# Using Spark and Watson on IBM Bluemix for Sentiment Analytics

## Miracle Summer of Code Virtual Labs Series

## Mounika Chirukuri

Lead Researcher – Big Data/Analytics
Miracle Software Systems, Inc.

April 26th, 2016

# Miracle's SoC Series : Using Spark and Watson on IBM Bluemix for Sentiment Analytics

## Overview

In this lab we will be using Twitter Apps, Spark and Watson from within Bluemix to analyze a twitter data stream. Through that stream we will also plot graphs by using GraphX within the Spark service.

## Prerequisites

You will need the following to complete this lab successfully,

- Browser for accessing IBM Bluemix and Twitter
- Active email ID for registering with Bluemix
- Registered and active accounts with Bluemix

## Technology Involved

The following technologies will be covered in this lab,

- IBM Bluemix(PaaS)
- Twitter Apps
- Apache Spark(Service in Bluemix)
- Scala and Python Commands

## Lab Steps

So, let us get started with the lab!
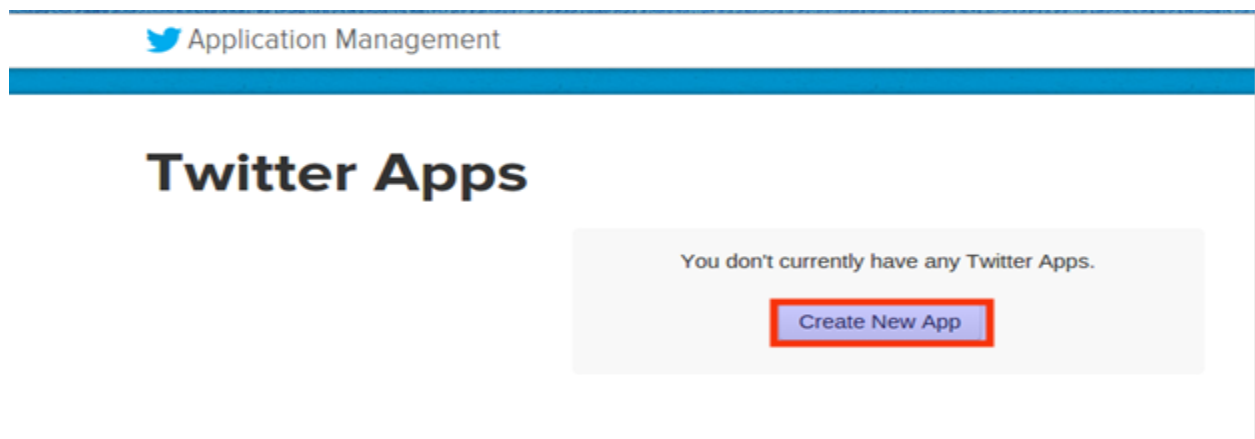
## #1 | Create an application in Twitter

The first step will be to make sure that we have access to Twitter.

Login to Twitter Apps at **https://apps.twitter.com/** (or) Register today at **https://twitter.com/signup**

**Note** : Twitter account should be registered with mobile in order to create applications



After you login, Click on **"Create New App"**.

Scroll down and check "Yes, I agree" box and click on **"Create Your Twitter Application"**.

Consumer Secret and Consumer Key will be generated. Go to **keys and access tokens** tab.



Scroll down and click on **"Create my access token"**.

**Your Access Token**

*You haven't authorized this application for your own account yet.*

*By creating your access token here, you will have everything you need to make API calls right away. The access token generated will be assigned your application's current permission level.*

Token Actions

Create my access token

Token keys will be generated as follows,

**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 | ch_mounika |
| Owner ID | 331125272 |

Token Actions

Regenerate My Access Token and Token Secret    Revoke Token Access

Please note down the Consumer keys and Access tokens for later purpose.

## #2 | Access Bluemix

Make sure that we have access to the IBM Bluemix Console with either the free trial option (or) the paid subscription option.

