

## MSDS610 Week 8 Spark Continued Assignment - Nathan Worsham

I was very interested in this weeks assignment as it pertains to what I already work in - IT Security. I started by getting the download link and getting the data.

```
[hadoop@week1 ~]$ wget http://kdd.ics.uci.edu/databases/kddcup99/kddcup.data.gz
--2016-03-01 11:41:30-- http://kdd.ics.uci.edu/databases/kddcup99/kddcup.data.gz
Resolving kdd.ics.uci.edu (kdd.ics.uci.edu)... 128.195.1.95
Connecting to kdd.ics.uci.edu (kdd.ics.uci.edu)[128.195.1.95]:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 18115902 (17M) [application/x-gzip]
Saving to: 'kddcup.data.gz'

46% [=====] 8,400,464 682KB/s eta 12s 46% [=====] 8,503,272 687KB/s eta 47% [=====] 8,529,832 47% [=====] 8,556,792 692KB/s eta 11s
51% [=====] 8,583,752 697KB/s eta 10s 51% [=====] 8,610,712 702KB/s eta 9s
eta 11s1% [=====] 8,637,672 707KB/s eta 8s 1% [=====] 8,664,632 712KB/s eta 7s
100% [=====] 18,115,902 721KB/s in 23s

2016-03-01 11:41:53 (786 KB/s) - 'kddcup.data.gz' saved [18115902/18115902]
```

When I tried to unpack it I received errors that `tar: This does not look like a tar archive`, I had even tried providing the `z` option but received the same problem and tried redownloading the file incase it was corrupted. After looking at a superuser.com (2014) question, I realized that the file likely is just a Gzip file and not a TAR. So I then used `gzip` to unpack it.

```
[hadoop@week1 week0_data]$ gzip -d kddcup.data.gz
[hadoop@week1 week0_data]$ ls
kddcup.data kddcup.data.gz.old
[hadoop@week1 week0_data]$ rm kddcup.data.gz.old
[hadoop@week1 week0_data]$ ls -l
total 725176
-rw-r--r-- 1 hadoop hadoop 742579829 Jun 26 2007 kddcup.data
[hadoop@week1 week0_data]$
```

I went ahead and created a new folder in the HDFS and then copied the kdd data to it.

```
[hadoop@week1 ~]$ hadoop fs -mkdir week8
[hadoop@week1 ~]$ hadoop fs -ls
Found 7 items
-rw-r--r-- 1 hadoop supergroup 10408717 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 2016-02-10 09:56 script2-hadoop-results
drwxr-xr-x - hadoop supergroup 0 2016-01-12 20:03 shakespeare
drwxr-xr-x - hadoop supergroup 0 2016-02-15 22:34 user
drwxr-xr-x - hadoop supergroup 0 2016-03-01 11:55 week8
[hadoop@week1 ~]$ hadoop fs -copyFromLocal ~/week8_data/kddcup.data week8
[hadoop@week1 ~]$ hadoop fs -ls week8
Found 1 items
-rw-r--r-- 1 hadoop supergroup 742579829 2016-03-01 11:56 week8/kddcup.data
[hadoop@week1 ~]$
```

The spark shell is opened and the csv file is loaded as an RDD.

```
scala> val rawData = sc.textFile("hdfs://localhost:9000/user/hadoop/kddcup.data")
16/03/01 16:38:18 INFO storage.MemoryStore: Block broadcast_0 stored as values in memory (estimated size 61.8 KB, free 61.8 KB)
16/03/01 16:38:18 INFO storage.MemoryStore: Block broadcast_0_piece0 stored as bytes in memory (estimated size 19.4 KB, free 81.2 KB)
16/03/01 16:38:18 INFO spark.BlockManagerInfo: Added broadcast_0_piece0 in memory on localhost:60370 (size: 19.4 KB, free: 517.4 MB)
16/03/01 16:38:18 INFO spark.SparkContext: Created broadcast 0 from textFile at <console>:27
rawData: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[1] at textFile at <console>:27
```

I stumbled a bit on the location for the CSV file as I didn't realize the `hadoop fs` command starts off in the user's home directory similar to Linux as the command kept telling me it could not find the file. Once I realized the full path was `/user/hadoop/` I moved the file in the HDFS and then was able to run the command to see what labels where in the data.

Next the categorical values are removed as K-means clustering requires numerical data.

