# PRINCIPLES OF BIG DATA MANAGEMENT

# PROJECT – I

# SPRING 2017

Submitted by: TEAM 7

Sri Sai Anusha Gandu (sgr43)

Sushma Mitta (smgp6)

Sadanand Kallakuri (sk789)

Abhilash Reddy Gaddam (aggg6)

# CONTENTS

# 1. INTRODUCTION

## 1.1 About Twitter

Twitter is an online news and social networking service where users post and interact with messages, "tweets," restricted to 140 characters. Registered users can post tweets, but those who are unregistered can only read them. Twitter Inc. is based in San Francisco, California, United States, and has more than 25 offices around the world. Twitter was created in March 2006 by Jack Dorsey, Noah Glass, Biz Stone, and Evan Williams and launched in July, whereby the service rapidly gained worldwide popularity. As of 2016, Twitter had more than 319 million monthly active users.

## 1.2 About the Project

Here, we have collected the tweets using twitter API through tweepy using the keywords Python, JavaScript and Ruby in JSON (JavaScript Object Notation) format. The tweets then collected have been analyzed and different SQL queries are written to obtain the result.

# 2. REQUIREMENTS

## 2.1 Languages

1. Python
2. Scala
3. SQL
4. Java
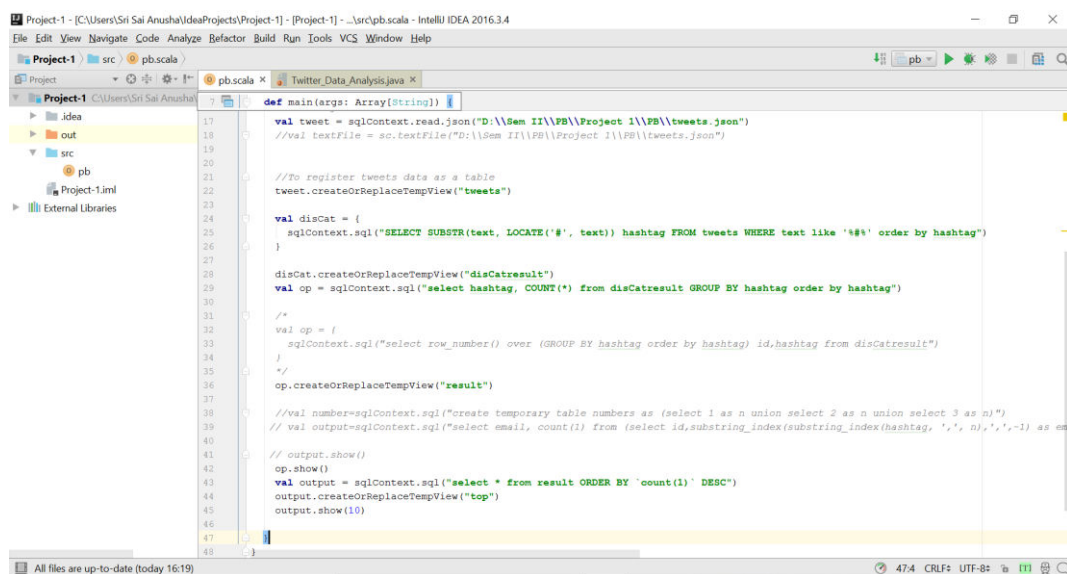
## 2.2 Software

1. IntelliJ IDEA 3.4 (IDE)
2. Python 3.6
3. JDK 1.8
4. Scala 2.12.1
5. Spark 2.1
6. Virtual Box (Cloudera)

## 3. TASKS

1. To collect the tweets in JSON (JavaScript Object Notation) format.
2. Find the list of top ten hash-tags.
3. To create HDFS directories for each of top ten hash-tags along with two other directories named Others and None.
4. Segregate the tweets into the respective HDFS directories.
5. To implement a function that can count the frequency of keywords in the directories.

## 4. TOP TEN HASHTAGS



We have collected the data as tables and from the count thus obtained is sorted in descending manner. The output obtained is shown below.