Login to Bluemix at **http://bluemix.net**  (or) Register today at **https://console.ng.bluemix.net/registration/**
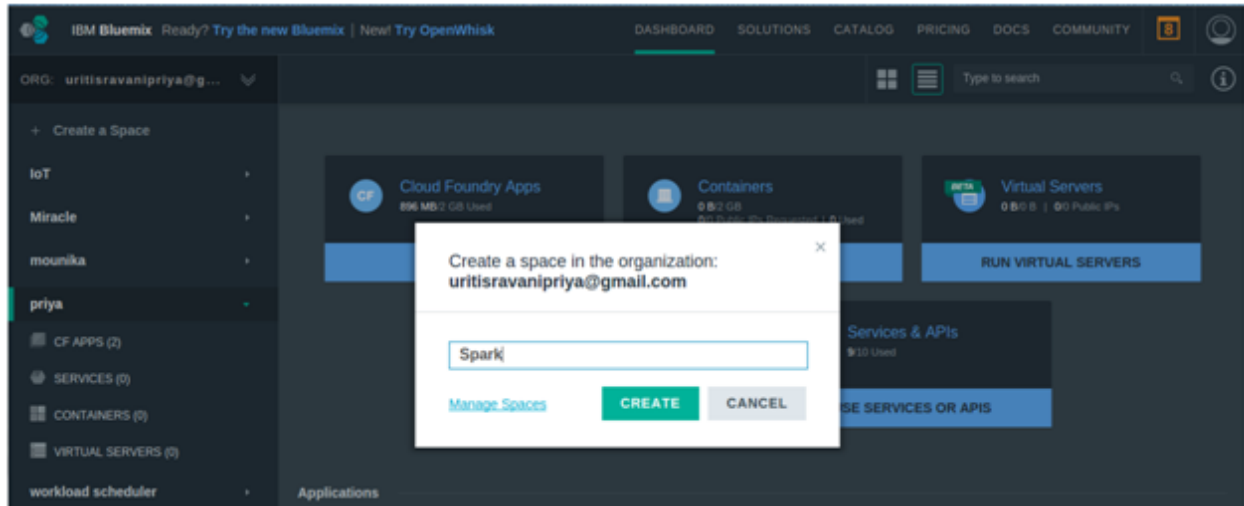
After you login, you can see the dashboard where you can take a look at your applications and services. You can also go to the profile icon at the top and change which Cloud Availability Region you are working in. Make sure you select **US South** region, as **Apache Spark Service** is only available in that region.
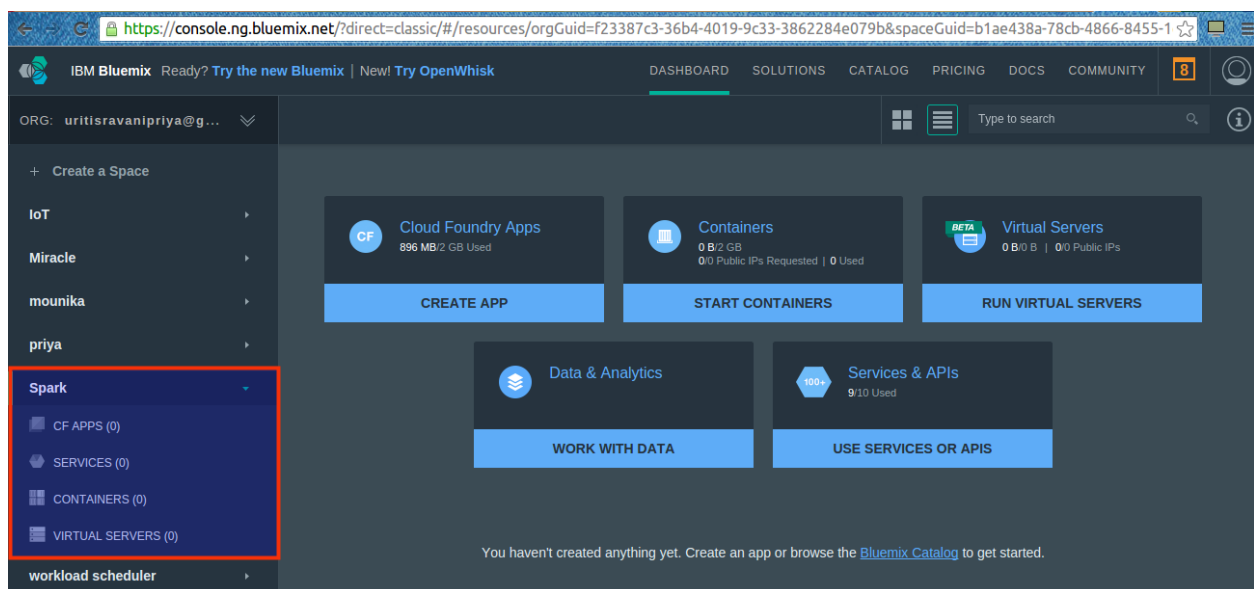
## #3 | Create Space

