

Assignment 11.3:

Problem Statement:

Create a flume agent that streams data from Twitter and stores in the HDFS.

Solution/Steps:

Streaming Twitter Data

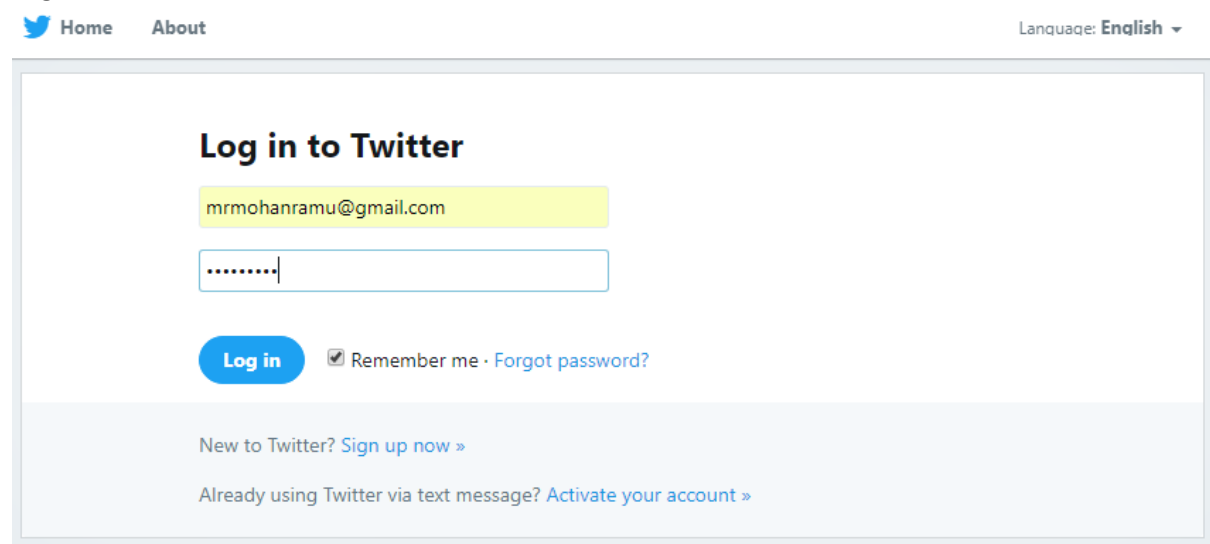
To stream data to our database from twitter we should have the following pre-requisites.

- Twitter account
- Hadoop cluster

If both prerequisites are available we can move to our further step.

Step 1:

Login to the twitter account

A screenshot of the Twitter login interface. At the top, there are links for 'Home' and 'About' next to the Twitter bird icon, and a language selector set to 'English'. The main heading is 'Log in to Twitter'. Below it, there is a text input field containing the email 'mrmohanramu@gmail.com'. Underneath the email field is a password input field with masked characters '.....'. To the left of the password field is a blue 'Log in' button. To the right of the password field is a checkbox labeled 'Remember me' followed by a link 'Forgot password?'. At the bottom of the login area, there are two links: 'New to Twitter? Sign up now »' and 'Already using Twitter via text message? Activate your account »'.

Step 2:

Go to the following link and click the 'create new app' button.

<https://apps.twitter.com/app>

Twitter Apps

Create New App



mohanAcadgildApp

This app will help me to do analysis in flume

Tweet

Step 3:

Enter the necessary details.

Create an application

Application Details

Name *

mohanAcadgildApp

Your application name. This is used to attribute the source of a tweet and in user-facing authorization screens. 32 characters max.

Description *

This app will help me to do analysis in flume

Your application description, which will be shown in user-facing authorization screens. Between 10 and 200 characters max.

Website *

http://www.yahoo.com

Your application's publicly accessible home page, where users can go to download, make use of, or find out more information about your application. This fully-qualified URL is used in the source attribution for tweets created by your application and will be shown in user-facing authorization screens.
(If you don't have a URL yet, just put a placeholder here but remember to change it later.)

Callback URL

Where should we return after successfully authenticating? OAuth 1.0a applications should explicitly specify their oauth_callback URL on the request token step, regardless of the value given here. To restrict your application from using callbacks, leave this field blank.

Step 4:

Accept the developer agreement and select the 'create your Twitter application' button.

Developer Agreement

☒ Yes, I have read and agree to the [Twitter Developer Agreement](#).

Create your Twitter application

Step 5:

Select the 'Keys and Access Token' tab.

