

MSDS610 Week 5 Pig Assignment - Nathan Worsham

Similar to previous weeks, I began this assignment by powering up my week 1 VM and then getting the download and unpacking the archive. And again similar to previous weeks I installed Pig into the home directory of my hadoop user, created a symbolic link, updated the environment variables for Pig, and then ran the test command successfully that the instructions called for.

```
[hadoop@week1 ~]$ wget http://apache.arvixe.com/pig/latest/pig-0.15.0.tar.gz
--2016-02-08 21:44:58-- http://apache.arvixe.com/pig/latest/pig-0.15.0.tar.gz
Resolving apache.arvixe.com (apache.arvixe.com)... 198.58.67.82
Connecting to apache.arvixe.com (apache.arvixe.com)[198.58.67.82]:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 128917625 (115M) [application/x-gzip]
Saving to: 'pig-0.15.0.tar.gz'

100%[=====] 128,917,625  (6.85 MB/s) -> 'pig-0.15.0.tar.gz' saved [128917625/128917625]

[hadoop@week1 ~]$ tar -xvf pig-0.15.0.tar.gz
pig-0.15.0/
pig-0.15.0/bin/
pig-0.15.0/conf/
pig-0.15.0/contrib/
pig-0.15.0/contrib/piggybank/
pig-0.15.0/contrib/piggybank/java/
pig-0.15.0/contrib/piggybank/java/build/
pig-0.15.0/contrib/piggybank/java/build/classes/
pig-0.15.0/contrib/piggybank/java/build/classes/org/
pig-0.15.0/contrib/piggybank/java/build/classes/org/apache/
pig-0.15.0/contrib/piggybank/java/build/classes/org/apache/pig/
pig-0.15.0/contrib/piggybank/java/build/classes/org/apache/pig/piggybank/
pig-0.15.0/contrib/piggybank/java/build/classes/org/apache/pig/piggybank/evaluation/
pig-0.15.0/contrib/piggybank/java/build/classes/org/apache/pig/piggybank/evaluation/datetime/

[hadoop@week1 ~]$ vi .bashrc
[hadoop@week1 ~]$ tail -2 .bashrc
export PIG=/home/hadoop/pig
export PATH=$PATH:$PIG/bin
[hadoop@week1 ~]$ . .bashrc
[hadoop@week1 ~]$ /pig/bin/pig -help

Apache Pig version 0.15.0 (r1682971)
compiled Jun 01 2015, 11:44:35

USAGE: Pig [options] [-] : Run interactively in grunt shell.
      Pig [options] -e[xcute] cmd [cmd ...] : Run cmd(s).
      Pig [options] [-f[file]] file : Run cmds found in file.
options include:
  -d, --logd[conf - Logd] configuration file, overrides log.conf
```

I went ahead and copied the interactive mode example with the `/etc/passwd` file in local mode just to confirm it was working as expected. It did work as expected with the exception that it had some warnings of deprecated features.

```
[hadoop@week1 ~]$ pig -x local
16/02/09 02:06:41 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
16/02/09 02:06:41 INFO pig.ExecTypeProvider: Picked LOCAL as the ExecType
2016-02-09 02:06:41,204 [main] INFO org.apache.pig.Main - Apache Pig version 0.15.0 (
2016-02-09 02:06:41,207 [main] INFO org.apache.pig.Main - Logging error messages to:
2016-02-09 02:06:41,227 [main] INFO org.apache.pig.impl.util.Utils - Default bootup f
2016-02-09 02:06:41,542 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
2016-02-09 02:06:41,543 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
er.address
2016-02-09 02:06:41,543 [main] INFO org.apache.pig.backend.hadoop.executionengine.HEx
2016-02-09 02:06:41,657 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
checksum
grunt> A = load '/etc/passwd' using PigStorage(':');
2016-02-09 02:07:21,588 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
checksum
grunt> B = foreach A generate $0 as id;
2016-02-09 02:07:21,589 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
grunt> dump B;
```

```
(root)
(bin)
(daemon)
(admin)
(lp)
(sync)
(shutdown)
(halt)
(mail)
(operator)
(games)
(ftp)
(nobody)
(avahi-autoipd)
(systemd-bus-proxy)
(systemd-network)
(dbus)
(polkitd)
(tss)
(postfix)
(chrony)
(sshd)
(nworsham)
(hadoop)
grunt>
```