The Kmeans implementation is imported and then run on the "data" which took some time to run. Then the centroids are printed.

The previous model only did a fitting of 2 clusters on the data. So seeing how the existing labels lands in each cluster, all land in one cluster except for 1 "portsweep" which is a poor model then.

2

Create a function to define a Euclidean distance and another to return the distance for the data point to its cluster's centroid. This will help determine the number of clusters to use.

```
scala> def distance(a: Vector, b: Vector) =
  | math.sqrt(a.toArray.zip(b.toArray).
  | map(p => p._1 - p._2).map(d => d * d).sum)
distance: (a: org.apache.spark.mllib.linalg.Vector, b: org.apache.spark.mllib.linalg.Vector)Double

scala> def distToCentroid(datum: Vector, model: KMeansModel) = {
  | val cluster = model.predict(datum)
  | val centroid = model.clusterCenters(cluster)
  | distance(centroid, datum)
  | }
distToCentroid: (datum: org.apache.spark.mllib.linalg.Vector, model: org.apache.spark.mllib.clustering.KMeansModel)Double
```

Now another function is created to measure the average distance from the data points to the centroids.

```
scala> import org.apache.spark.rdd._
import org.apache.spark.rdd._

scala> def clusteringScore(data: RDD[Vector], k: Int) = {
  | val kmeans = new KMeans()
  | kmeans.setK(k)
  | val model = kmeans.run(data)
  | data.map(datum => distToCentroid(datum, model)).mean()
  | }
clusteringScore: (data: org.apache.spark.rdd.RDD[org.apache.spark.mllib.linalg.Vector], k: Int)Double
```

Now use this function to evaluate values of k from 5 to 40, counting by 5s. Like the example, the value actually increases for 35, though they received a better output for 30 than I did.

```
16/03/02 01:43:00 INFO executor.Executor: Finished task 5.0 in stage 340.0 (TID 2040): 2789 bytes result sent to driver
16/03/02 01:43:00 INFO scheduler.TaskSetManager: Finished task 5.0 in stage 340.0 (TID 2040) in 7434 ms on localhost (6/6)
16/03/02 01:43:00 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 340.0, whose tasks have all completed, from pool
16/03/02 01:43:00 INFO scheduler.DAGScheduler: ResultStage 340 (mean at <console>:54) finished in 72.451 s
16/03/02 01:43:00 INFO scheduler.DAGScheduler: Job 220 finished: mean at <console>:54, took 72.467009 s
(5,1938.8583418859188)
(10,1592.2768468128086)
(15,1407.9929596925513)
(20,1300.2743912730707)
(25,1041.6914354420257)
(30,866.4511133122714)
(35,883.4160185558521)
(40,521.4234074648682)

scala> 16/03/02 01:59:41 INFO storage.BlockManagerInfo: Removed broadcast_497_piece0 on localhost:37837 in memory (size: 12.5 KB, free: 74.1 MB)
16/03/02 01:59:41 INFO spark.ContextCleaner: Cleaned accumulator 443
```

Running the iteration longer and decrease the threshold for the minimum amount considered significant and then finally run in parallel with `par.map`. Took longer to run than the previous but in my results the output starts to increase at 100 with 90 being the best.

```
scala> kmeans.setRuns(10)
warning: there were 1 deprecation warning(s); re-run with -deprecation for details
res5: kmeans.type = org.apache.spark.mllib.clustering.KMeans@7aeb3664

scala> kmeans.setEpsilon(1.0e-6)
res6: kmeans.type = org.apache.spark.mllib.clustering.KMeans@7aeb3664
```

```
(30,936.6630608255432)
(40,404.50119276950556)
(50,325.4491402759626)
(60,294.09423933023567)
(70,308.4794767407293)
(80,314.0387545764995)
(90,185.08670187535128)
(100,227.5759505959566)
```

Changing each feature into a Z-score or standard score can help normalize the output. This is done by subtracting the mean and dividing by the standard deviation. Various outputs leading up to the final output.