 Application Management



mohanAcadgildApp

Test OAuth

Details

Settings

Keys and Access Tokens

Permissions



This app will help me to do analysis in flume

<http://www.yahoo.com>

Organization

Information about the organization or company associated with your application. This information is optional.

Organization None

Organization website None

Step 6:

Copy the consumer key and the consumer secret code.

 Application Management



mohanAcadgildApp

Test OAuth

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) 09ieFYOZ3gltT80CcxhkAmcHN

Consumer Secret (API Secret) uy5NKjsyP19d2gujFitZIdTit25MT59NKTizkwBIFm9vUNPoJ5

Access Level Read and write (modify app permissions)

Owner mrmohanramu

Owner ID 875303803464974337

Step 7:

Scroll down further and select the 'create my access token' button.

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

Now, you will receive a message stating “that you have successfully generated your application access token”.

Status

Your application access token has been successfully generated. It may take a moment for changes you've made to reflect.

[Refresh](#) if your changes are not yet indicated.

Step 8:

Copy the Access Token and Access token Secret code.

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	875303803464974337- ViWylha876bX1k2IhlqzRLF0KsfkMBb
Access Token Secret	XDkzifEAcodZNxNY2URvygLNkkH8IglO0fECjyMPpgB9Q
Access Level	Read and write
Owner	mrmohanramu
Owner ID	875303803464974337

Follow Step 9 and Step 10 to install Apache flume

Step 9: Download flume tar file from below link and extract it.

wget <http://archive.apache.org/dist/flume/1.8.0/apache-flume-1.8.0-bin.tar.gz>

```
[acadgild@localhost ~]$ wget http://archive.apache.org/dist/flume/1.8.0/apache-flume-1.8.0-bin.tar.gz
--2017-12-05 07:00:57-- http://archive.apache.org/dist/flume/1.8.0/apache-flume-1.8.0-bin.tar.gz
Resolving archive.apache.org... 163.172.17.199
Connecting to archive.apache.org|163.172.17.199|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 58688757 (56M) [application/x-gzip]
Saving to: "apache-flume-1.8.0-bin.tar.gz"

100%[=====>] 58,688,757 293K/s in 7m 15s

2017-12-05 07:08:27 (132 KB/s) - "apache-flume-1.8.0-bin.tar.gz" saved [58688757/58688757]

[acadgild@localhost ~]$ ls
apache-flume-1.8.0-bin.tar.gz      hadoop      Public
Assignment3.2-0.0.1-SNAPSHOT.jar hdfs:      sample_temperature_dataset.csv
Assignment3.3-0.0.1-SNAPSHOT.jar hive        spark
assignment11.txt                 hive-site.xml spark-2.2.0-bin-hadoop2.7
company.java                     mapreduce-0.0.1-SNAPSHOT.jar spark-2.2.0-bin-hadoop2.7.tgz
customer.java                    max-temp.txt Sqoop_try.txt
derby.log                       metastore_db sqoop.zip
Desktop                          Music       television.txt
Documents                       mysql-connector-java-5.1.44.zip Templates
Downloads                       Pictures    Videos
eclipse                         pig         workspace
eclipse-jee-neon-M3-linux-gtk-x86_64.tar.gz pig_1510851976774.log
employee.java                   pig_1511973662738.log
[acadgild@localhost ~]$ tar xzvf apache-flume-1.8.0-bin.tar.gz
apache-flume-1.8.0-bin/lib/flume-ng-confiuration-1.8.0.jar
apache-flume-1.8.0-bin/lib/slf4j-api-1.6.1.jar
apache-flume-1.8.0-bin/lib/slf4j-log4j12-1.6.1.jar
apache-flume-1.8.0-bin/lib/log4j-1.2.17.jar
```

Update the path of extracted flume directory in the .bashrc file as mentioned in the below image.
NOTE: keep the path same as where the extracted file exists.

```
[acadgild@localhost ~]$ cat .bashrc
# .bashrc

# Source global definitions
if [ -f /etc/bashrc ]; then
    . /etc/bashrc
fi

# User specific aliases and functions

export JAVA_HOME=/usr/local/java
export PATH=$PATH:$JAVA_HOME/bin

export HADOOP_HOME=/usr/local/hadoop-2.6.0
export PATH=$PATH:$HADOOP_HOME/bin

export FLUME_HOME=/usr/local/flume

export PIG_INSTALL=/usr/local/pig

export HIVE_HOME=/usr/local/hive
export PATH=$HIVE_HOME/bin:$PATH
export HADOOP_USER_CLASSPATH_FIRST=true

export HBASE_HOME=/usr/local/hbase

export SQOOP_HOME=/usr/local/sqoop

export OOZIE_HOME=/usr/local/oozie-4.1.0

export FLUME_HOME=/home/acadgild/apache-flume-1.8.0-bin

export PATH=$PATH:$FLUME_HOME/bin:$PIG_INSTALL/bin:$HIVE_HOME/bin:$HBASE_HOME/bin:$SQOOP_HOME/bin:$HADOOP_HOME/sbin:$OOZIE_HOME/bin:$FLUME_HOME/bin
```