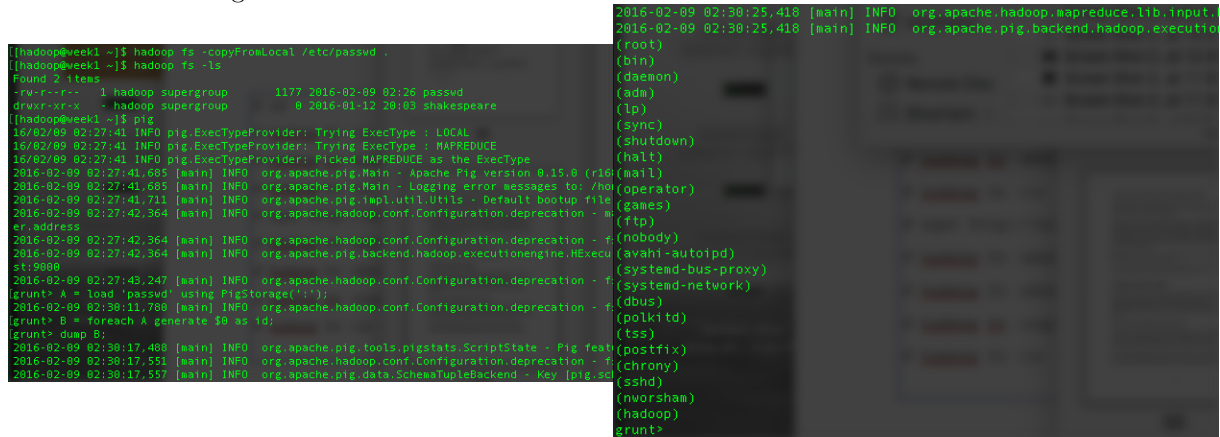
Next I wanted to do the exact same test but in Hadoop. Looking at the hadoop exercise, there needs to be an environment variable set for `PIG_CLASSPATH` which it gives an example as `"/mycluster/conf"`. I went ahead and looked for a `conf` directory at the root of the hadoop install but could find no such folder, so then I used the `find` command. None of the results seemed right and after carefully reading I noticed the instructions state "the directory that contains the `core-site.xml`". So re-running my `find` command targeting the `core-site.xml`—a file I was definitely familiar with from past weeks—I found the directory was `hadoop/etc/hadoop`. So I set my environment variable accordingly.

```
[hadoop@week1 ~]$ find hadoop -type d -iname "conf"
hadoop/share/hadoop/httpfs/localhost/conf
hadoop/share/doc/hadoop/api/src-html/org/apache/hadoop/yarn/conf
hadoop/share/doc/hadoop/api/src-html/org/apache/hadoop/conf
hadoop/share/doc/hadoop/api/org/apache/hadoop/yarn/sls/conf
hadoop/share/doc/hadoop/api/org/apache/hadoop/yarn/conf
hadoop/share/doc/hadoop/api/org/apache/hadoop/conf
[hadoop@week1 ~]$ find hadoop -type f -iname "core-site.xml"
hadoop/share/hadoop/common/templates/core-site.xml
hadoop/etc/hadoop/core-site.xml
[hadoop@week1 ~]$ vi .bashrc
[hadoop@week1 ~]$ tail .bashrc
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native
export PATH=$PATH:$HADOOP_HOME/sbin:$HADOOP_HOME/bin
export JAVA_LIBRARY_PATH=$HADOOP_HOME/lib/native:$JAVA_LIBRARY_PATH

export HIVE=/home/hadoop/hive
export PATH=$PATH:HIVE/bin

export PIG=/home/hadoop/pig
export PIG_CLASSPATH=$HADOOP_HOME/etc/hadoop
export PATH=$PATH:$PIG/bin
[hadoop@week1 ~]$ source .bashrc
```

