

Analysis of G16HW1.py

Description of MRPrintStatistics()

Some general notations:

L = number of partitions; K = set of centroids; N = number of points; k_j = centroid j in K .

Round 1

- Map phase: for every $0 \leq j < N$ and $0 \leq i < |K|$ separately $(\text{coord}_j, g_j) \rightarrow (k_i, (\text{coord}_j, g_j))$
where coord_j are coordinates of point in RDD, g_j is the group of the point of index j .
- Reduce phase: for every $0 \leq p < L$ separately $(k_i, (\text{coord}_j, g_j)) \rightarrow (j \bmod L, (k_{ip}, (\text{SA}_{ip}, \text{SB}_{ip})))$
where SA_{ip} and SB_{ip} are respectively the number of points with label A and B and related to k_{ip} .

Round 2

- Map phase: -
- Reduce phase: $(k_i, S_i) \rightarrow (k_i, (\text{NA}_i, \text{NB}_i))$ where S_i is a list of pairs of values $\text{SA}_{ip}, \text{SB}_{ip}$ related to k_i .

$$\text{NA}_i = \sum_{p=0}^{L-1} \text{SA}_{ip} \text{ and } \text{NB}_i = \sum_{p=0}^{L-1} \text{SB}_{ip} .$$

Analysis of MRPrintStatistics()

Rounds = 2

Round 1

- Map phase : $O(K)$
In the first map phase we have used `flatMap()`.
`parse_coors_group()` is a conversion of points with complexity $O(1)$.
`calc_center()` computes the distance from each point to all centroids. The complexity is $O(K)$.
- Reduce phase: $O(N*K/L)$
The complexity of `mapPartitions()` is $O(N/L)$.
For each centroid, we compute the sum of points A , B in the partition and in the worst case the complexity is $O(K)$.

Round 2

- Map phase: $O(1)$
- Reduce phase: $O(L)$
In input of the last reduce phase we can have at most L tuples, one for each partition.

$$M_L = O(\max\{O(K), O(N*K/L), O(1), O(L)\}) = O(N*K/L)$$