

ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ

ΠΡΟΧΩΡΗΜΕΝΑ ΘΕΜΑΤΑ ΒΑΣΕΩΝ ΔΕΔΟΜΕΝΩΝ

Ακ. έτος 2019-2020

Εξαμηνιαία Εργασία

ΚΟΝΤΟΓΙΑΝΝΗΣ ΑΝΔΡΕΑΣ

03115187

Θέμα 3ο: Machine Learning - Ομαδοποίηση δεδομένων με εκτέλεση του KMeans αλγόριθμου

Χρησιμοποιώντας τα δεδομένα της εκφώνησης, σκοπό της μελέτης αποτελούσε η εύρεση των 5 κορυφαίων σημείων επιβίβασης με την χρήση του KMeans αλγορίθμου. Επειδή τα δεδομένα ήταν αρκετά μεγάλα, κατασκευάσαμε τον KMeans με Map-Reduce method, αφού προηγουμένως είχαμε κάνει όλα τα απαραίτητα, που μας παρείχαν οι οδηγίες του μαθημάτος, για την εγκατάσταση του hdfs και του spark. Χρησιμοποιήσαμε το Module της python2, pyspark, μία τελευταίας τεχνολογίας δωρεάν βιβλιοθήκη, που παρέχει όλες τις απαραίτητες λειτουργίες του big data, συμπεριλαμβανομένων και map-reduce υπηρεσιών.

Ας ξεκινήσουμε από τα βασικά. Τι είναι ο KMeans; Ο KMeans είναι ένας κλασικός αλγόριθμος μηχανικής μάθησης, και μάλιστα unsupervised learning, που επιτελεί την συσταδοποίηση των δεδομένων σε K συστάδες (clusters), με το K να αποτελεί υπερπαραμέτρο που δίνεται εξωτερικά. Αποτελείται από δύο στάδια: το E step και το M step. Στο E step (Expectation) βρίσκει κάθε σημείο του δείγματος σε ποια συστάδα ανήκει σύμφωνα με μία μετρική απόστασης (εδώ Haversine) από το prototype της καθεμιάς. Στο M step (Maximization) προσδιορίζει ξανά τα prototypes των συστάδων, ως τον μέσο όρο των σημείων που τις αποτελούν. Ο αλγόριθμος τερματίζει ύστερα από σύγκλιση είτε μετά από κάποιες επαναλήψεις οι οποίες έχουν ληφθεί εξωτερικά ως υπερπαραμέτρος. Σημειώνουμε πως στην αρχή έχουμε αρχικοποιήσει τα prototypes σε κάποιες τιμές πχ. στις πρώτες K εγγραφές του συνόλου δεδομένων μας.

Αν είχαμε να κάνουμε με ένα αρκετά μικρότερο dataset, τότε το task μας θα ήταν απλά να προγραμματίσουμε τον παραπάνω κλασικό αλγόριθμο και να κάνουμε fit τα δεδομένα μας.

Σημείωση: Ο κώδικας για τον κλασικό KMeans χωρίς map-reduce παρέχεται extra μέσα στο .zip στον ξεχωριστό φάκελο classic_kmeans.

Συνεχίζουμε με το map-reduce πρόβλημα. Αρχικά, βλέπουμε πως είναι το jps στον master και στο slave, για να βεβαιωθούμε ότι έχουμε στήσει σωστά το hdfs και το spark:

```
user@master:~$ jps
816 Jps
32738 NameNode
500 SecondaryNameNode
655 Master
```

```
user@slave:~$ jps
4457 Jps
25497 Worker
25358 DataNode
```

Στο σημείο αυτό, μπορούμε να σχεδιάσουμε την map-reduce υλοποίηση. Παρακάτω, παρατίθεται ο ψευδοκώδικας για τον map-reduce KMeans. Θα ξεκινήσουμε με το σπλιτάρισμα των γραμμών, την προσάρτηση των επιθυμητών συντεταγμένων και το φιλτράρισμά τους για την αποφυγή invalid συντεταγμένων.

```
1 map(key, line):
2     spl = split(line, ",") # cut ","
3     emit ( float(spl[3]), float(spl[4]) ) # return (x,y)
4
5 map(x,y):
6     if x <> 0 and y <> 0 :
7         emit (x,y) # return valid (x,y)
```

