPERCAEYOSS,8)
(#MTVEXPOPREYES,8)
(#MTVINSTAMXCANTU,8)
(#MTVFRESCAKIKA,7)
(#MTVCRUSHREYES,7)
...
```

455M of 1057M 🗑

**Problems** | **Tasks** | **Console** ⌷

PopularHashtags [Scala Application] /Library/Java/JavaVirtualMachines/jdk1.8.0_131.jdk/Contents/Home/bin/java (May 27, 2017, 6:59:42 PM)
```
(#YoungAwardsMX17,9)
(#DiscoPreferido,9)
(#CantantePreferidaDelAño,9)
...

-----------------------------------------
Time: 1495926026000 ms
-----------------------------------------
(#PREMIOSMTVMIAW,19)
(#MTVACTAPPCAELI,12)
(#MTVSUPERCAEYOSS,12)
(#MTVEXPOPREYES,12)
(#MTVINSTAMXCANTU,12)
(#MTVFRESCAKIKA,11)
(#MTVCRUSHREYES,11)
(#YoungAwardsMX17,9)
(#DiscoPreferido,9)
(#CantantePreferidaDelAño,9)
...
```

463M of 1057M 🗑

We can see from the output of Twitter Streaming that some hashtags of MTV are trending quite a lot now.

This project helped me realize that Spark Streaming is a great tool not just to store streaming data, but also to perform real-time analytics on the data and give insights.