n

```
16/03/02 07:39:59 INFO scheduler.DAGS
n: Long = 4898431
```

sums



```
scala> val kmeans = new KMeans()
kmeans: org.apache.spark.mllib.clustering.KMeans = org.apache.spark.mllib.clustering.KMeans@3689b897

scala> kmeans.setK(150)
res5: kmeans.type = org.apache.spark.mllib.clustering.KMeans@3689b897

scala> kmeans.setRuns(10)
warning: there were 1 deprecation warning(s): re-run with -deprecation for details
res6: kmeans.type = org.apache.spark.mllib.clustering.KMeans@3689b897

scala> kmeans.setEpsilon(1.0e-6)
res7: kmeans.type = org.apache.spark.mllib.clustering.KMeans@3689b897

scala> val parseFunction = buildCategoricalAndLabelFunction(rawData)
```

Again I was trying my best at following along with their code from git, but this time I was resolved to run it in pieces if I could. I also shutdown my VM and gave it another cpu and another GB of RAM to try to help the number crunching and time to complete.

```
16/03/04 09:42:57 INFO scheduler.DAGScheduler: Job 6 finished: aggregate at <console>:46, took 54.153574 s
normalizeFunction: org.apache.spark.mllib.linalg.Vector => org.apache.spark.mllib.linalg.Vector => <function1>

scala> val normalizedData = data.map(normalizeFunction)
normalizedData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.linalg.Vector] = MapPartitionsRDD[10] at map at <console>:54

scala> normalizedData.cache()
res8: normalizedData.type = MapPartitionsRDD[10] at map at <console>:54

scala> val model = kmeans.run(normalizedData)
16/03/04 09:46:18 INFO spark.SparkContext: Starting job: takeSample at KMeans.scala:376
16/03/04 09:46:18 INFO scheduler.DAGScheduler: Got job 7 (takeSample at KMeans.scala:376) with 6 output partitions
16/03/04 09:46:18 INFO scheduler.DAGScheduler: Final stage: ResultStage 10 (takeSample at KMeans.scala:376)
```

But as it turns out the only portion that took time to compute was the `val model = kmeans.run(normalizedData)`. As expected there are many more centroids now.

```
scala> model.clusterCenters.foreach(println)
[-0.06683318542426188,1.2722595694683434,-1.1714318212179429,-0.28322762881782588,-0.01483648485872775,
7426275E-4,-0.01467828452095886,-4.518263416747377E-4,-6.389798854583881E-4,-0.014677243787549586,-0.0,
21018,-0.1892557967881207,-0.014586515531105643,-0.10980862368986825,-0.0052498194968180655,-0.0201141,
-0.014677243787549586,-0.014649387142574893,1.8468268527162669,-0.01468758277633199,-0.05784987536945,
4328429728805,-0.014705047547807803,-0.014684199598836933,-0.005496787161547559,-0.014579513074798895,
58692662,-0.02628501806798148,-0.0148819588658975561,-0.014663322634485447,-0.00156517411733783,-0.01,
62862,-0.13848809597977588,-0.0148819588658975561,-0.03315281719839832,-0.014579513074798895,-0.01476804,
-0.01471199838464653,-0.01468758277633199,-0.014656356243181384,-0.014558485514561297,-0.014572587258,
4642414727862962,-0.014705047547807803,-0.014656356243181384,-0.014705047547807803,-0.0018103146265166,
74328429728805,-0.014656356243181384,-0.014815741589747728,-0.01468758277633199,-1.1687641998916775,-0.0,
6270721,-0.014684199598836933,-6.389798854583881E-4,-0.027984871691862938,-0.014691152199563575,-0.014,
1E-4,-0.0013554881318982463,-0.017956958181665682,-0.0144387514875846,-0.014684199598836933,-0.0147674,
-0.005733127284811937,-0.048682928369985,-0.00341123357221592,-0.00319491868898051,-1.081211279010686,
29169395,2.152080863754422,-0.004990649788216317,-0.01457258725872132,-0.0019487588281734776,-0.00169,
88,-0.015139173434892635,-0.0011834846181561417,-0.026528075953889122,-0.00459881558475935,-0.4893676,
-0.008257788397449333,-0.0045464613866548836,-0.0032845891671842145,-0.009572339215774545,-0.008584578,
-6.389798854583881E-4,-0.028911383443148115,-0.42157255579177785,-1.1536488443325582,2.15225511853138,
32962357,-0.24887165481126185,-1.0111988618780855,0.4685337881289249,-0.28845236399883286,0.34199278889,
28814158997,0.32418619392437714,-1.25789689659774,-0.1566684879543488,-0.1522486852786887,2.15113227454,
319657773663]
[-0.05886851958868315,1.2722595694697427,-1.171431821217578,-0.28322762881730253,-0.01483648485872775,
432887E-4,-0.014678284520992564,-4.5182634167587613E-4,-6.389798854583881E-4,-0.014677243787549586,-0.0,
538978,-0.18925579678808198,-0.014586515531108886,-0.1269958954286929,-0.005249819496818084,-0.020114,
-0.014677243787563112,-0.014649387142687782,-0.5384888243418645,-0.01468758277632066,-0.0578498753695,
774328429777144,-0.01470504754779885,-0.014684199598857688,-0.005496787161547477,-0.014579513074832853,
2622882,-0.026285018067981577,-0.014881958865899901,-0.014663322634448174,-0.0015651741173529268,-0.01,
77435,-0.141881863966876,-0.014881958865899901,-0.03315281719868387,-0.014579513074832853,-0.014768049,
-0.014711998384657857,-0.01468758277632066,-0.014656356243221253,-0.01455848551461846,-0.014572587258,
464241472777635,-0.01470504754779885,-0.014656356243221253,-0.01470504754779885,-0.001810314626517658]
```