Now I started up HDFS and yarn then ran just simply `pig` since the default mode is mapreduce mode. It started up with again the previous deprecation warnings. I was pleasantly surprised to learn that `pig` retains a command history that I could access with the up and down arrows from my keyboard. I was equally as pleased that running the same commands in mapreduce mode returned the same output, except that the return took a little longer.



```

(hadoop@week1 ~)$ hadoop fs -copyFromLocal etc/passwd .
(hadoop@week1 ~)$ hadoop fs -ls
Found 2 items
-rw-r--r-- 1 hadoop supergroup 1177 2016-02-09 02:26 passwd
drwxr-xr-x 2 hadoop supergroup 0 2016-01-12 20:03 shakespeare
(hadoop@week1 ~)$ pig
2016-02-09 02:27:41 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
2016-02-09 02:27:41 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
2016-02-09 02:27:41 INFO pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2016-02-09 02:27:41,685 [main] INFO org.apache.pig.Main - Apache Pig version 0.15.0 (r16
2016-02-09 02:27:41,685 [main] INFO org.apache.pig.Main - Logging error messages to: /ho
2016-02-09 02:27:41,711 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file
2016-02-09 02:27:42,364 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - W
ar address
2016-02-09 02:27:42,364 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - f
2016-02-09 02:27:42,364 [main] INFO org.apache.pig.backend.hadoop.executionengine.MExecu
st:9000
2016-02-09 02:27:43,247 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - f
(grunt> A = load 'passwd' using PigStorage('');
2016-02-09 02:30:11,700 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - f
(grunt> B = foreach A generate $0 as 'id';
(grunt> dump B;
2016-02-09 02:30:17,480 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig feat
2016-02-09 02:30:17,551 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - f
2016-02-09 02:30:17,557 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig so
(grunt>

```

I decided the next step would be good to do the `pig` tutorial for local and mapreduce mode. The first instruction was to make sure that `JAVA_HOME` was set, which I checked and it was, but then it also asked for another environment variable `PIG_HOME` to be set to the same as the previous environment variable `PIG` which I went ahead and did. Next I needed to run the `"ant"` command in the tutorial directory, a command I had never heard of. I had to install the command first, looks like from its description it is a "Build tool for java". However after running the command, I did not receive the expected output, I received a couple of errors.

```

[echo] *** Compiling Tutorial files ***
[javac] /home/hadoop/pig-0.15.0/tutorial/build.xml:66: warning: 'includeantruntime'
was not set, defaulting to build.sysclasspath=last; set to false for repeatable builds
[javac] Compiling 7 source files to /home/hadoop/pig-0.15.0/tutorial/build/classes

```

BUILD FAILED

```

/home/hadoop/pig-0.15.0/tutorial/build.xml:66: /home/hadoop/pig-0.15.0/build/ivy/lib/
Pig does not exist.

