

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. Carefully specify the input and output pairs of each round. You may assume that N is a global variable known to the algorithm.
- 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 generalize the previous point to attain $O(N^\epsilon)$, for any $\epsilon \in (0, 1/2)$?

① INPUT : $\{(i, x_i) : 0 \leq i < N, x_i \text{ integer}\}$
OUTPUT : $(0, \text{mean}(S) = (1/N) \cdot \sum_{i=0}^{N-1} x_i)$

Round 1

Map Phase : $\forall 0 \leq i < N \quad (i, x_i) \rightarrow (i \bmod \sqrt{N}, x_i)$

Reduce Phase : for each $0 \leq j < \sqrt{N}$ separately, let

$L_j = \text{list of integers } x \text{ from intermediate pairs } (j, x)$

$(j, L_j) \rightarrow (0, \text{sum}_j = \sum_{x \in L_j} x)$

Round 2

Map Phase : empty

Reduce Phase : let $L_0 = \text{sum}_0, \text{sum}_1, \dots, \text{sum}_{\sqrt{N}-1}$ be
the list of pairs produced by R1

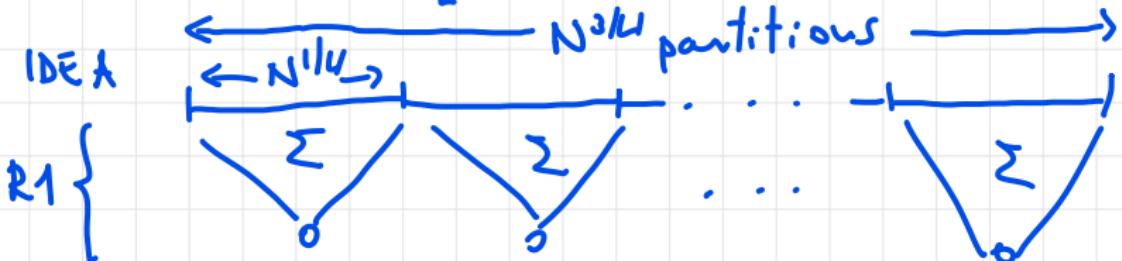
$$(0, L_0) \rightarrow (0, \text{mean}(S) = \frac{1}{N} \sum_{\substack{\text{sum}_j \in L_0 \\ \text{Global variable}}} \text{sum}_j)$$

Obs. If we want to avoid possible overflows due to very large numbers, we can compute the mean locally in each partition in R1, together with the size of the partition, and then, in R2, compute the "weighted" mean of the \sqrt{N} local means

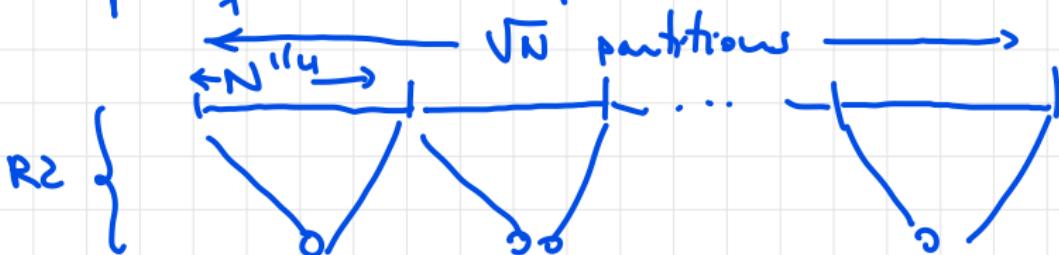
$M_L = O(\sqrt{N})$ since in the reduce phases of the 2 rounds,
each reducer adds up to \sqrt{N} integers

$M_A = O(N)$ since the number of input, output,
and intermediate pairs in each
round is $O(N)$

(2) OBJECTIVE : $M_L = O(N^{1/4})$



Output of R1 : $N^{3/4}$ partial sums



Output of R2 : $N^{1/2}$ partial sums



output of $R3$: $N^{1/4}$ partial sums

$R4$: Compute mean (S) by aggregating
the $N^{1/4}$ partial sums produced by $R3$

