

MapReduce: Exercises

Exercise (See ex. 2.3.1.(b) of [LRU14])

Let S be a set of N integers represented by pairs (i, x_i) , with $0 \leq i < N$, where i is the key of the pair and x_i is an arbitrary integer.

- ① Design a 2-round MR algorithm to compute the arithmetic mean of the integers in S , using $O(\sqrt{N})$ local space and $O(N)$ aggregate space.
- ② Show how to modify the algorithm of the previous point to reduce the local space bound to $O(N^{1/4})$, at the expense of an increased number of rounds.
- ③ Can you attain $M_L = O(M)$, for any $M < \sqrt{N}$?

input: pairs (i, x_i) , $i \in [0, N]$ key, $x_i \in \mathbb{R}$

output: $(0, \frac{1}{N} \sum_{i=0}^{N-1} x_i)$

① 2-round solution, $M_1 = O(\sqrt{N})$, $M_2 = O(N)$

Round 1:

ℓ (parameter)

- Map: \forall input pair $(i, x_i) \rightarrow (i \bmod \sqrt{N}, x_i)$
- Reduce: \forall key $j \in [0, \sqrt{N}]$, let L_j be the set of values with key $j \Rightarrow$ emit pair $(0, \sum_{x \in L_j} x)$ add $|L_j|$ as value

Round 2:

- Map: empty
- Reduce: let L_0 be the set of values with key 0 \Rightarrow emit $(0, \frac{1}{|L_0|} \sum_{y \in L_0} y)$

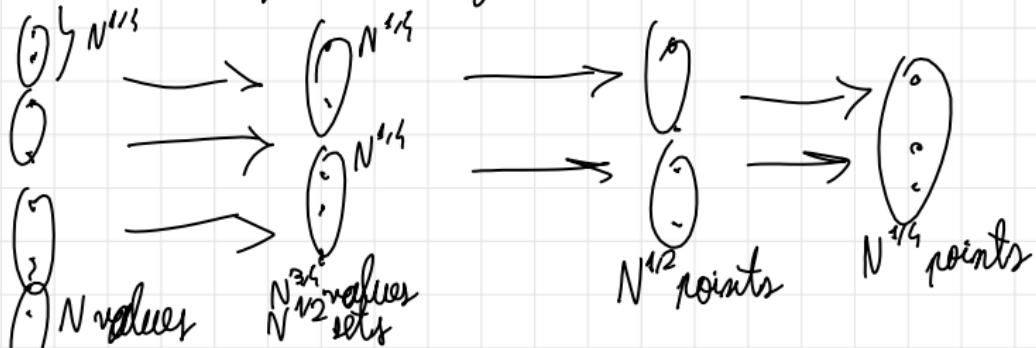
We assume N is known \Rightarrow possible in Spark

We can modify the algorithm for the case we don't know N

$$M_L = \max \{ O(1), O(\sqrt{N}), O(1), O(\sqrt{N}) \} = O(\sqrt{N})$$
$$M_A = \max \{ O(N), O(N), O(\sqrt{N}), O(\sqrt{N}) \} = O(N)$$

② goal: $M_L = O(\sqrt[4]{N})$, $M_A = O(N)$

We could use $f = N^{3/4} \Rightarrow N^{1/4}$ el. for partition \Rightarrow
 $\Rightarrow N^{3/4}$ partial sums
on top of this, repeat partitioning



ROUND 1:

- Map: $(i, x_i) \rightarrow (i \bmod N^{3/4}, x_i)$
- Reduce: $\forall \text{key } j \in [0, N^{3/4}], L_j: \text{list of values with key } j \Rightarrow$
 $\Rightarrow \text{note } |L_i| = N^{1/4}, \text{ emit } (j, \sum_{x \in L_i} x)$

ROUND 2:

- Map: $(j, y_j) \rightarrow (j \bmod N^{1/2}, y_j)$
- Reduce: $\forall \text{key } k, \dots \Rightarrow \text{emit } (k, \sum_{z \in L_k} z) \quad (k \in [0, N^{1/2}])$

ROUND 3:

- Map: $(k, z_k) \quad k \in [0, N^{1/2}] \rightarrow \text{emit } (z \bmod N^{1/4}, z_k)$
- Reduce: $\forall \text{key } h \in [0, N^{1/4}], \text{ let } L_h = \dots \Rightarrow$
 $\Rightarrow \text{emit } (0, \sum_{z \in L_h} z)$