After setting the path of flume directory, save and close the .bashrc file. And then in the terminal type the below command to update the .bashrc file.

```
[acadgild@localhost ~]$ source .bashrc
[acadgild@localhost ~]$ cat .bashrc
```

Step 10:

Note: Make sure you have below jars placed in your \$FLUME_HOME/lib directory:

1. twitter4j-core-X.XX.jar
2. twitter4j-stream-X.XX.jar
3. twitter4j-media-support-X.XX.jar

```
[acadgild@localhost apache-flume-1.8.0-bin]$ ls lib | grep twitter
flume-twitter-source-1.8.0.jar
twitter4j-core-3.0.3.jar
twitter4j-media-support-3.0.3.jar
twitter4j-stream-3.0.3.jar
```

Step 11:

Create a new file inside the conf directory inside the Flume-extracted directory.
Copy the Flume configuration code from the below link and paste it in the newly created file.

<https://drive.google.com/open?id=0B1QaXx7tpw3Sb3U4LW9SWINidkk>

```
[acagild@localhost apache-flume-1.8.0-bin]$ cd conf/
[acagild@localhost conf]$ pwd
/home/acagild/apache-flume-1.8.0-bin/conf
[acagild@localhost conf]$ ls
flume-conf.properties.template  flume-env.ps1.template  flume-env.sh.template  flume_twitter.conf  log4j.properties
[acagild@localhost conf]$ cat flume_twitter.conf
TwitterAgent.sources = Twitter
TwitterAgent.channels = MemChannel
TwitterAgent.sinks = HDFS

# Describing/Configuring the source
TwitterAgent.sources.Twitter.type = org.apache.flume.source.twitter.TwitterSource
TwitterAgent.sources.Twitter.consumerKey=09ieFY0Z3glT800cxhkAmcHN
TwitterAgent.sources.Twitter.consumerSecret=uy5NKjsyP19d2gujFitZfdTit25MT59NKTizkwBiFm9vUNPoJ5
TwitterAgent.sources.Twitter.accessToken=875303803464974337-ViWylha876bX1k2IhIqzRLf0KsfkMBb
TwitterAgent.sources.Twitter.accessTokenSecret=XDkzifEAcodZNxNY2URvvgLNkkH8Iql00fECiyMPpgB9Q
TwitterAgent.sources.Twitter.keywords=hadoop, bigdata, mapreduce, mahout, hbase, nosql
# Describing/Configuring the sink

TwitterAgent.sources.Twitter.keywords= hadoop,election,sports, cricket,Big data

TwitterAgent.sinks.HDFS.channel=MemChannel
TwitterAgent.sinks.HDFS.type=hdfs
TwitterAgent.sinks.HDFS.hdfs.path=hdfs://localhost:9000/user/flume/tweets
TwitterAgent.sinks.HDFS.hdfs.fileType=DataStream
TwitterAgent.sinks.HDFS.hdfs.writeformat=Text
TwitterAgent.sinks.HDFS.hdfs.batchSize=1000
TwitterAgent.sinks.HDFS.hdfs.rollSize=0
TwitterAgent.sinks.HDFS.hdfs.rollCount=10000
TwitterAgent.sinks.HDFS.hdfs.rollInterval=600

TwitterAgent.channels.MemChannel.type=memory
TwitterAgent.channels.MemChannel.capacity=10000
TwitterAgent.channels.MemChannel.transactionCapacity=1000

TwitterAgent.sources.Twitter.channels = MemChannel
TwitterAgent.sinks.HDFS.channel = MemChannel[acagild@localhost conf]$
```

Step 12:

We have to decide which keywords tweet data to be collected from the twitter application. So, you can change the keywords in the TwitterAgent.sources.Twitter.keywords command.

In our example, we are fetching tweet data related to Hadoop, election, sports, cricket and Big data.

Step 13:

Open a new terminal and start all the Hadoop daemons, before running the flume command to fetch the twitter data.