Click on **"Create a Space"** on left side of the Dashboard to create a space in Bluemix Cloud. Give any name and click on "create".
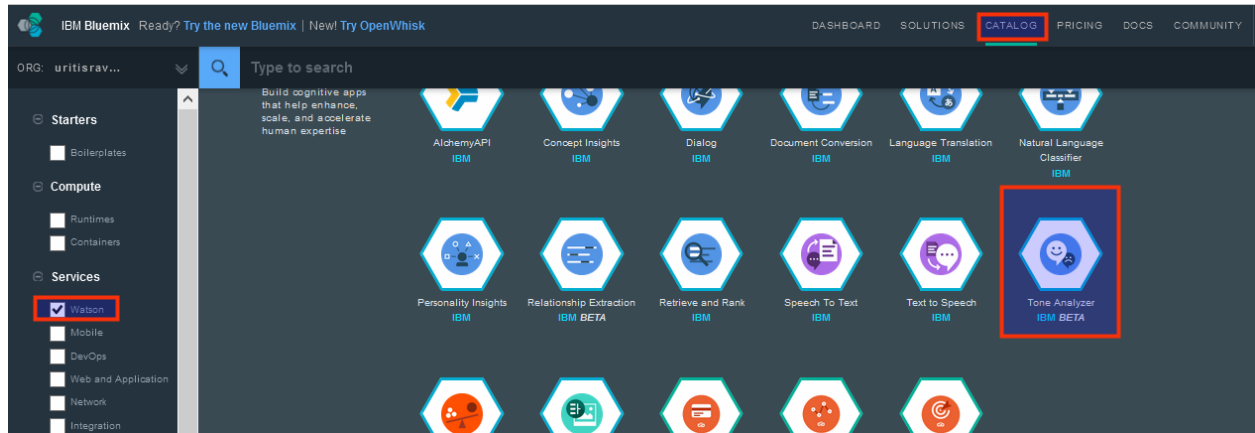


Once the space is created, you can find it at the left side menu of the Dashboard.
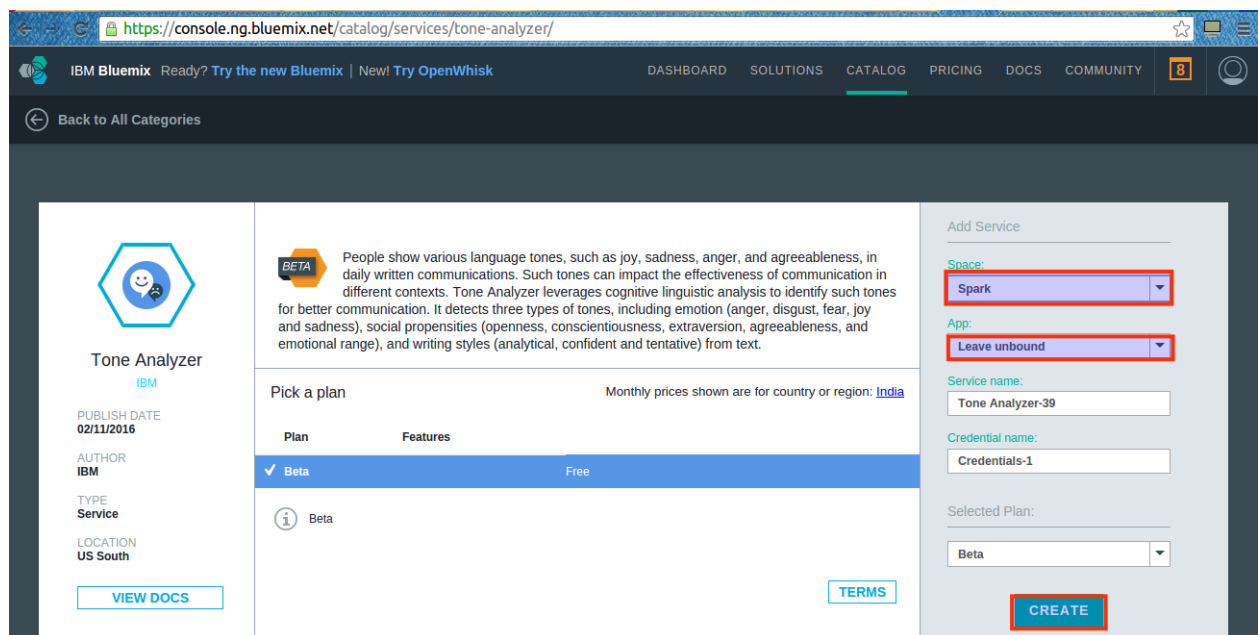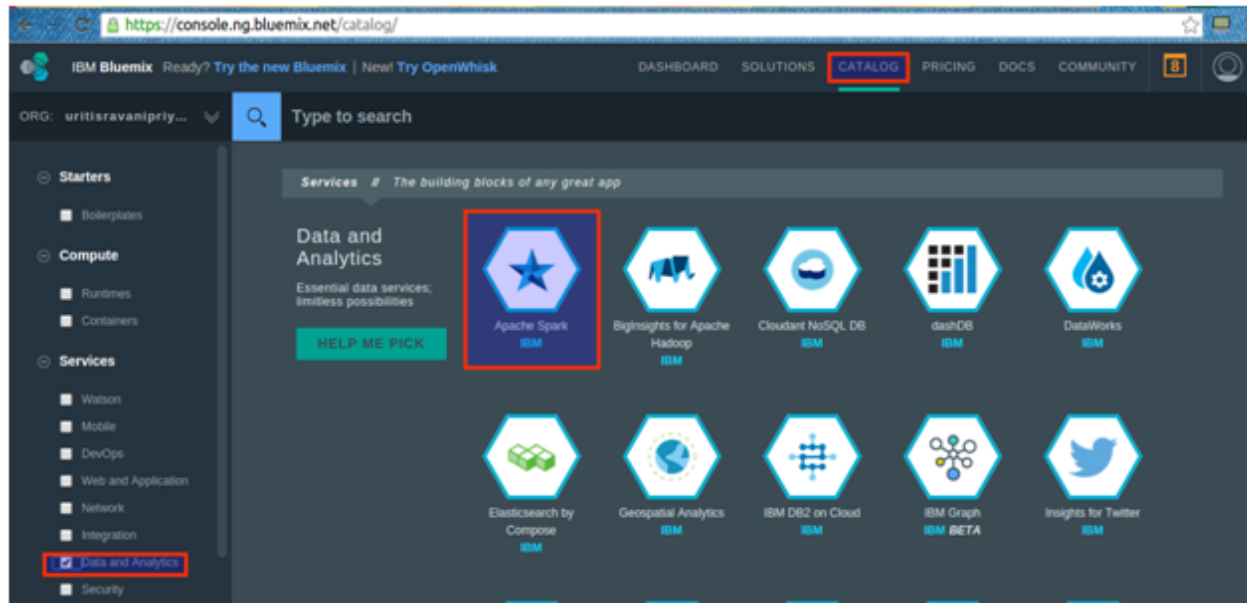
## #4 | Create Watson Tone Analyzer Service