ROUND 4:

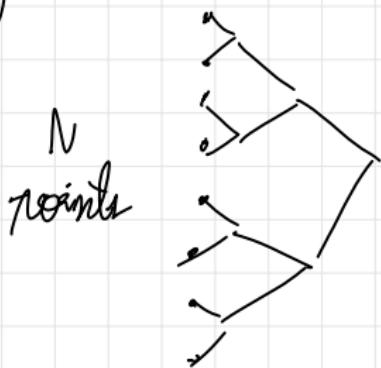
- Map: empty
- Reduce: let $L_0 = \text{list of values with key } 0 \Rightarrow \text{emit } (0, \frac{1}{N} \sum_{z \in L_0} z)$

$$M_L = \max \{ O(N^{1/4}), O(1) \} = O(N^{1/4})$$

all Reduce all Map

$$M_R = \max \{ O(N), O(N^{3/4}), O(N^{1/2}), O(N^{1/4}) \} = O(N)$$

③



in general, $\log_2 N$ reduces,
 $M_L = O(1)$ by running 2 el. per time

running M el. per time:

$$N \rightarrow \frac{N}{M} \rightarrow \frac{N}{M^2} = \dots \rightarrow \frac{N}{M^i} \rightarrow \dots \rightarrow M \rightarrow$$

$$\Rightarrow i = \frac{\log_2 N/M}{\log_2 M} + 1 = \frac{\log N}{\log M} = \log_M N$$

Exercise (See ex. 2.3.1.(d) of [LRU14])

Let S be a set of N integers represented by pairs (i, x_i) , with $0 \leq i < N$, where i is the key of the pair and x_i is an arbitrary integer. We want to design an efficient MR algorithm to compute the number D of distinct integers in S .

- ① Design a simple 2-round MR algorithm to compute D and analyze its local and aggregate space requirements. Under what circumstances is the local space proportional to N , thus not meeting the design goals of MR algorithms?
- ② Design a better MR algorithm to compute D in $O(1)$ rounds, using $O(\sqrt{N})$ local space and $O(N)$ aggregate space.

① problem: $O(N)$ impossible

ROUND 1:

- Map: $(i, x_i) \rightarrow (x_i, 1)$
- Reduce: $\forall \text{key } x, \text{ let } L_x: \text{values with key } x \Rightarrow \text{emit } (0, x_i)$

ROUND 2:

- Map: empty
- Reduce: let $L_0 = \{\text{list of distinct elements}\} \Rightarrow \text{emit } (0, |L_0|)$

There could be one very frequent element or all different elements

② part 1: remove duplicates \rightarrow 2 rounds
part 2: count number of distinct elements $\Rightarrow //$

PART 1:

ROUND 1:

- Map: $(i, x_i) \rightarrow (i \bmod \sqrt{N}, x_i)$
- Reduce: $\forall \text{key } j \in [0, \sqrt{N}]$ let L_j = list of values with key j
 \forall distinct value x in L_j , emit (x, j)

ROUND 2:

- Map: empty
- Reduce: \forall integer x , $L_x = \{\text{list of values in } [0, \sqrt{N}]\} \Rightarrow$
 \Rightarrow emit (x, j) , $j \in L_x$ (take whatever from L_x)

PART 2:

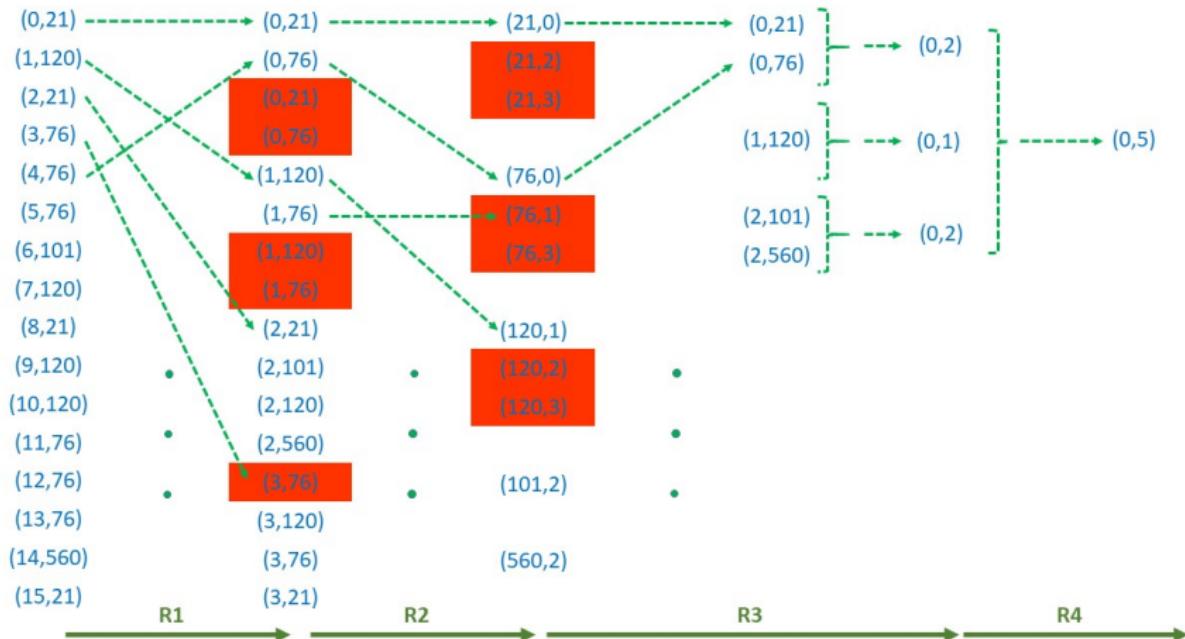
ROUND 1:

- Map: $(x, j) \rightarrow (j, x)$ (~~cannot have $> \sqrt{N}$ different keys~~)
- Reduce: $\forall j \in [0, \sqrt{N}-1], L_j = \dots \Rightarrow \text{emit } (0, \sum_{i=0}^j x_i)$

ROUND 2:

- Map: empty
- Reduce: $L_0 = \{\text{partial counts}\} \Rightarrow \text{emit } (0, \sum_{x \in L_0} x)$

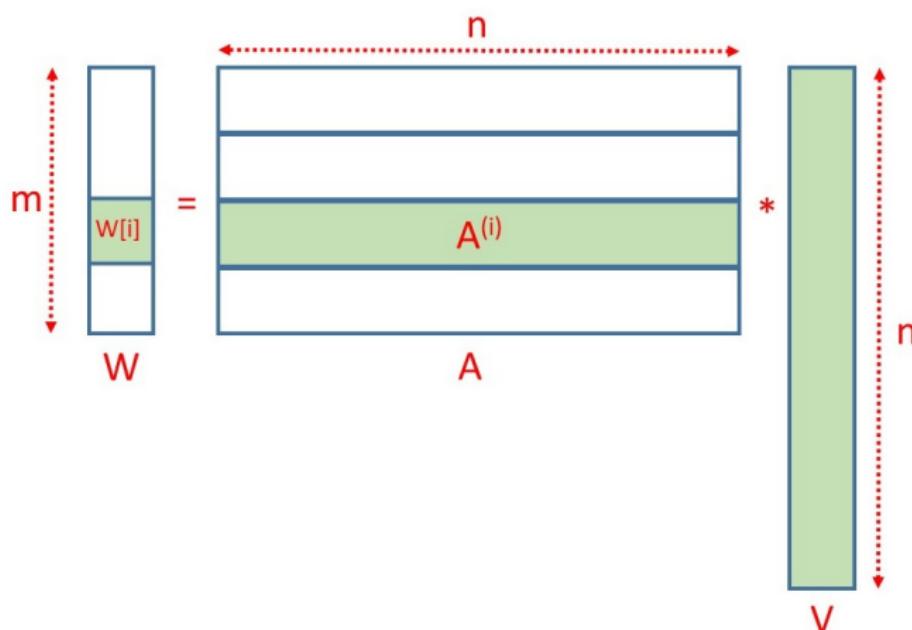
Example ($N = 16$)



Exercise

Design an $O(1)$ -round MR algorithm for computing a matrix-vector product $W = A \cdot V$, where A is an $m \times n$ matrix, V is an n -vector, and $m \leq \sqrt{n}$. Your algorithm must use $o(n)$ local space and linear (i.e., $O(mn)$) aggregate space.

How would your solution change if m were larger?



Exercise

Design a 2-round MR algorithm for the word count problem, using random partitioning, and analyze its local and aggregate space requirements.

