

# CODE REVIEW

## Assignment A8 [A7]

*November 16, 2017*

**Code Author:** Shabbir Hussain

**Review Author:** Vineet Trivedi

**Commit Id:** 47eb2e529d9d4d408d7bb33014df750552b51e9b

### **Report:**

1. For convergence of KMeans Subtask 1 the report states that it takes 5 minutes for convergence per task in subtask 1. It would be a good idea to mention the number of iterations for convergence as well.
2. One of the requirements of the assignments is:  
“Observe whether:  
A song’s loudness, length, or tempo predict its hotness  
A song’s loudness, length, tempo, or hotness predict its combined hotness”.  
The above requirement/observation was not reported.
3. System and algorithm design for agglomerative clustering is reported.  
However, cluster and performance reporting for the following tasks in agglomerative clustering is missing:
  - a. Fuzzy Length
  - b. Fuzzy Tempo
  - c. Fuzzy hotness
  - d. Combined hotness
4. As aforementioned reporting is missing the following requirement is not completely met:  
“Use each method to perform each clustering task. Compare the results, as well as the performance of the solutions”.
5. Performance reporting for Subtask 2 is missing. Neither the number of iterations nor the time taken for clustering/convergence is reported.

6. One of the requirements of the assignment is  
"Evaluate your solution using your a local machine and AWS".  
AWS evaluation and reporting is missing for both Subtasks.

**Run:**

1. It takes 54 iterations to converge for fuzzy length on the subset and 47 iterations on the full dataset which is a bit odd as convergence generally happens around 10-20 iterations for all other tasks in Shabbir's assignment as well as all tasks in my assignment.
2. On running the makefile it is observed that agglomerative clustering is performed only for fuzzy loudness.
3. According to my understanding agglomerative clustering combines two closest clusters at every iteration. Hence the total number of clusters should decrease by 1 at every iteration. In the given implementation, the total number of clusters decreases drastically and arbitrarily at every iteration:

```
- Start time:2017/00/16 19:00:03
Num Clusters = 999056
Num Clusters = 32828
Num Clusters = 27571
Num Clusters = 20563
Num Clusters = 18720
Num Clusters = 18562
Num Clusters = 18503
Num Clusters = 17511
Num Clusters = 12984
Num Clusters = 12819
Num Clusters = 12679
Num Clusters = 11863
Num Clusters = 11348
Num Clusters = 11347
Num Clusters = 11346
Num Clusters = 11345
Num Clusters = 10980
Num Clusters = 10978
Num Clusters = 10977
Num Clusters = 10715
Num Clusters = 10712
Num Clusters = 10709
Num Clusters = 10691
Num Clusters = 10690
Num Clusters = 10689
Num Clusters = 10628
Num Clusters = 4300
Num Clusters = 4177
Num Clusters = 4176
Num Clusters = 4175
Num Clusters = 3660
Num Clusters = 3577
```

4. The 3 files can only be accepted in '.gz' format. This is because the final part of the path has been hardcoded. It would be better for the main program to accept the files in an unzipped format and leave the compressing and decompressing part to the makefile.

```
object ArtistRecord {  
  def loadCSV(sc: SparkContext, path: String): RDD[ArtistRecord] = {  
    sc.textFile(path + "artist_terms.csv.gz")  
  }  
  
object SimilarArtist {  
  def loadCSV(sc: SparkContext, path: String): RDD[SimilarArtist] = {  
    sc.textFile(path + "similar_artists.csv.gz")  
  }  
  
object SongRecord {  
  def loadCSV(sc: SparkContext, path: String): RDD[SongRecord] = {  
    sc.textFile(path + "song_info.csv.gz")  
  }  
}
```

### **Conclusion:**

Overall the assignment is well done, the report is pristine, concise and informative. The use of heat maps is brilliant. However, some parts are missing mainly performance data of each task, performance on AWS and agglomerative clustering implementation of fuzzy length, fuzzy tempo, fuzzy hotness and combined hotness.