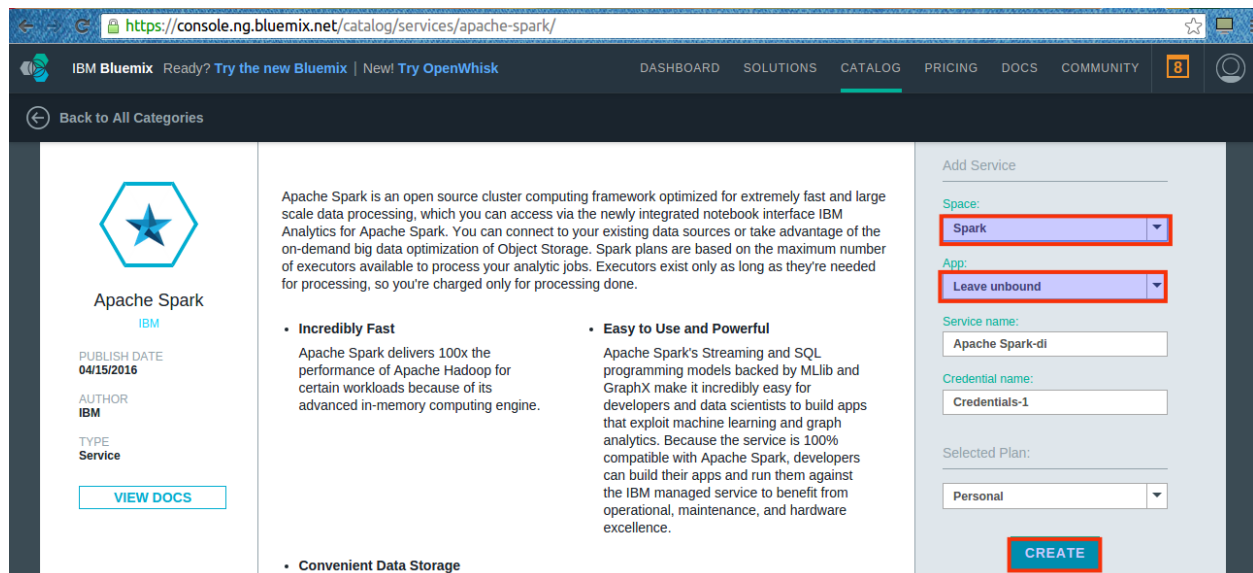Go to "Catalog" and select "Watson" in the left side menu. Click on **"Tone Analyzer"** service.



Choose your created space in the "Space" field. Choose "Leave Unbound" in the "App" field and click on "Create".

Click on "Service Credentials".



Copy the credentials for later purpose.