Στην συνέχεια, παραθέτουμε τον ψευδοκώδικα για τον KMeans, ξεκινώντας από το E step και καταλήγοντας στο M step του, που έχουμε περιγράψει παραπάνω.

```
1 KMeans_e_step(metric, (x,y), centroids):
2     # Returns the index of closest centroid to point (x,y) using metric distance
3     r := 0
4     i := 0
5     min := +Inf
6     for all centroid in centroids:
7         distance := metric( (x,y), centroid)
8         if distance < min:
9             min := distance
10            r := i
11            i := i + 1
```

```
1 # KMeans #
2 centroids := first k records of dataset
3 MAX_ITER := 3
4 it := 0
5 for all it: iter < MAX_ITER
6 {
7     # E Step #
8     map(x,y): # find closest cluster (r) and returns (r, (x,y), 1) in order to count cluster size
9         emit ( KMeans_e_step(haversine, (x, y), centroids), ((x,y), 1) )
10
11     # M Step #
12     reduce(r, list_of_xy_and_1):
13         emit ( r, ( sum( list_of_xy(x)(y)), sum(list_of_1) ) )
14
15     map(r, (sum_xy, sum_r)): # update centroids
16         emit ( r, ( sum_xy(x) / sum_r, sum_xy(y) / um_r ) )
17
18     store centroids
19
20     it := it + 1
21 }
```

Σημειώνουμε για το M step ότι στο reduce βρίσκουμε για κάθε centroid το άθροισμα πάνω σε όλα τα x , το άθροισμα πάνω σε όλα τα y και το άθροισμα πάνω στους άσους, ενώ στο map βρίσκουμε το μέσο όρο των δυο πρώτων αθροισμάτων με την διαίρεση τους με το τρίτο.

Τέλος, τα αποτελέσματα όταν τρέξαμε τον αλγόριθμο ήταν τα παρακάτω:

```
KMeans results
(0, (-74.0191832401856, 40.71055443942774))
(1, (-73.94113734517788, 40.76328399383515))
(2, (-73.99772401348547, 40.731550302890184))
(3, (-74.00521692085603, 40.73256456305599))
(4, (-73.98429681172031, 40.75122648858301))
```

Αξίζει να παραθέσουμε και τα logs, τόσο του master όσο και του slave:

master logs

```
user@master:~/spark-2.4.4-bin-hadoop2.7/logs$ ls
spark-user-org.apache.spark.deploy.master.Master-1-master.out  spark-user-org.apache.spark.deploy.master.Master-1-master.out.4
spark-user-org.apache.spark.deploy.master.Master-1-master.out.1  spark-user-org.apache.spark.deploy.master.Master-1-master.out.5
spark-user-org.apache.spark.deploy.master.Master-1-master.out.2  spark-user-org.apache.spark.deploy.worker.Worker-1-master.out
spark-user-org.apache.spark.deploy.master.Master-1-master.out.3
```

slave logs

```
user@slave:~/spark-2.4.4-bin-hadoop2.7/logs$ ls
spark-user-org.apache.spark.deploy.master.Master-1-slave.out  spark-user-org.apache.spark.deploy.worker.Worker-1-slave.out.3
spark-user-org.apache.spark.deploy.worker.Worker-1-slave.out  spark-user-org.apache.spark.deploy.worker.Worker-1-slave.out.4
spark-user-org.apache.spark.deploy.worker.Worker-1-slave.out.1  spark-user-org.apache.spark.deploy.worker.Worker-1-slave.out.5
spark-user-org.apache.spark.deploy.worker.Worker-1-slave.out.2
```

dataset link: <http://83.212.76.198:50070/explorer.html#/>