Round 1

Map Phase: $\forall 0 \leq i < N : (i, x_i) \rightarrow (i \bmod N^{3/4}, x_i)$

Reduce Phase: $\forall 0 \leq j < N^{3/4}$ separately let

L_j = list of integers x from intermediate pairs
 (j, x)

$$(j, L_j) \rightarrow (j, \sum_{x \in L_j} x)$$

Round 2

Map Phase : $\nabla_{0 \leq j < N^{3/4}} (j, s) \rightarrow (j \bmod N^{1/2}, s)$

Reduce Phase : $\nabla_{0 \leq j < N^{1/2}}$ separately, let

L_j = list of integers s in intermediate pairs (j, s)

$$(j, L_j) \rightarrow (j, \sum_{s \in L_j} s)$$

Round 3 }
Round 4 }

Exercise

Analysis

$$R = 4$$

M_L = for each round

* Map Phase: $O(1)$

* Reduce Phase: $O(N^{1/4})$

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

$M_A = O(N)$ since the number of input,
output and intermediate pieces in
each round is $O(N)$

$$\textcircled{3} \quad N \rightarrow N^{1-\varepsilon} \rightarrow N^{1-2\varepsilon} \dots N^{1-i\varepsilon} = N^{\varepsilon} \quad \begin{matrix} i+1 \text{ round} \\ i+1 = \frac{1}{\varepsilon} \end{matrix}$$

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 INPUT: $\{(i, x_i) : 0 \leq i < N, x_i \text{ integer}\}$
OUTPUT $(0, \text{count})$ count = # distinct integers x_i 's

IDEA

R1: eliminate duplicates

R2: Count how many integers (all distinct)
are left

Round 1

Map Phase: $\forall 0 \leq i < N \quad (i, x_i) \rightarrow (x_i, i)$

Reduce Phase: $\forall \text{integer } x \text{ separately let}$
 $L_x = \text{list of } i\text{'s from intermediate pairs } (x_i, i)$
with $x_i = x$

$$(x, L_x) \rightarrow (0, x)$$

Round 2

Map phase: empty

Reduce phase : let l_0 be the list of integers from pairs $(0, x)$ produced by R1

$$(0, l_0) \rightarrow (0, |l_0|) \xrightarrow{\# \text{ distinct integers}}$$

Analysis

$M_A = O(N)$ since the number of input, output and intermediate pairs in each round is $O(N)$

$M_L =$

* N_1 = Max # of copies per integer

* N_2 = # distinct integers

$$M_L = O(\max(N_1, N_2)) = O(N_1 + N_2)$$

and both N_1 and N_2 can be proportional to N in the worst case.

② OBJECTIVE $M_L = O(\sqrt{N})$

IDEA: * Duplicate removals : 2 rounds

* Count distinct : 2 rounds

Round 1

Map Phase: $\forall 0 \leq i < N \quad (i, x_i) \rightarrow (i \bmod \sqrt{N}, x_i)$

Reduce Phase: $\forall 0 \leq j < \sqrt{N}$ separately, let

$L_j = \text{list of integers } x \text{ from intermediate pairs } (j, x)$
 $(j, L_j) \rightarrow \{(j, x) : \# \text{ distinct integer } x \in L_j\}$

AT THIS POINT: there are $\leq \sqrt{N}$ copies of
each integer

Round 2

Map Phase: ∇ pair (j, x) from R1 $(j, x) \rightarrow (x, j)$

Reduce Phase: ∇ integer x , let

$L_x = \text{list of indices } j \text{ of intermediate pairs } (x, j)$
 $(x, L_x) \rightarrow (x, j)$ with j an arbitrary index
from L_x

AT THIS POINT

- * Only 1 copy of each integer is left
- * $\leq \sqrt{N}$ pairs with index j

Round 3

Map phase $\forall \text{par}(x, j) : (x, j) \rightarrow (j, x)$

Reduce phase $\forall 0 \leq j < \sqrt{N}$ let

L_j : list of integers from intermediate pairs (j, x)