## #5 | Create Apache Spark Service

Go to "Catalog" and select "Data and Analytics" in the left side menu. Click on **"Apache Spark"** service.
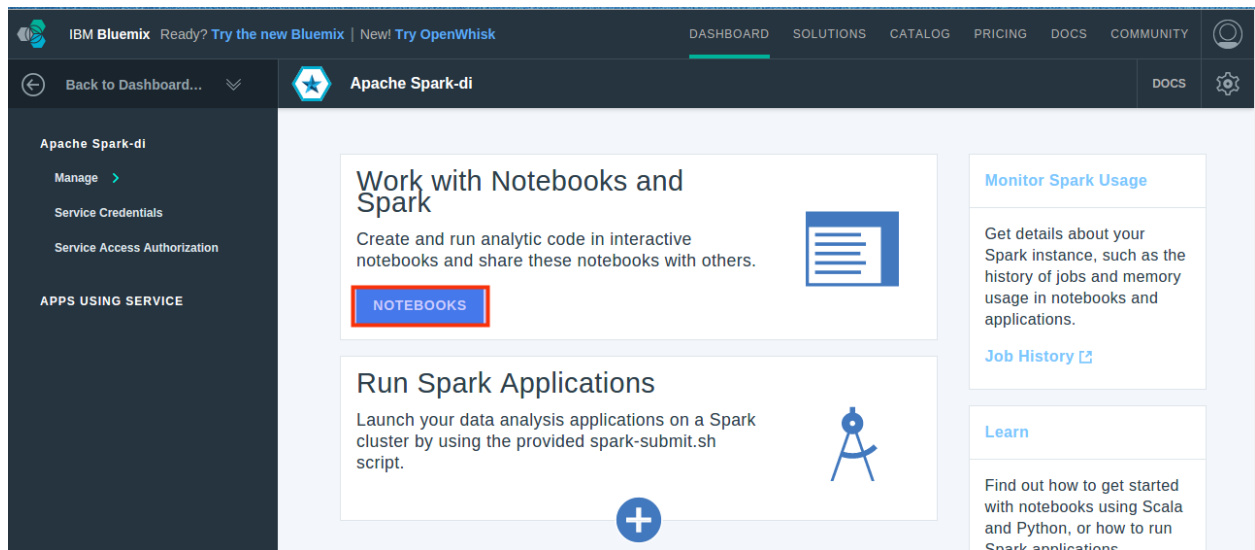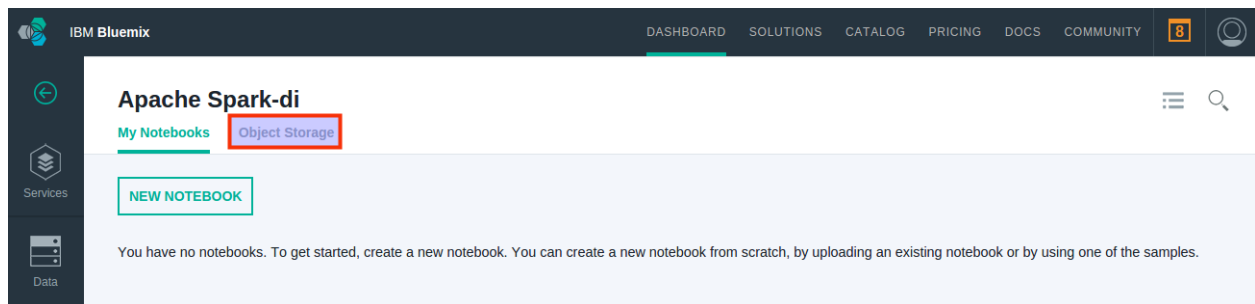
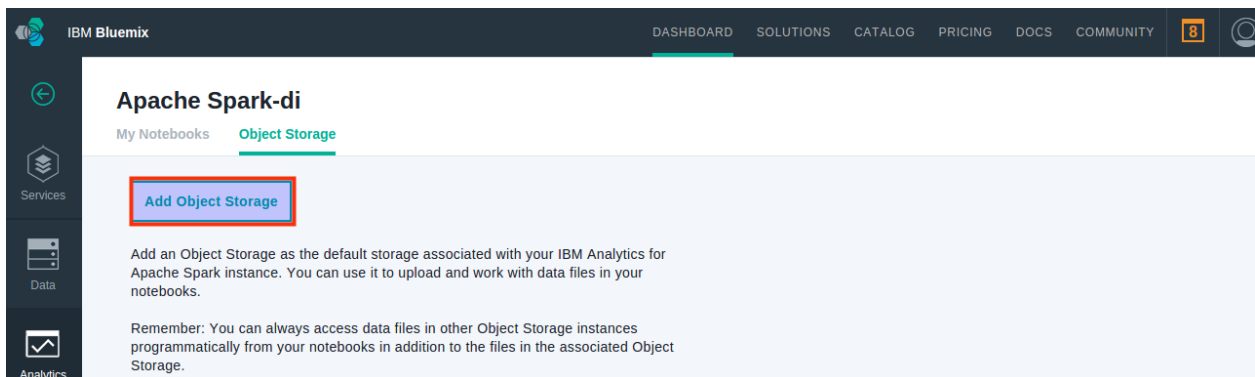Choose your created space in the "Space" field. Choose "Leave Unbound" in the "App" field and click on "Create".