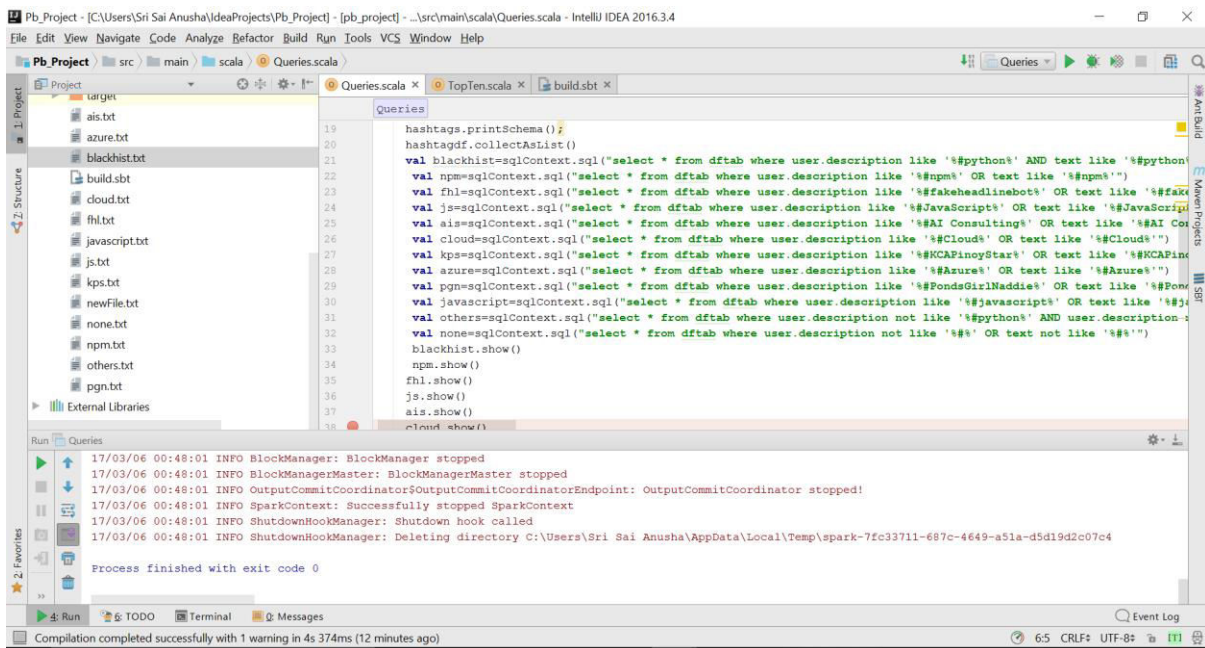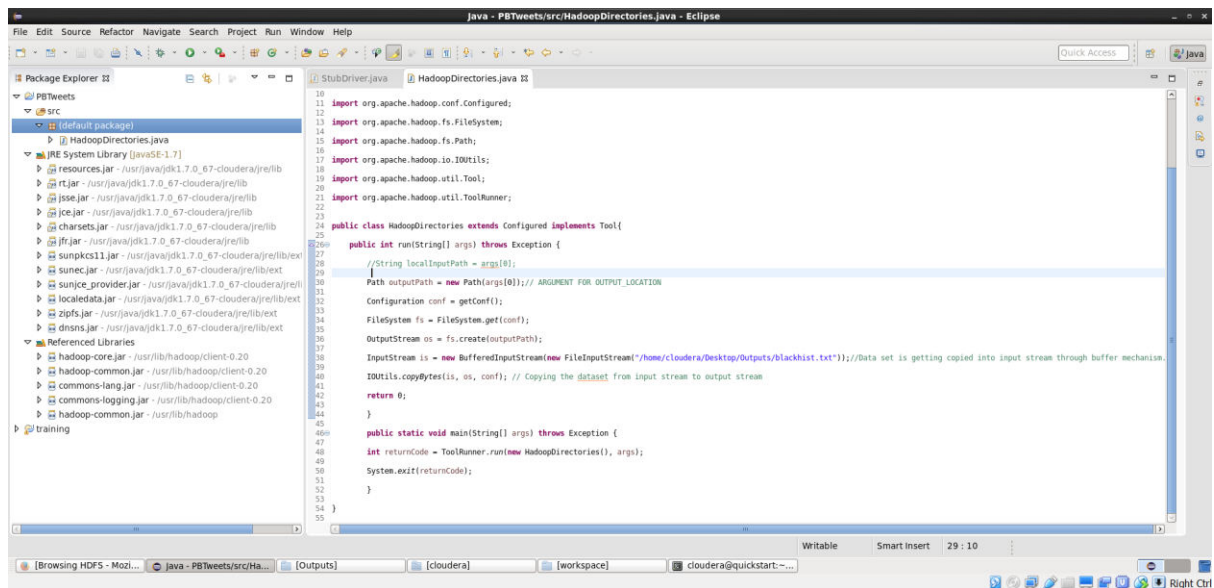# 5. HADOOP FILE DIRECTORIES

The following figure shows the queries written to obtain the top ten directories.