After that I went ahead and tried printing the clusters with labels as before like the book showed using their `val clusterLabelCount` but I received several errors and a driver stacktrace that I wasn't sure about so I moved onto the anomaly detection. Here in the git code the line

```
val anomalyDetector = buildAnomalyDetector(data, normalizeFunction)
```

calls the entire previous function, so I ran the rest of that function and then assumed that the value `anomalyDetector` needed to equal the output of the function `buildAnomalyDetector`.

```
scala> val anomalyDetector = buildAnomalyDetector(data, normalizeFunction)
<console>:57: error: not found: value buildAnomalyDetector
    val anomalyDetector = buildAnomalyDetector(data, normalizeFunction)
    ^

scala> val anomalyDetector = (datum: Vector) => distToCentroid(normalizeFunction(datum), model) > threshold
anomalyDetector: org.apache.spark.mllib.linalg.Vector => Boolean = <function1>

scala> val anomalies = originalAndData.filter {
  |   case (original, datum) => anomalyDetector(datum)
  |   }
anomalies: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[89] at keys at <console>:75
```

But then I realized the book's code spells this out differently than the git code and takes this into account.

```
scala> val anomalies = originalAndData.filter { case (original, datum) =>
|         val normalized = normalizeFunction(datum)
|         distToCentroid(normalized, model) > threshold
|       }.keys
anomalies: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[91] at keys
```