Click on "Notebooks".

Click on "Object Storage" tab.



Click on "Add Object Storage".



Select your space for "Space" field. Click on "Create".

A new Object Storage will be created.

## #6 | Create Scala Notebook

Click on **"My Notebooks"** to create a new notebook.



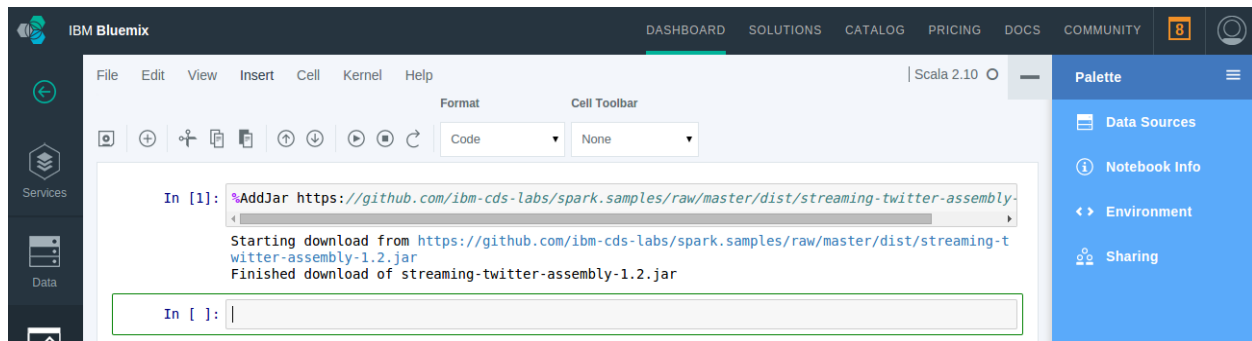Give any name for "Name" field. Select "Scala" radio button for "Language" field. Click on "Create Notebook".

A Scala shell will be opened. In the first cell, enter the following command to install the application jar.

**%AddJar https://github.com/ibm-cds-labs/spark.samples/raw/master/dist/streaming-twitter-assembly-1.2.jar -f**



Click on run button to execute the cell.

Once the jar is downloaded, configure the credentials needed to connect to Twitter and Watson Tone Analyzer service.
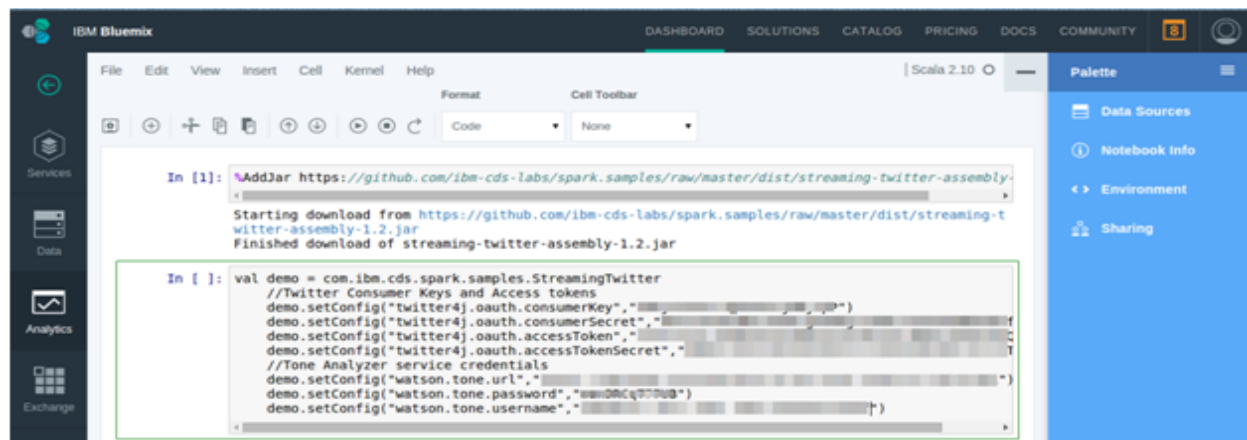
```scala
val demo = com.ibm.cds.spark.samples.StreamingTwitter

//Twitter Consumer Keys and Access tokens
demo.setConfig("twitter4j.oauth.consumerKey","<Consumer-Key>")
demo.setConfig("twitter4j.oauth.consumerSecret","<Consumer-Secret>")
demo.setConfig("twitter4j.oauth.accessToken","<access-Token>")
demo.setConfig("twitter4j.oauth.accessTokenSecret","<access-Token-Secret>")

//Tone Analyzer service credentials
demo.setConfig("watson.tone.url","<URL>")
demo.setConfig("watson.tone.password","<password>")
demo.setConfig("watson.tone.username","<user-name>")
```

