# Lab Tutorial

By Sridhar reddy Puli (spuli@calstatela.edu), Ram Dharan Donda (rdonda@calstatela.edu), Goutham kumar Pola (gpola@calstatela.edu) and Vinay Chennupati (vchennu@calstatela.edu)
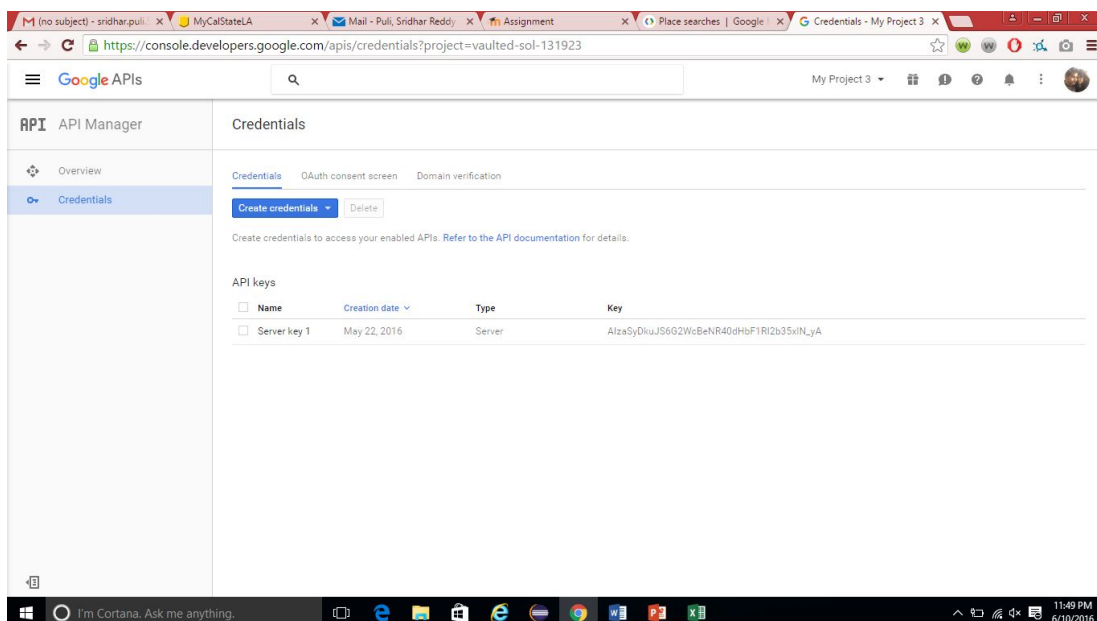
06/10/2016