## Sample JAVA Code to create the directories



## Workspace

## Hadoop File Directories



## Browsing Each Directory Including Others and None

**Screenshot 1:**

Cloudera [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

Sun Mar 5, 23:04 cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/blackhist

Search

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/blackhist | Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 44.83 KB | Sun Mar 05 22:52:18 -0800 2017 | 1 | 128 MB | blackhist |

Hadoop, 2016.

**Screenshot 2:**

Cloudera [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

Sun Mar 5, 23:04 cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/cloud

Search

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/cloud | Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 802.48 KB | Sun Mar 05 22:32:40 -0800 2017 | 1 | 128 MB | cloud |

Hadoop, 2016.

**Screenshot 3:**

Cloudera [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

Sun Mar 5, 23:05 cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/fhl

Search

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/fhl | Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 361.75 KB | Sun Mar 05 22:33:49 -0800 2017 | 1 | 128 MB | fhl |

Hadoop, 2016.

**Screenshot 4:**

Cloudera [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

Sun Mar 5, 23:06 cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/javascript

Search

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/javascript | Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 9.15 MB | Sun Mar 05 22:35:00 -0800 2017 | 1 | 128 MB | javascript |

Hadoop, 2016.

Cloudera [Running] - Oracle VM VirtualBox

File  Machine  View  Input  Devices  Help

Applications  Places  System  Sun Mar 5, 23:06  cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/js

Search

Cloudera  Hue  Hadoop  HBase  Impala  Spark  Solr  Oozie  Cloudera Manager  Getting Started

Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/js — Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 4.91 MB | Sun Mar 05 22:35:39 -0800 2017 | 1 | 128 MB | js |

Hadoop, 2016.

---

Cloudera [Running] - Oracle VM VirtualBox

File  Machine  View  Input  Devices  Help

Applications  Places  System  Sun Mar 5, 23:06  cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/kps

Search

Cloudera  Hue  Hadoop  HBase  Impala  Spark  Solr  Oozie  Cloudera Manager  Getting Started

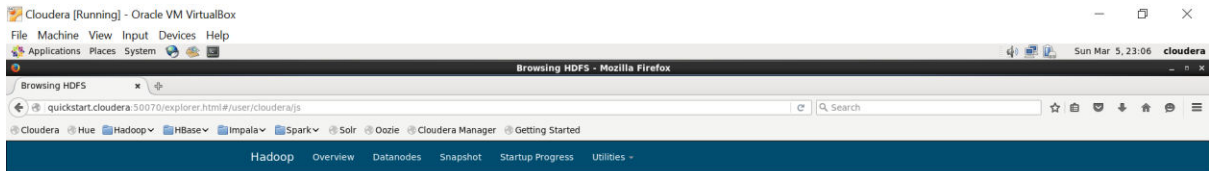Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/kps — Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 671.42 KB | Sun Mar 05 22:36:26 -0800 2017 | 1 | 128 MB | kps |

Hadoop, 2016.

---

Cloudera [Running] - Oracle VM VirtualBox

File  Machine  View  Input  Devices  Help

Applications  Places  System  Sun Mar 5, 23:07  cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/none

Search

Cloudera  Hue  Hadoop  HBase  Impala  Spark  Solr  Oozie  Cloudera Manager  Getting Started

Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/none — Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 118.43 MB | Sun Mar 05 22:37:23 -0800 2017 | 1 | 128 MB | none |

Hadoop, 2016.

---

Cloudera [Running] - Oracle VM VirtualBox

File  Machine  View  Input  Devices  Help

Applications  Places  System  Sun Mar 5, 23:07  cloudera

Browsing HDFS - Mozilla Firefox

Browsing HDFS

quickstart.cloudera:50070/explorer.html#/user/cloudera/npm

Search

Cloudera  Hue  Hadoop  HBase  Impala  Spark  Solr  Oozie  Cloudera Manager  Getting Started

Hadoop  Overview  Datanodes  Snapshot  Startup Progress  Utilities

## Browse Directory

/user/cloudera/npm — Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | cloudera | cloudera | 844.62 KB | Sun Mar 05 22:37:52 -0800 2017 | 1 | 128 MB | npm |

Hadoop, 2016.