$(j, L_j) \rightarrow (0, c_j = |L_j|)$

Round 4

Map Phase: empty

Reduce Phase: let $L_0 = c_0 c_1 \dots \in \mathbb{C}^{\sqrt{N}-1}$

$$(0, L_0) \rightarrow (0, \sum_{c_j \in L_0} c_j)$$

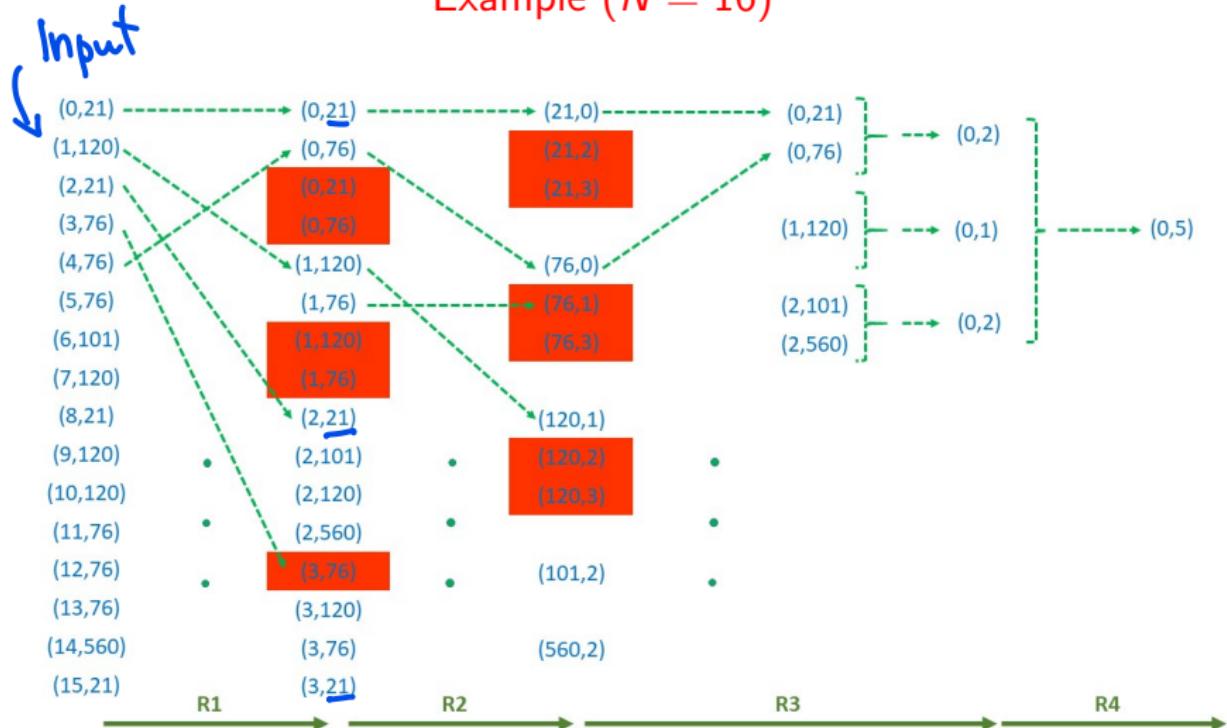
→ # distinct integers
in the input set

Analysis

$M_A = O(N)$ usual argument

$M_L = \Theta(\sqrt{N})$ because of the initial partition and of the fact that after $R_1 \leq \sqrt{N}$ copies of each integer were left