```

I was able to fix the `includeantruntime` warning by following a [Stackoverflow.com](http://stackoverflow.com) (2011) thread and adding `includeantruntime="false"` inside of the `<javac>` options. The final error though was about a build directory not existing at the root level of `pig`. Taking a look at that area, indeed I saw no build directory but I did see a `"build.xml"` file just like there was in the tutorial directory. Given this info I assumed I would first need to run `ant` at the root level but was a bit worried that might break something as the installation of `pig` I chose was the already built version and not the source. I went ahead and copied the entire `pig` directory to a test directory and ran the `ant` command on the test instance. Sure enough it built the missing directories and running `ant` inside of the tutorial directory now worked and I now had my `pigtutorial.tar.gz` file.

```

jar:
[echo] svnString : unknown
[jar] Building jar: /home/hadoop/pig_test/build/pig-0.15.0-SNAPSHOT.jar
[echo] svnString : unknown
[jar] Building jar: /home/hadoop/pig_test/build/pig-0.15.0-SNAPSHOT-withouthadoop.jar
[jar] META-INF/ASL2.0 already added, skipping
[jar] META-INF/LICENSE already added, skipping
[jar] META-INF/NOTICE already added, skipping

copyCommonDependencies:

copyH1Dependencies:
[copy] Copying 1 file to /home/hadoop/pig_test
[move] Moving 1 file to /home/hadoop/pig_test/legacy

copyH2Dependencies:

BUILD SUCCESSFUL
Total time: 3 minutes 34 seconds
[[hadoop@week1 pig_test]$ cd tutorial/
[[hadoop@week1 tutorial]$ ls
build build.xml data scripts src
[[hadoop@week1 tutorial]$ ant
Buildfile: /home/hadoop/pig_test/tutorial/build.xml

init:

compile:
[echo] *** Compiling Tutorial files ***
[javac] Compiling 7 source files to /home/hadoop/pig_test/tutorial/build/classes
[javac] warning: [options] bootstrap class path not set in conjunction with -source 1.5
[javac] warning: [options] source value 1.5 is obsolete and will be removed in a future release
[javac] warning: [options] target value 1.5 is obsolete and will be removed in a future release
[javac] warning: [options] To suppress warnings about obsolete options, use -Xlint:options.
[javac] 4 warnings

jar:
[echo] *** Creating tutorial.jar ***
[jar] Building jar: /home/hadoop/pig_test/tutorial/build/output/pigtap/tutorial.jar

cp:
[echo] *** Preparing tar creation ***
[copy] Copying 6 files to /home/hadoop/pig_test/tutorial/build/output/pigtmp

tar:
[echo] *** Creating tutorial.jar ***
[tar] Building tar: /home/hadoop/pig_test/tutorial/build/pigtutorial.tar
[gzip] Building: /home/hadoop/pig_test/tutorial/pigtutorial.tar.gz

BUILD SUCCESSFUL
Total time: 2 seconds
[[hadoop@week1 tutorial]$ ls
build build.xml data pigtutorial.tar.gz scripts src

```

After extracting the tar file, I had the remaining files I needed. First running the local test—`script1-local.pig`—which appears to read the “excite-small.log”, a log file of the Excite search engine and find the popular query phrases as they relate to times of the day. However the output from the local mode ended up with only single popular words, not phrases, and it seems to include stop words that probably should have been filtered out to make a better analysis but regardless shows the ability of the pig system.

```

[[hadoop@week1 pigtmp]$ cd script1-local-results.txt/
[[hadoop@week1 script1-local-results.txt]$ ls
part-r-000000 _SUCCESS
[[hadoop@week1 script1-local-results.txt]$ ls -l
total 4
-rw-r--r-- 1 hadoop hadoop 838 Feb  9 07:29 part-r-000000
-rw-r--r-- 1 hadoop hadoop  0 Feb  9 07:29 _SUCCESS
[[hadoop@week1 script1-local-results.txt]$ cat part-r-000000
07 new 2.4494897427831788 2 1.1428571428571426
08 pictures 2.04939015319192 3 1.4999999999999998
08 computer 2.4494897427831788 2 1.1428571428571426
08 s 2.545584412271571 3 1.3636363636363635
10 free 2.2657896674810685 4 1.923076923076923
10 to 2.6457513110645983 2 1.125
10 pics 2.794802794804192 3 1.3076923076923075
10 school 2.828427124746188 2 1.1111111111111114
11 pictures 2.04939015319192 3 1.4999999999999998
11 in 2.1572774865200244 3 1.4285714285714284
13 the 3.1309398305840723 6 1.9375
14 music 2.1105794120443453 4 1.6666666666666667
14 city 2.2368679774997982 2 1.6666666666666665
14 university 2.412090756622189 3 1.4000000000000001
15 adult 2.8284271247461883 2 1.1111111111111112
17 chat 2.9184275004359965 3 1.2857142857142854
19 in 2.1572774865200244 3 1.4285714285714284
19 car 2.23686797749979 3 1.3333333333333333
[[hadoop@week1 script1-local-results.txt]$