Use the 'jps' command to see the running Hadoop daemons.

```
[acagild@localhost ~]$ jps
3088 NameNode
3526 ResourceManager
3192 DataNode
3627 NodeManager
3388 SecondaryNameNode
4764 Jps
```

Step 14:

Create a new directory inside HDFS path, where the Twitter tweet data should be stored.

```
hadoop fs -mkdir -p /user/flume/tweets
```

```
[acadgild@localhost ~]$ hadoop fs -mkdir -p /user/flume/tweets
17/12/05 07:31:26 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
[acadgild@localhost ~]$ hadoop fs -ls /user/flume/tweets
17/12/05 07:31:37 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
```

Step 15:

For fetching data from Twitter, Use the below command to fetch the twitter tweet data into the HDFS cluster path.

flume-ng agent -n TwitterAgent -f <location of created/edited conf file>

flume-ng agent -n TwitterAgent -f /home/acadgild/apache-flume-1.8.0-bin/conf/flume_twitter.conf

The above command will start fetching data from Twitter and steams it into the HDFS given path.

```
[acadgild@localhost ~]$ flume-ng agent -n TwitterAgent -f /home/acadgild/apache-flume-1.8.0-bin/conf/flume_twitter.conf
Warning: No configuration directory set! Use --conf <dir> to override.
Info: Including Hadoop libraries found via (/usr/local/hadoop-2.6.0/bin/hadoop) for HDFS access
Info: Excluding /usr/local/hadoop-2.6.0/share/hadoop/common/lib/slf4j-api-1.7.5.jar from classpath
Info: Excluding /usr/local/hadoop-2.6.0/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar from classpath
Info: Including HBASE libraries found via (/usr/local/hbase/bin/hbase) for HBASE access
Info: Excluding /usr/local/hbase/lib/slf4j-api-1.6.4.jar from classpath
Info: Excluding /usr/local/hbase/lib/slf4j-log4j12-1.6.4.jar from classpath
Info: Excluding /usr/local/hadoop-2.6.0/share/hadoop/common/lib/slf4j-api-1.7.5.jar from classpath
Info: Excluding /usr/local/hadoop-2.6.0/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar from classpath
Info: Including Hive libraries found via (/usr/local/hive) for Hive access
17/12/05 07:33:56 INFO sink.DefaultSinkFactory: Creating instance of sink: HDFS, type: hdfs
17/12/05 07:33:56 INFO node.AbstractConfigurationProvider: Channel MemChannel connected to [Twitter, HDFS]
17/12/05 07:33:56 INFO node.Application: Starting new configuration: { sourceRunners: {Twitter=EventDrivenSourceRunner: { source:org.apache.flume.source.twitter.TwitterSource{name:Twitter,state:IDLE} }} sinkRunners: {HDFS=SinkRunner: { policy:org.apache.flume.sink.DefaultSinkProcessor@7f8e26af counterGroup:{ name:null counters:{} } }} channels: {MemChannel=org.apache.flume.channel.MemoryChannel{name: MemChannel} } }
17/12/05 07:33:56 INFO node.Application: Starting Channel MemChannel
17/12/05 07:33:56 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: CHANNEL, name: MemChannel: Successfully registered new MBean.
17/12/05 07:33:56 INFO instrumentation.MonitoredCounterGroup: Component type: CHANNEL, name: MemChannel started
17/12/05 07:33:56 INFO node.Application: Starting Sink HDFS
17/12/05 07:33:56 INFO node.Application: Starting Source Twitter
17/12/05 07:33:56 INFO twitter.TwitterSource: Starting twitter source org.apache.flume.source.twitter.TwitterSource{name:Twitter,state:IDLE}
17/12/05 07:33:56 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: SINK, name: HDFS: Successfully registered new MBean.
17/12/05 07:33:56 INFO instrumentation.MonitoredCounterGroup: Component type: SINK, name: HDFS started
17/12/05 07:33:56 INFO twitter.TwitterSource: Twitter source Twitter started.
17/12/05 07:33:56 INFO twitter4j.TwitterStreamImpl: Establishing connection.
17/12/05 07:34:14 INFO twitter4j.TwitterStreamImpl: Connection established.
17/12/05 07:34:14 INFO twitter4j.TwitterStreamImpl: Receiving status stream.
17/12/05 07:34:18 INFO hdfs.HDFSDataStream: Serializer = TEXT, UseRawLocalFileSystem = false
17/12/05 07:34:19 INFO hdfs.BucketWriter: Creating hdfs://localhost:9000/user/flume/tweets/FlumeData.1512439458728.tmp
17/12/05 07:34:19 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
17/12/05 07:34:20 INFO twitter.TwitterSource: Processed 100 docs
17/12/05 07:34:28 INFO twitter.TwitterSource: Processed 200 docs
17/12/05 07:34:39 INFO twitter.TwitterSource: Processed 300 docs
17/12/05 07:34:48 INFO twitter.TwitterSource: Processed 400 docs
17/12/05 07:34:50 INFO twitter.TwitterSource: Processed 500 docs
```

Once, the tweet data started streaming it into the given HDFS path we can use 'Ctrl+c' command to stop the streaming process.