# Yelp Data Analysis using Spark

## Objectives

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

## Exercise 1: Get data manually using REST API
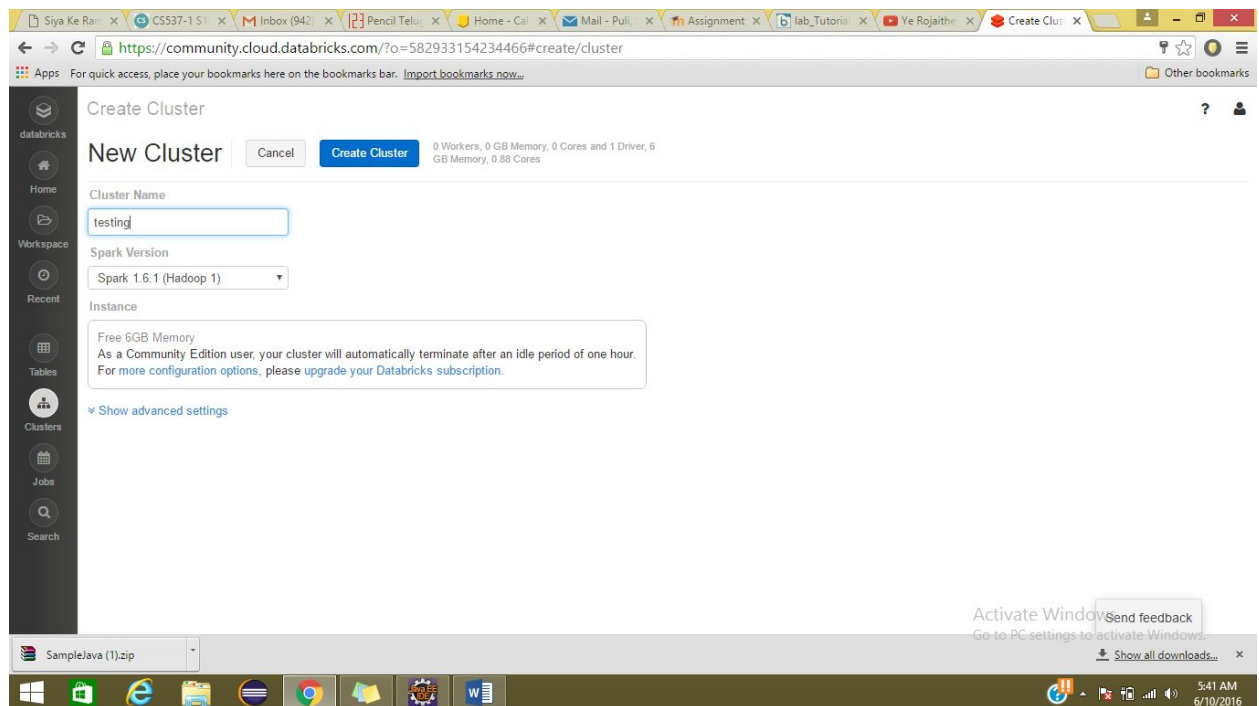
1. Create Google API keys at https://developers.google.com/places/web-service/get-api-key
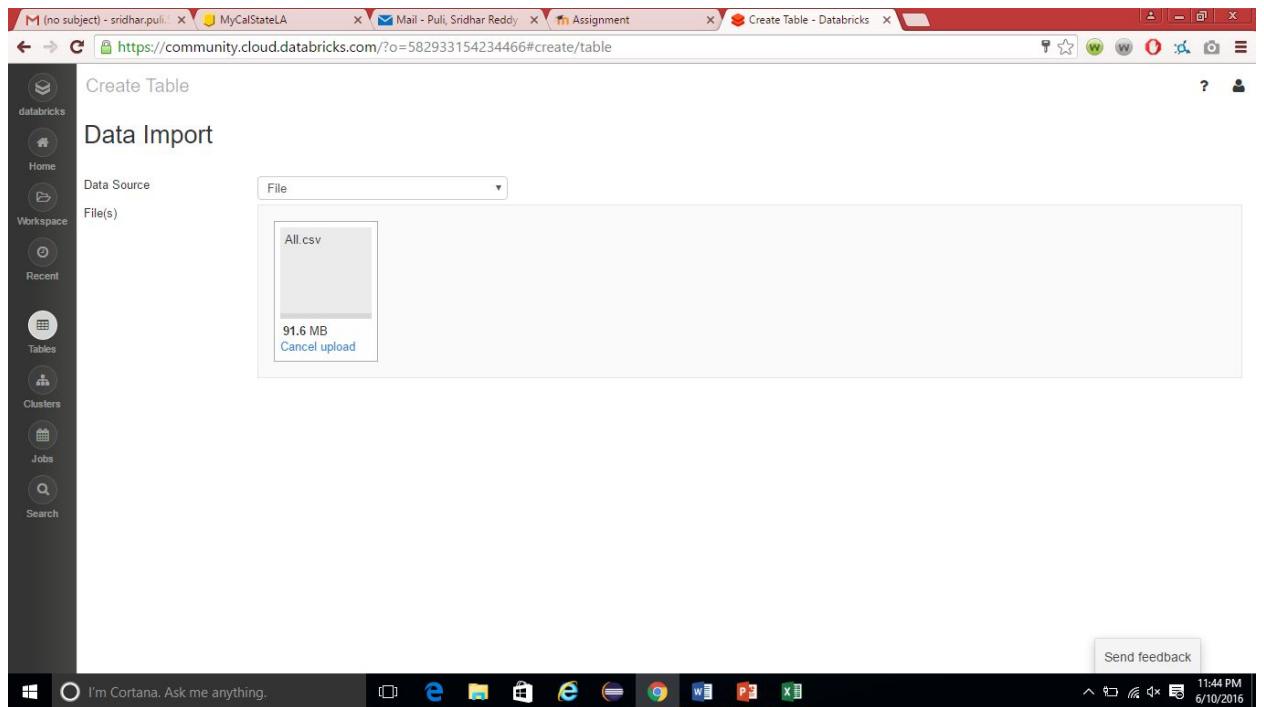
2. Run the given Google places java code by replacing key values with your keys in the code.

3. Also run the given java code for Yelp data.

4. Result data will be stored in a file in your workspace.

5. Convert the Json data into csv using online tools and merge data into one file.

# Exercise 2: Create Spark cluster and load data

1. Sign into your databricks account.

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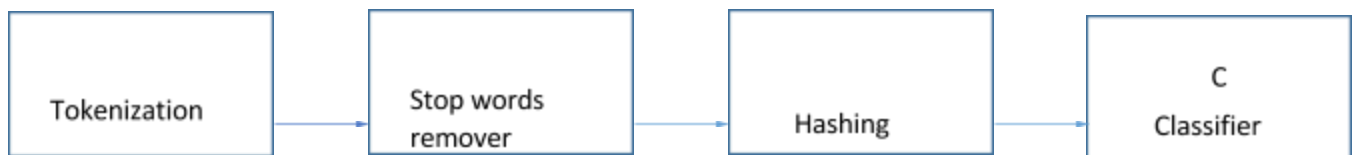
3. Give the cluster name and click create cluster.

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



# Exercise 3: Train NLP

**Pipeline stages:**

Spark machine learning API makes it easier to combine multiple steps into a single workflow called as pipeline. The figure below shows the pipeline stages and data is transformed at each stage before passing it to the next.