```

```

Input(s):
Successfully read 4581 records from: "file:///home/hadoop/pigtap/excite-small.log"

Output(s):
Successfully stored 18 records in: "file:///home/hadoop/pigtmp/script1-local-results.txt"

Counters:
Total records written : 18
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:

```

Next I could try the mapreduce version of the exercise. I started by copying a bz2 version of the log file to the HDFS. Looking at the size of the file compared to the local version, it is considerably bigger—10 mb versus 204 kb. I had already taken care of the environment variable `PIG_CLASSPATH` but now I needed to make `HADOOP_CONF_DIR` equivalent to the same path.

```

[hadoop@week1 ~]$ cd pigtmp/
[hadoop@week1 pigtmp]$ ls
excite.log.bz2 excite-small.log script1-hadoop.pig script1-local.pig scrip
[hadoop@week1 pigtmp]$ hadoop fs -copyFromLocal ~/pigtmp/excite.log.bz2 .
[hadoop@week1 pigtmp]$ ls -l
total 10480
-rw-r--r-- 1 hadoop hadoop 10488717 Feb  9 07:01 excite.log.bz2
-rw-r--r-- 1 hadoop hadoop 208348 Feb  9 07:01 excite-small.log
-rw-r--r-- 1 hadoop hadoop 3835 Feb  9 07:01 script1-hadoop.pig
-rw-r--r-- 1 hadoop hadoop 3820 Feb  9 07:01 script1-local.pig
drwxrwxr-x 2 hadoop hadoop  84 Feb  9 07:29 script1-local-results.txt
-rw-r--r-- 1 hadoop hadoop 3489 Feb  9 07:01 script2-hadoop.pig
-rw-r--r-- 1 hadoop hadoop 3488 Feb  9 07:01 script2-local.pig
-rw-r--r-- 1 hadoop hadoop 10703 Feb  9 07:01 tutorial.jar
[hadoop@week1 pigtmp]$ ls -lh
total 11M
-rw-r--r-- 1 hadoop hadoop 10M Feb  9 07:01 excite.log.bz2
-rw-r--r-- 1 hadoop hadoop 204K Feb  9 07:01 excite-small.log
-rw-r--r-- 1 hadoop hadoop 3.8K Feb  9 07:01 script1-hadoop.pig
-rw-r--r-- 1 hadoop hadoop 3.8K Feb  9 07:01 script1-local.pig
drwxrwxr-x 2 hadoop hadoop  84 Feb  9 07:29 script1-local-results.txt
-rw-r--r-- 1 hadoop hadoop 3.5K Feb  9 07:01 script2-hadoop.pig
-rw-r--r-- 1 hadoop hadoop 3.4K Feb  9 07:01 script2-local.pig
-rw-r--r-- 1 hadoop hadoop 11K Feb  9 07:01 tutorial.jar
[hadoop@week1 pigtmp]$ vi ~/.bashrc
[hadoop@week1 pigtmp]$ tail ~/.bashrc
export JAVA_LIBRARY_PATH=$HADOOP_HOME/lib/native:$JAVA_LIBRARY_PATH

export HIVE=/home/hadoop/hive
export PATH=$PATH:$HIVE/bin

export PIG=/home/hadoop/pig
export PIG_HOME=/home/hadoop/pig
export PIG_CLASSPATH=$HADOOP_HOME/etc/hadoop
export HADOOP_CONF_DIR=$PIG_CLASSPATH
export PATH=$PATH:$PIG/bin
[hadoop@week1 pigtmp]$

```

I was now ready to run the mapreduce version of the script. As expected it took considerably longer considering the much bigger file to analyze and also being within the HDFS.

```

Input(s):
Successfully read 944954 records (31338918 bytes) from: "hdfs://localhost:9000/user/hadoop/excite.log.bz2"

Output(s):
Successfully stored 13530 records (206471324 bytes) in: "hdfs://localhost:9000/user/hadoop/script1-hadoop-results"