Either way they both end up with the same output on anomalies.

```

06/03/04 12:47:33 INFO scheduler.DAGScheduler: Job 42 finished. take at console=74, took 8.83249 s
0.00, tcp, telnet, 51, 167, 19827, 0.0, 0.0, 0.1, 83, 1, 2, 91, 0.0, 0.0, 1, 1, 1, 1.00, 1.00, 0.00, 0.00, 1.00, 0.00, 20, 0.1, 8, 0.29, 0.14, 0.05, 0.25, 0.05, 0.12, 0.00, 0.00, normal.
0.00, tcp, telnet, 51, 145, 13236, 0.0, 0.0, 0.1, 31, 1, 2, 38, 0.0, 0.0, 0.0, 1, 1, 1, 1.00, 1.00, 0.00, 0.00, 1.00, 0.00, 29, 0.10, 0.28, 0.10, 0.03, 0.20, 0.07, 0.22, 0.00, 0.00, normal.
0.00, tcp, telnet, 51, 307, 2374, 0.0, 0.1, 0.0, 1, 0, 1, 0, 1, 3, 1, 0, 0, 0.0, 1, 1, 0.00, 0.00, 0.00, 0.00, 1.00, 0.00, 69, 4, 0.03, 0.04, 0.01, 0.75, 0.00, 0.00, 0.00, 0.00, normal.
11565, tcp, telnet, SF, 565, 111864, 0.0, 0.0, 0.0, 1, 217, 1, 2, 247, 0.0, 4, 0.0, 0.1, 1.00, 0.00, 0.00, 0.00, 0.00, 1.00, 0.00, 0.00, 0.3, 0.38, 0.25, 0.12, 0.08, 0.00, 0.00, 0.12, 0.33, normal.
0.00, tcp, telnet, 51, 142, 26316, 0.0, 0.0, 0.1, 31, 1, 2, 36, 0.0, 0.0, 0.0, 0.0, 1, 1, 1, 1.00, 1.00, 0.00, 0.00, 1.00, 0.00, 0.00, 94, 22, 0.22, 0.04, 0.01, 0.09, 0.19, 0.02, 0.00, 0.00, normal.
0.00, tcp, telnet, 51, 399, 26200, 0.0, 0.1, 0.1, 0, 1, 0, 1, 0, 0, 0.0, 0.0, 0.15, 0.07, 0.06, 0.00, 0.00, 1.00, 0.00, 0.00, 12, 231, 255, 1.00, 0.00, 0.00, 0.01, 0.01, 0.01, 0.01, 0.00, 0.00, normal.
0.00, tcp, telnet, 51, 2895, 14208, 0.0, 0.0, 0.1, 0.0, 0.0, 0.13, 0.0, 0.0, 0.0, 1, 1, 1, 1.00, 1.00, 0.00, 0.00, 1.00, 0.00, 0.00, 21, 2, 0.18, 0.18, 0.05, 0.88, 0.05, 0.58, 0.00, 0.00, normal.
23, tcp, telnet, SF, 104, 276, 0.0, 0.0, 0.5, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1, 1, 1, 0.00, 0.00, 0.00, 1.00, 0.00, 0.00, 1.00, 0.00, 1, 2, 1, 0.00, 0.00, 1.00, 0.00, 0.00, 0.00, 0.00, 0.00, guess passed.
0.12, tcp, telnet, SF, 246, 11938, 0.0, 0.0, 4, 1, 0, 0, 0, 2, 0.0, 0.0, 0.1, 1, 0.00, 0.00, 0.00, 0.00, 1.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, normal.
12249, tcp, telnet, SF, 3943, 44466, 0.0, 0.1, 0.1, 1, 13, 1, 0, 0, 12, 0.0, 0.0, 0.1, 1, 0.00, 0.00, 0.00, 0.00, 1.00, 0.00, 0.00, 61, 0.0, 13.0, 0.05, 0.02, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, normal.

```

The most anomalous event the book had was contained in my top 10 output, but is the 6th listing. My first listing is a telnet connection. I was curious about the connections—more specifically my top listing, so I went back to the KDD cup data for what the columns represented. Also using this site—[ita.ee.lbl.gov](http://ita.ee.lbl.gov) (Paxson, nd) to understand the "flags". It would seem that an "S1" state is for when 2 sides have established a connection (TCP handshake) but nothing further has been "seen". Though according to the listing, the destination has received 19827 bytes. Telnet connections nowadays themselves are anomalous if they are seen over the internet because they are not a secure connection—anyone listening could see exactly what transpired, in 1999 that probably would have been common to see but today it is replaced by SSH and other more secure communication methods. My guess would be that the connection is anomalous because there are multiple 100% error rates (3x) happening on this connection.

When I had shutdown my VM to give it more resources I also exported a copy of it. While I was running the computations for the anomaly detector I started up the cloned VM on another computer and ran the clustering take-4 functions from the git code—the labels with entropy. This time I ended up receiving a similar output to the book that pointed to  $k=150$  was a good choice.

```
scala> def clusteringTake4(rawData: RDD[String]): Unit = {
  val parseFunction = buildCategoricalAndLabelFunction(rawData)
  val labelsAndData = rawData.map(parseFunction)
  val normalizedLabelsAndData =
    labelsAndData.mapValues(buildNormalizationFunction(labelsAndData.values)).cache()
  (80 to 160 by 10).map(k =>
    (k, clusteringScore3(normalizedLabelsAndData, k))).toList.foreach(println)
  normalizedLabelsAndData.unpersist()
}
clusteringTake4: (rawData: org.apache.spark.rdd.RDD[String])Unit
scala> clusteringTake4(rawData)
```

## References

- superuser.com, 2014. Retrieved from <http://superuser.com/questions/841865/extracting-a-tar-gz-file-returns-this-does-not-look-like-a-tar-archive>
- UCI KDD Archive, 1999. Retrieved from <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>
- Paxson, Vern. n.d. Retrieved from <http://ita.ee.lbl.gov/html/contrib/tcp-reduce-doc.html>