To start streaming tweets from Twitter, run the following command,

```
import org.apache.spark.streaming._
demo.startTwitterStreaming(sc, Seconds(20))
```
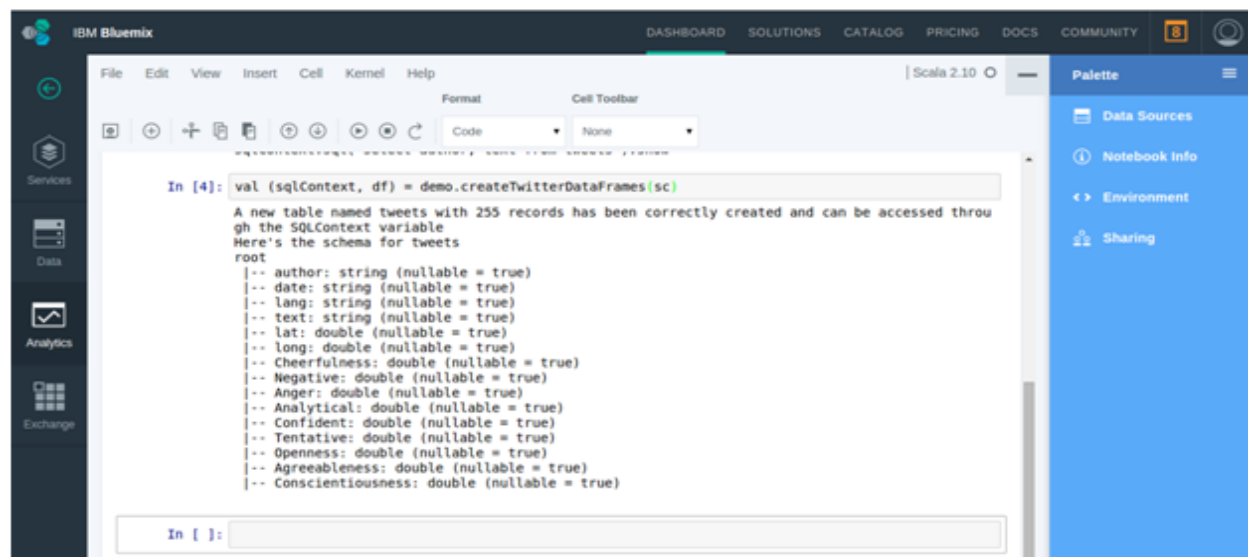


Once the stream stops, create a DataFrame using Spark SQL.

```
val(sqlContext,df) = demo.createTwitterDataFrames(sc)
```



Run the following command to display a sample of data.

```
val fullSet = sqlContext.sql("select * from tweets limit 100000")
fullSet.show
```

To save the dataset into a parquet file on Object Storage,

```
fullSet.repartition(1).saveAsParquetFile("swift://notebooks.spark/Tweets.parquet")
```

## #7 | Analyzing data using Python Notebook
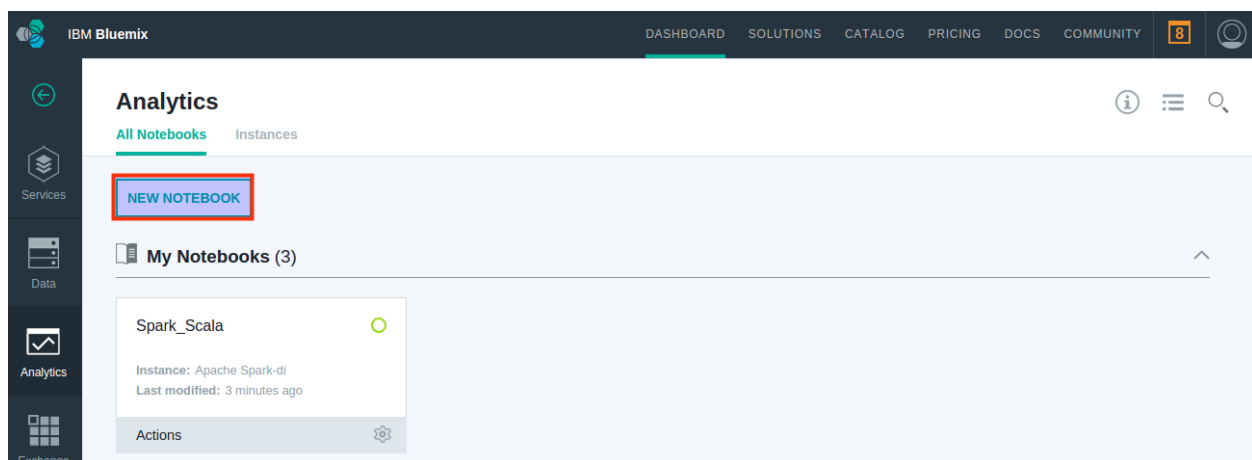
Go to Notebook page by clicking on back button on the upper left corner.



