# CIS5560 Term Project Tutorial

**Authors:** Tom Cruise; Mel Gibson

**Instructor:** Jongwook Woo

**Date: 05/18/2017**

# Lab Tutorial

yourname (yourname@calstatela.edu)

06/10/2016

# Yelp Data Analysis using Spark (your Title)

## Objectives

**List what your objectives are.** In this hands-on lab, you will learn how to:

- Get data manually using REST API

- Create Spark cluster

- Train NLP system

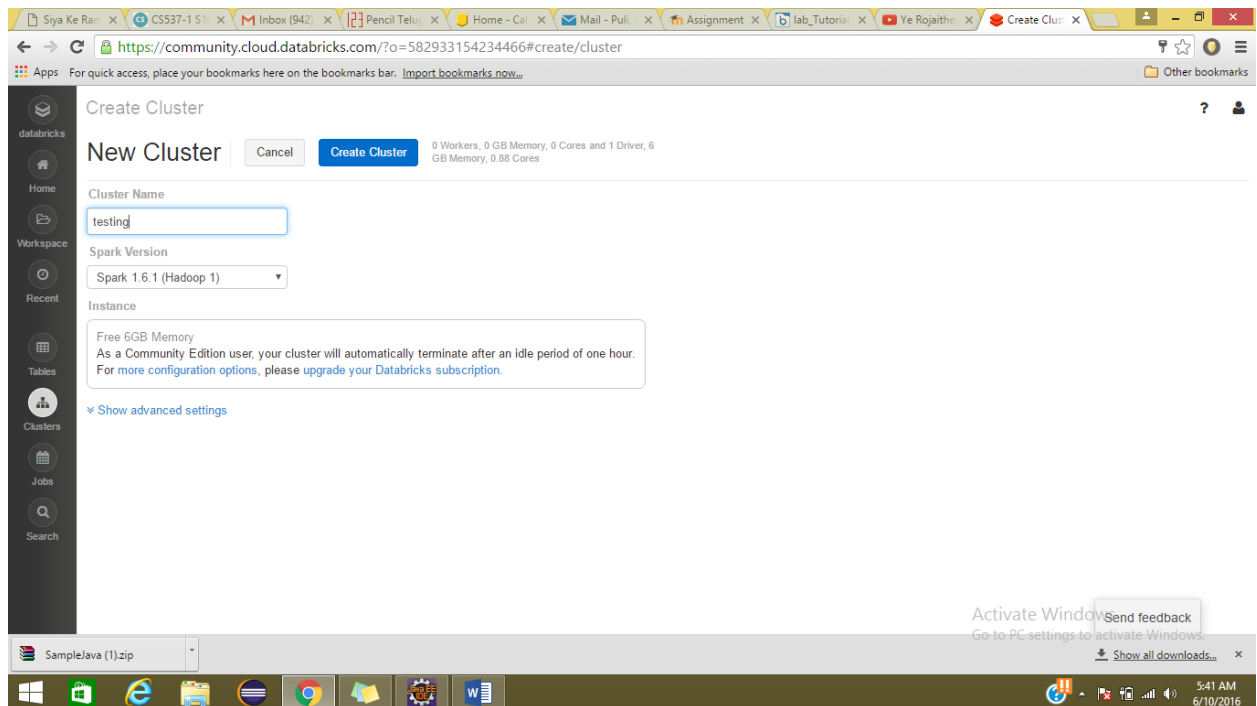- SQL commands to perform the analysis.

- Visualization

## Platform Spec

- IBM Bluemix BigInsights
- CPU Speed: ?

- # of CPU cores: ?
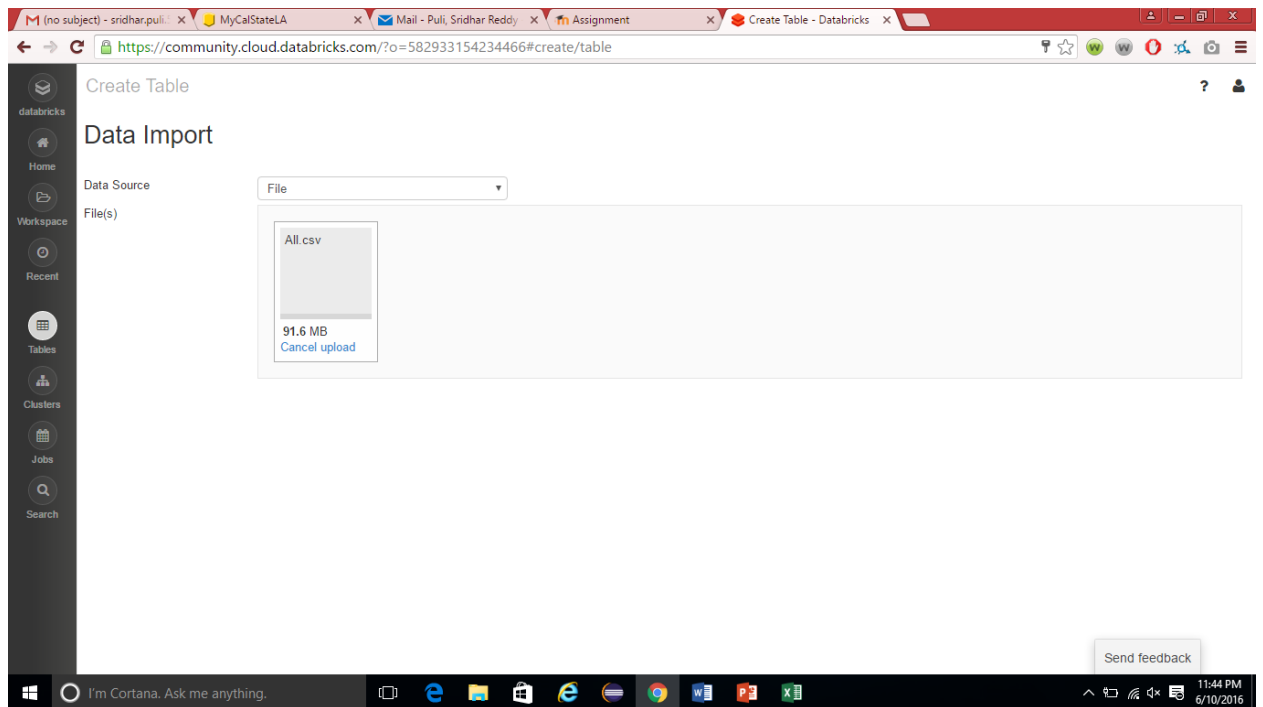- # of nodes: ?
- Total Memory Size: ?