Step 16:

To check the contents of the tweet data we can use the following command:

hadoop dfs -ls /user/flume/tweets

```
[acadgild@localhost ~]$ hadoop fs -ls /user/flume/tweets
17/12/05 07:36:22 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 1 items
-rw-r--r-- 1 acadgild supergroup 235283 2017-12-05 07:34 /user/flume/tweets/FlumeData.1512439458728.tmp
```

Step 17:

We can use the 'cat' command to display the tweet data inside the /user/flume/tweets/ path.

hadoop dfs -cat /user/flume/tweets/<flumeData file name>


```

17/12/05 07:36:41 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
{
  "type": "record", "name": "Doc", "doc": "adoc", "fields": [
    { "name": "id", "type": "string", { "name": "user_friends count", "type": [ "int", "null" ] }, { "name": "user location", "type": [ "string", "null" ] }, { "name": "user description", "type": [ "string", "null" ] }, { "name": "user statuses count", "type": [ "int", "null" ] }, { "name": "user followers count", "type": [ "int", "null" ] }, { "name": "user name", "type": [ "string", "null" ] }, { "name": "user screen name", "type": [ "string", "null" ] }, { "name": "created at", "type": [ "string", "null" ] }, { "name": "text", "type": [ "string", "null" ] }, { "name": "retweet count", "type": [ "long", "null" ] }, { "name": "retweeted", "type": [ "boolean", "null" ] }, { "name": "in_reply_to_user_id", "type": [ "long", "null" ] }, { "name": "source", "type": [ "string", "null" ] }, { "name": "in_reply_to_status_id", "type": [ "long", "null" ] }, { "name": "media url https", "type": [ "string", "null" ] }, { "name": "expanded url", "type": [ "string", "null" ] } ] }
}

J5jii100(2017-12-05T07:34:15Z) Fuck it, I got nothing to lose. <a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>

{
  "type": "record", "name": "Doc", "doc": "adoc", "fields": [
    { "name": "id", "type": "string", { "name": "user_friends count", "type": [ "int", "null" ] }, { "name": "user location", "type": [ "string", "null" ] }, { "name": "user description", "type": [ "string", "null" ] }, { "name": "user statuses count", "type": [ "int", "null" ] }, { "name": "user followers count", "type": [ "int", "null" ] }, { "name": "user name", "type": [ "string", "null" ] }, { "name": "user screen name", "type": [ "string", "null" ] }, { "name": "created at", "type": [ "string", "null" ] }, { "name": "text", "type": [ "string", "null" ] }, { "name": "retweet count", "type": [ "long", "null" ] }, { "name": "retweeted", "type": [ "boolean", "null" ] }, { "name": "in_reply_to_user_id", "type": [ "long", "null" ] }, { "name": "source", "type": [ "string", "null" ] }, { "name": "in_reply_to_status_id", "type": [ "long", "null" ] }, { "name": "media url https", "type": [ "string", "null" ] }, { "name": "expanded url", "type": [ "string", "null" ] } ] }

B59378651690908620928 @BTS twt I Love BTS 김태형 Kim Tae-Hyung

H026민 태형 > TaetaehyungJJK(2017-12-05T07:34:15Z) RT @vvsrkrk:

Fort Carson, CO I love when things are transparent, free and clear of all inhibition and judgement. Pharrell Williams

FullwoodTrey Fullwood(2017-12-05T07:34:15Z) RT @dream0genesis: Screamed! Did NeNe just say 'now Cynthia you and your wif need to calm down'?!

#noaa <a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>

0ツクツクボウシ botTKTK BOUHSHT(2017-12-05T07:34:15Z)

<a href="http://twittbot.net/" rel="nofollow">twittbot.net</a>

이나즈마 프리파라를 좋아하는 집력
프리파라 우왕(빅출몰)
재물, 가족, 보석 전문장인 /
미국계는 이쪽에서 놀라더라고요!
日本語, English ok/
出さ by 70キビ

<a href="http://twitter.com/download/android" rel="nofollow">Twitter for Android</a>
937865169081851904_Director, Training Services @CannabisCCI
Excited to be building quality training built by industry for industry. Opinions expressed here are my own. Shannon KloetShannonKloet

```