Click on **"New Notebook".**



Give any name for **"Name"** field. Select the **"Python"** radio button and click on **"Create Notebook"**.

A **Python shell will appear**. Load the data from Object Storage and create a DataFrame with the entire dataset by using the following command,

```python
# Import SQLContext and data types
from pyspark.sql import SQLContext
from pyspark.sql.types import *
# sc is an existing SparkContext.
sqlContext = SQLContext(sc)
parquetFile = sqlContext.read.parquet("swift://notebooks.spark/Tweets.parquet")
print parquetFile
parquetFile.registerTempTable("tweets");
sqlContext.cacheTable("tweets")
tweets = sqlContext.sql("SELECT * FROM tweets")
print tweets.count()
tweets.cache()
```

Now we can start analyzing the data by create an array that will hold the count for each sentiment. For each sentiment, run a SQL query that counts the number of tweets for which the sentiment score is greater than 60%.

```python
sentimentDistribution=[0] * 9
for i, sentiment in enumerate(tweets.columns[-9:]):
    sentimentDistribution[i]=sqlContext.sql("SELECT count(*) as sentCount FROM tweets where " + sentiment + " > 60")\
        .collect()[0].sentCount
```

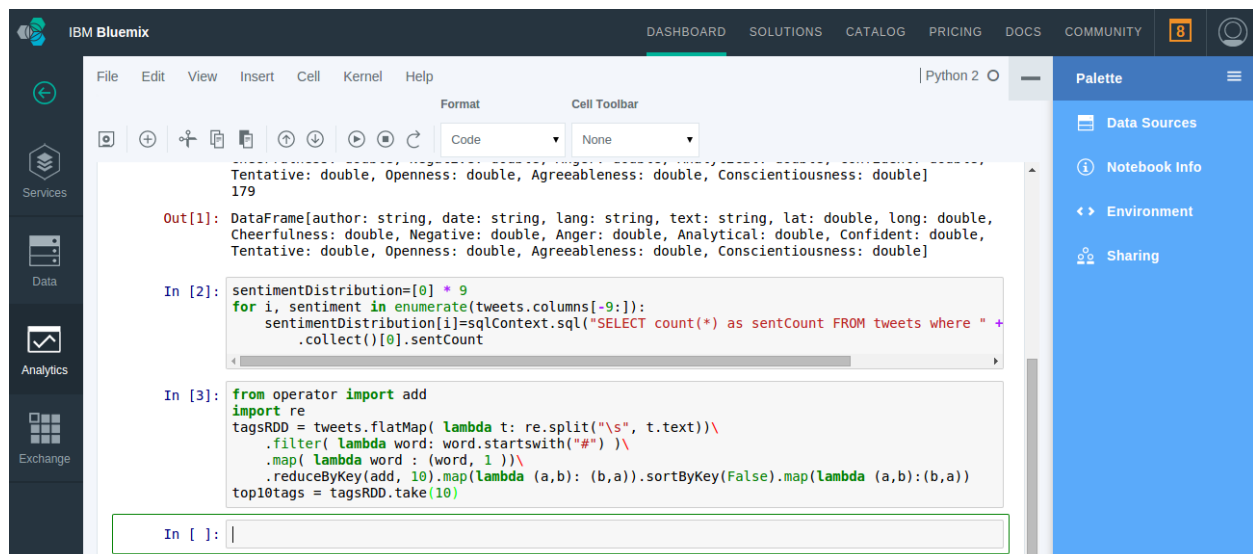Compute the top 10 hashtags contained in the tweets by using the following command,

```python
from operator import add
import re
tagsRDD = tweets.flatMap( lambda t: re.split("\s", t.text))\
    .filter( lambda word: word.startswith("#") )\
    .map( lambda word : (word, 1 ))\
    .reduceByKey(add, 10).map(lambda (a,b):
(b,a)).sortByKey(False).map(lambda (a,b):(b,a))
top10tags = tagsRDD.take(10)
```



Plot the RDD as a pie chart by using the following code,

```python
%matplotlib inline
import matplotlib
import matplotlib.pyplot as plt
params = plt.gcf()
plSize = params.get_size_inches()
params.set_size_inches( (plSize[0]*2, plSize[1]*2) )
labels = [i[0] for i in top10tags]
sizes = [int(i[1]) for i in top10tags]
colors = ['yellowgreen', 'gold', 'lightskyblue', 'lightcoral',
```

```
"beige", "paleturquoise", "pink", "lightyellow", "coral"]
plt.pie(sizes, labels=labels, colors=colors,autopct='%1.1f%%',
shadow=True, startangle=90)
plt.axis('equal')
plt.show()
```



The top 10 Hashtags sentiment percentage will appear as a pie chart.