# Step 1: Get data manually using REST API

**Explain what this step is for.** This step is to get data manually….

1. Create Google API keys at https://develop:

2. Sign into your databricks account.

3. Go to Clusters option on the left and click on create cluster.

4. Give the cluster name and click create cluster.



5. Under tables section click on create table and select the file to upload.

# Step 2: Train NLP

**Explain what this step is for.** This step is to …
**Code should be in the following format and indent:**

```
import org.apache.spark.ml.feature.RegexTokenizer
val tokenizer = new RegexTokenizer()
  .setPattern("\\p{L}+").setMinTokenLength(3)
.setGaps(false)
  .setInputCol("text")
  .setOutputCol("words")

val tokenized_df=tokenizer.transform(splits(0))

vi) Use the below code to remove stop words
Run them in separate cells for better understanding

%sh wget
http://ir.dcs.gla.ac.uk/resources/linguistic_utils/stop_words -O
/tmp/stopwords
%fs cp file:/tmp/stopwords dbfs:/tmp/stopwords
val stopwords = sc.textFile("/tmp/stopwords").collect()

import org.apache.spark.ml.feature.StopWordsRemover
// Set params for StopWordsRemover
val remover = new StopWordsRemover()
```

```
      .setStopWords(stopwords) // This parameter is optional
      .setInputCol("words")
      .setOutputCol("filtered")

  // Create new DF with Stopwords removed
  val filtered_df = remover.transform(tokenized_df)
```

1. To show top ten categories

```
sqlContext.sql("Select categories__001,count(*) as count1 from

business_data13 group by categories__001 order by count1

desc").show(10)
```
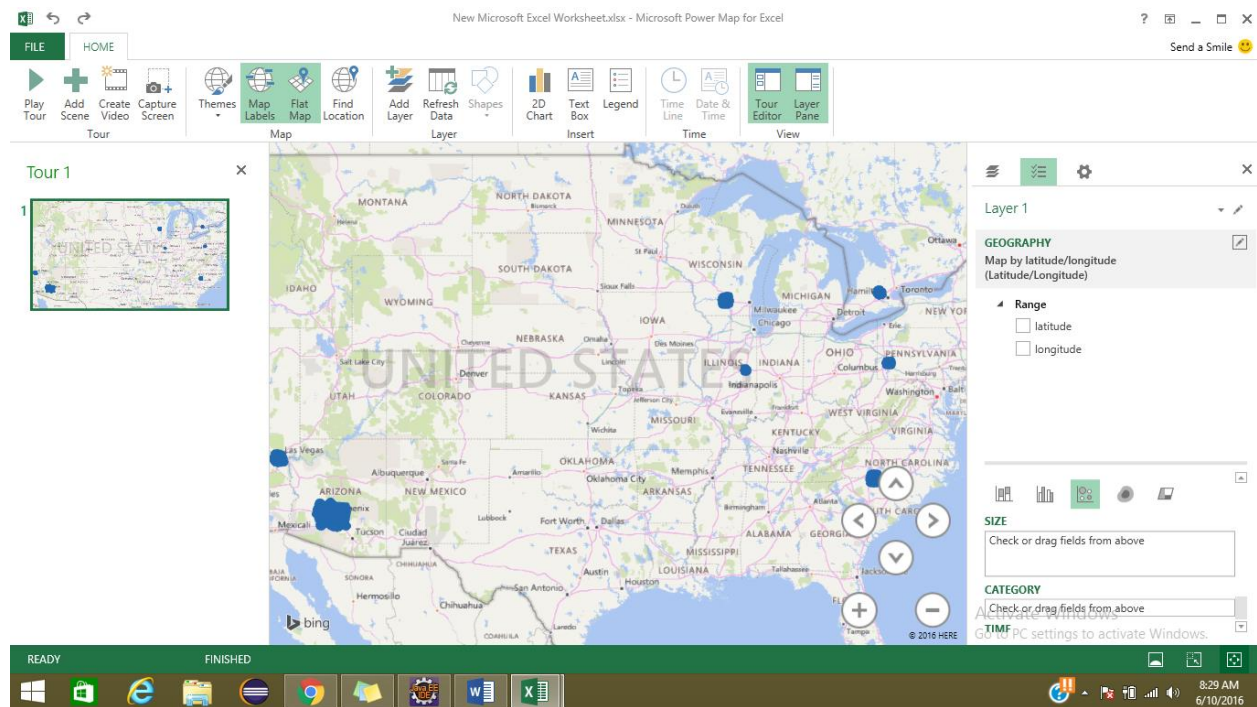


## Step 3: Visualization

**Explain what this step is for.** This step is to…



1. To visualize location type of results on map, convert csv file to excel and click on map button under insert tab.

# References

1. URL of Data Source, http://www.calstatela.edu

2. URL of your Github

3. URL of References

4.