

BIG DATA COMPUTING 2021/22 - HOMEWORK 3

PYTHON VERSION

Run your algorithm on the cluster on CloudVeneto using the following datasets: **HIGGS-REDUCED-7D.txt** (about 1.2M points in 7 dimensions), and **artificial9000.txt** (9200 points in 2 dimensions). The datasets are in the **directory /data/BDC2122** of the HDFS. You must fill the two tables below, one for each dataset, where the headers of the rows indicate the values to report, and the headers of the columns indicate the configurations of parameters to be used.

The first table collects results aimed at assessing the **scalability** of the algorithm.

HIGGS-REDUCED-7D.txt	2 executors k=10, z=150, L=2	4 executors k=10, z=150, L=4	8 executors k=10, z=150, L=8	16 executors k=10, z=150, L=16
Time to read input from file (in ms)	11156.63170 8145142	7430.318355 560303	5485.337734 222412	5754.363536 834717
Time of ROUND 1 (in ms)	369891.9751 6441345	187617.8948 879242	96774.56498 146057	48003.89409 0652466
Time of ROUND 2 (in ms)	49.28803443 9086914	150.6471633 9111328	431.1201572 418213	1829.772233 9630127
Time to compute objective function (in ms)	2348.786592 4835205	1363.988161 0870361	634.9282264 709473	492.1143054 962158
Value of objective function	9.122390510 09	7.568193009 74	6.532160553 26	6.062788908 65

The second table collects results aimed at comparing the **accuracy** attained by the algorithm against the one attained by the sequential algorithm from Homework 2 on the entire dataset.

Artificial9000.txt	2 executors k=9, z=200, L=2	4 executors k=9, z=200, L=4	8 executors k=9, z=200, L=8	16 executors k=9, z=200, L=16	Sequential algorithm from Homework 2 with k=9, z=200
Value of objective function	13.24243818 18	12.85579161 31	11.84149167 12	11.71766060 27	11.57693970788481 2

Provide below a brief comment to justify the scalability and accuracy observed (your answer should be of at most 6 lines, font 12 points):

The scalability of the algorithm is shown by the fact that by increasing the number of partitions the total execution time decreases. The execution time of Round 1 halves accordingly by doubling the executors and the Round 2 increases because the sequential algorithm has to handle more points. By increasing the number of partitions the value of the objective function converges to the score of the sequential algorithm, that's because increasing the partitions leads to a better representation of the initial full pointset, which is where Sequential computes its solution.