



CAP - Developing with Spark and Hadoop:

Homework Assignment Guide for Students

Homework: Implement an Iterative Algorithm with Spark2

Homework: Implement an Iterative Algorithm with Spark

Files and Data Used in This Homework:

Exercise Directory: `$DEV1/exercises/spark-iterative`

Data files (HDFS): `/loudacre/devicestatus_etl/*`

Stubs (professor VM only):

`KMeansCoords.pyspark`

`KMeansCoords.scalaspark`

In this Exercise, you will practice implementing iterative algorithms in Spark by calculating k-means for a set of points.

Review the Data

In the optional section of the “Use RDDs to Explore and Transform a Dataset” exercise, you used Spark to extract the date, maker, device ID, latitude and longitude from the `devicestatus.txt` data file, and store the results in the HDFS directory `/loudacre/devicestatus_etl`.

If you did not have time to complete that optional homework, run the solution script now following the steps below (professor VM only, students who need it should request the solution). (If you have run the course catch-up script, this is already done for you.)

- Copy `$DEV1DATA/devicestatus.txt` to HDFS directory `/loudacre/`
- Run the Spark script `$DEV1/exercises/spark-etl/bonus/DeviceStatusETL` (either `.pyspark` or `.scalapark` depending on which language you are using)

Examine the data in the dataset. Note that the latitude and longitude are the 4th and 5th fields, respectively, e.g.

```
2014-03-15:10:10:20,Sorrento,8cc3b47e-bd01-4482-b500-  
28f2342679af,33.6894754264,-117.543308253  
2014-03-15:10:10:20,MeeToo,ef8c7564-0a1a-4650-a655-  
c8bbd5f8f943,37.4321088904,-121.485029632
```

Calculate k-means for device location

If you are already familiar with calculating k-means, try doing the exercise on your own. Otherwise, follow the step-by-step process below.

1. Start by copying the provided `KMeansCoords` stub file (professor VM only), which contains the following convenience functions used in calculating k-means:
 - `closestPoint`: given a (latitude/longitude) point and an array of current center points, returns the index in the array of the center closest to the given point
 - `addPoints`: given two points, return a point which is the sum of the two points – that is, $(x1+x2, y1+y2)$
 - `distanceSquared`: given two points, returns the squared distance of the two. This is a common calculation required in graph analysis.
2. Set the variable `K` (the number of means to calculate). For this use `K=5`.
3. Set the variable `convergeDist`. This will be used to decide when the k-means calculation is done – when the amount the locations of the means changes between iterations is less than `convergeDist`. A “perfect” solution would be 0; this number represents a “good enough” solution. For this exercise, use a value of 0.1.
4. Parse the input file, which is delimited by the character ‘,’ into (latitude, longitude) pairs (the 4th and 5th fields in each line). Only

include known locations (that is, filter out (0,0) locations). Be sure to persist (cache) the resulting RDD because you will access it each time through the iteration.

5. Create a K-length array called `kPoints` by taking a random sample of K location points from the RDD as starting means (center points). E.g.

```
data.takeSample(False, K, 42)
```

6. Iteratively calculate a new set of K means until the total distance between the means calculated for this iteration and the last is smaller than `convergeDist`. For each iteration:

- a. For each coordinate point, use the provided `closestPoint` function to map each point to the index in the `kPoints` array of the location closest to that point. The resulting RDD should be keyed by the index, and the value should be the pair: (point, 1). (The value '1' will later be used to count the number of points closest to a given mean.) E.g.

(1, ((37.43210, -121.48502), 1))
(4, ((33.11310, -111.33201), 1))
(0, ((39.36351, -119.40003), 1))
(1, ((40.00019, -116.44829), 1))
...

- b. Reduce the result: for each center in the `kPoints` array, sum the latitudes and longitudes, respectively, of all the points closest to that center, and the number of closest points. E.g.

(0,	((2638919.87653,-8895032.182481), 74693))
(1,	((3654635.24961,-12197518.55688), 101268))
(2,	((1863384.99784,-5839621.052003), 48620))
(3,	((4887181.82600,-14674125.94873), 126114))
(4,	((2866039.85637,-9608816.13682), 81162))

- c. The reduced RDD should have (at most) `K` members. Map each to a new center point by calculating the average latitude and longitude for each set of closest points: that is, map
`(index, (totalX, totalY), n)` to `(index, (totalX/n, totalY/n))`
- d. Collect these new points into a local map or array keyed by index.
- e. Use the provided `distanceSquared` method to calculate how much each center “moved” between the current iteration and the last. That is, for each center in `kPoints`, calculate the distance between that point and the corresponding new point, and sum those distances. That is the delta between iterations; when the delta is less than `convergeDist`, stop iterating.
- f. Copy the new center points to the `kPoints` array in preparation for the next iteration.

7. When the iteration is complete, display the final `K` center points.

This is the end of the Homework