Creating table for Natural Language processing:

i)  Download the data from yelp challenge
ii) As shown above create a table by using "review.json" file and create a **Scala** notebook.
iii) Query the data from the table as shown below.
 val data = sqlContext.sql("SELECT CAST(stars as DOUBLE) as label, text from TABLE NAME")
   val splits = data.randomSplit(Array(0.80, 0.20), seed = 10)
iv) We use the split data as
splits[0]→ Training data
   splits[1]→ Test data

**Don't copy this in the notebook, this is just for understanding.**

v)Use the below command to tokenize the data

```
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)
```

vii)Use the below code for hashing after removing the stop words

```
import org.apache.spark.ml.feature.{HashingTF, Tokenizer}

val hashingTF = new HashingTF()
  .setNumFeatures(1000)
  .setInputCol("filtered")
  .setOutputCol("features")
```

viii) Create a Naïve Bayes  model by using below code

```
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.PipelineStage
import org.apache.spark.ml.classification.NaiveBayes
```

```
val nb = new NaiveBayes()
nb.setModelType("multinomial")
```

ix) Set the piple line stages

```
val pipeline = new Pipeline().setStages(Array(tokenizer,remover,hashingTF,nb))
```

x)Fit the training data into pipeline
```
 val lrModel = pipeline.fit(splits(0))
```

xi)Test the model , with the test data
```
val dd = lrModel.transform(splits(1))
```

xii) To see the results of the test data
```
val res=dd.select("label","prediction")
display(res)
```

# Exercise 4: SQL commands to perform the analysis

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)
   ```

2. For showing total number of 4 and 5 star rating businesses for every area in los angeles.

   result=sqlContext.sql("SELECT count(stars) as total,city from business_datafinal where city

   IN('alhambra','Pasadena','Long beach','Santa monica','Beverly hills','burbank','West

   hollywood','arcadia','El monte','Monterey park','San gabriel','downey','baldwin

   park','Montebello','Los angeles') and stars IN(5,4) group by city order by total desc")
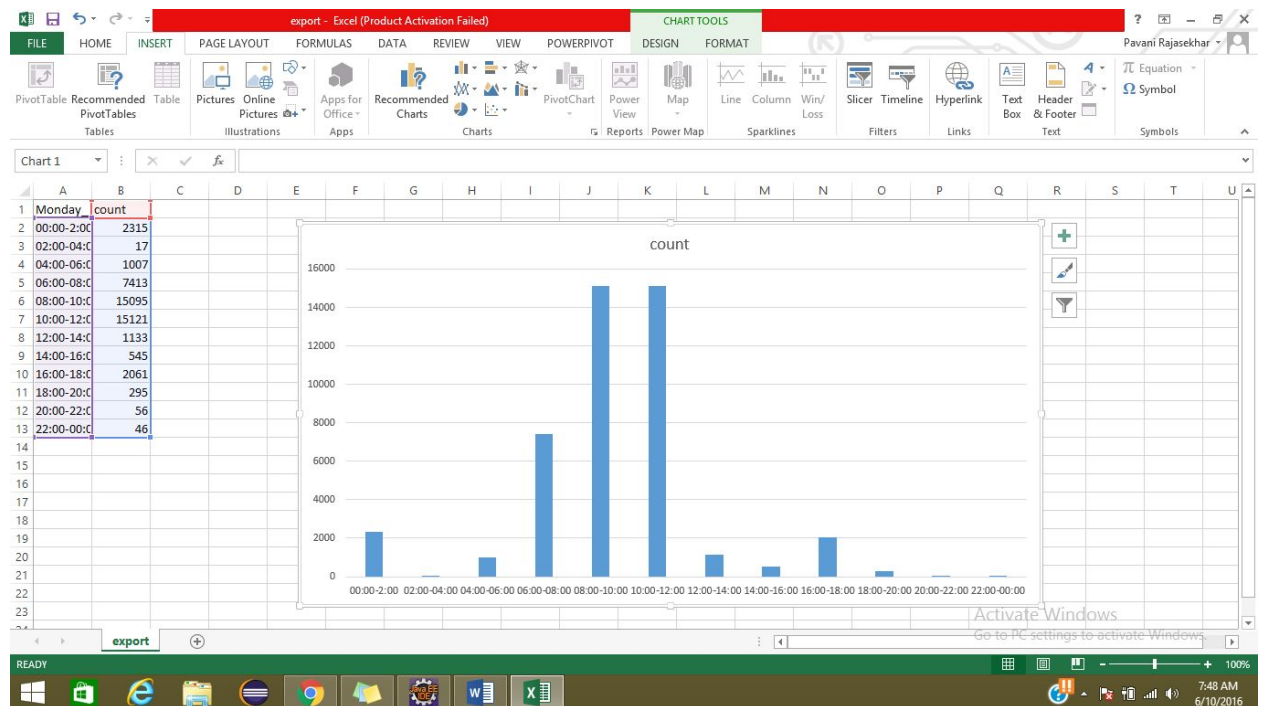
3. Similarly copy and paste the given SQL commands to find different kinds of analysis.

# Exercise 5: Visualization

1. Select the chart type below the result in databricks to instantly show the visualization.

2. Or, click the download button below the result.

3. Open the downloaded file, under insert tab select the chart best suited.

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