```

The results were much bigger than the local script test—13530 lines versus 18—and this time did have phrases in addition to single words. The results seemed to cover all hours of the day.

```

[hadoop@week1 pigtmp]$ hadoop fs -ls
Found 4 items
-rw-r--r-- 1 hadoop supergroup 10488717 2016-02-09 12:58 excite.log.bz2
-rw-r--r-- 1 hadoop supergroup 1177 2016-02-09 02:26 passwd
drwxr-xr-x - hadoop supergroup 0 2016-02-09 13:12 script1-hadoop-results
drwxr-xr-x - hadoop supergroup 0 2016-01-12 20:03 shakespeare
[hadoop@week1 pigtmp]$ hadoop fs -ls script1-hadoop-results
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2016-02-09 13:12 script1-hadoop-results/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 659954 2016-02-09 13:12 script1-hadoop-results/part-r-00000
[hadoop@week1 pigtmp]$ hadoop fs -tail -20 script1-hadoop-results/part-r-00000
-tail: Illegal option -20
Usage: hadoop fs [generic options] -tail [-f] <file>
[hadoop@week1 pigtmp]$ hadoop fs -tail script1-hadoop-results/part-r-00000
foto 2.691946385511004 3 1.3571428571428568
23 shauna 2.7714340304599967 4 1.625
23 bosch 2.7940079404491 3 1.3076923076923077
23 hottest 2.7950809018747373 3 1.3333333333333333
23 alanaak 2.81858908754796 9 2.782608695652174
23 chilton 2.82842712474619 3 1.2222222222222223
23 red hot 2.82842712474619 2 1.1111111111111112
23 jerk off 2.82842712474619 2 1.1111111111111112
23 live video 2.82842712474619 2 1.1111111111111112
23 brothels 2.82842712474619 2 1.1111111111111112
23 bubbles 2.82842712474619 2 1.1111111111111112
23 ping 2.82842712474619 5 1.8888888888888889
23 orgasmic 2.82842712474619 3 1.2222222222222223
23 highway patrol 2.82842712474619 2 1.1111111111111112
23 jenni cam 3.0 2 1.0999999999999999
23 cheryl bachman 3.162277660160379 2 1.0909090909090909
23 wallpaper 3.2068995827040717 6 1.9333333333333333
23 ap 3.316624790353599 2 1.0033333333333333
23 bachman 3.4641016151377557 2 1.0769230769230769
23 kinaberley 3.6055512754639896 2 1.0714285714285712
23 jerk 3.6055512754639896 2 1.0714285714285712
[hadoop@week1 pigtmp]$ hadoop fs -less script1-hadoop-results/part-r-00000
-less: Unknown command
[hadoop@week1 pigtmp]$ hadoop fs -cat script1-hadoop-results/part-r-00000|more
00 and shareware 2.112885836821291 3 1.5294117647058825
00 vcd 2.121328343194424 3 1.5
00 bluebird 2.1555530241167826 4 1.9285714285714282
00 cute 2.1650635894610955 4 1.0571428571428574
00 chested 2.182020625326997 4 1.75
00 diablo cheats 2.197401062294143 3 1.4705882352941178
00 psynosis 2.23686797749979 2 1.1666666666666667
00 vacancy 2.2360679774997982 2 1.1666666666666665
00 pennysive 2.2368679774997982 2 1.1666666666666665
00 worisette 2.2368679774997982 2 1.1666666666666665
00 labyrinth 2.2368679774997982 2 1.1666666666666665

```

Furthermore it is interesting to compare the results at different hours of the day, example the noon hour ngram results seem to be "cleaner" than the 11 and 12 am results.

References

pig.apache.org, 2016. Retrieved from <http://pig.apache.org/docs/r0.14.0/start.html>

Stackoverflow.com, 2011. Retrieve from <http://stackoverflow.com/questions/5103384/ant-warning-includeantruntime-was-not-set>

pig.apache.org, 2016. Retrieved from <https://pig.apache.org/docs/r0.10.0/udf.html#piggybank>