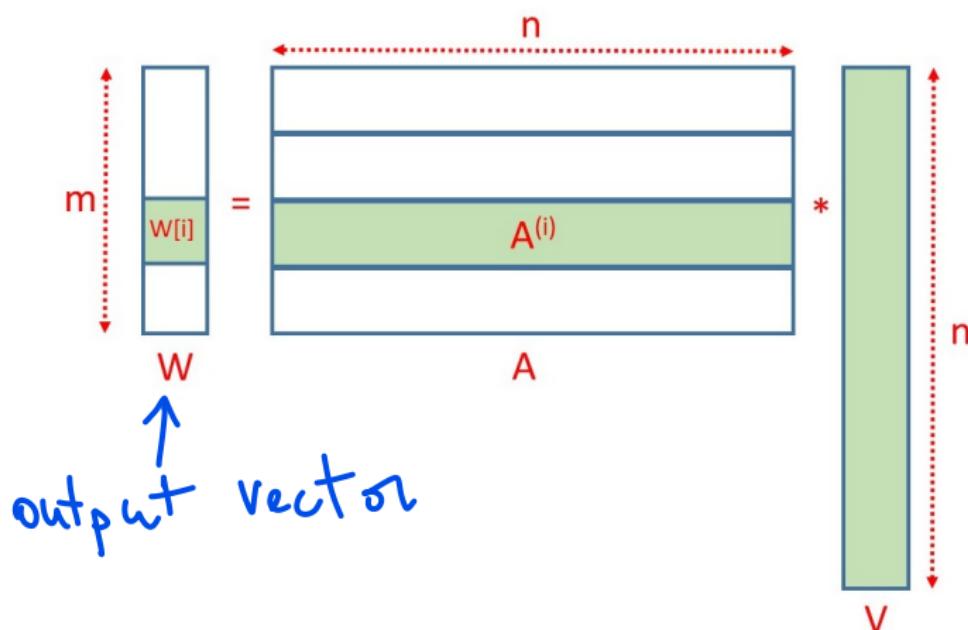
Example ($N = 16$)



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?



INPUT $A = \{((i, j), A[i, j]) : 0 \leq i < m \quad 0 \leq j < n\}$

$V = \{((j, -1), V[j]) : 0 \leq j < n\}$

OUTPUT $W = \{(i, W[i]) : 0 \leq i < m\}$

1-round algorithm (space inefficient)

2-round algorithm (space efficient)

$W[i] = A^{(i)} \circ V$ ↗ inner product

1-round algorithm

Map Phase

$$((i, j), A[i, j]) \rightarrow (i, ((i, j), A[i, j]))$$

$$((j, -1), V[j]) \rightarrow \{(i, ((j, -1), V[j])) : 0 \leq i < m\}$$

$\Rightarrow m$ copies of V

Notation: $V^{(i)} \equiv$ copy i of V $0 \leq i < m$

$$W[i] = A^{(i)} \circ V^{(i)}$$

\uparrow \nwarrow

i-th row of A i-th copy of V

Reduce Phase $\forall 0 \leq i < m$ separately, let

L_i = list of values from intermediate pairs

with key $i \Rightarrow L_i$ contains the representation
of $A^{(i)}$ and $V^{(i)}$

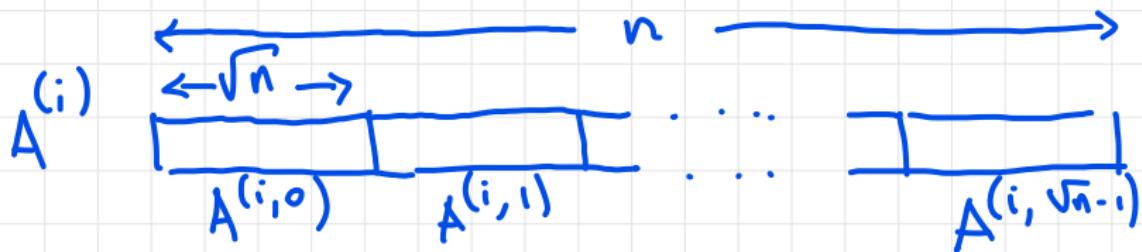
$$(i, L_i) \rightarrow (i, W[i]) = A^{(i)} \circ V^{(i)}$$

$$M_A = O(n \cdot m)$$

$$M_L = O(m+n) = O(n)$$

Map Phase Reduce Phase

OBJECTIVE : $M_L = o(n)$ in particular $M_L = O(\sqrt{n})$



$$\Rightarrow |A^{(i,s)}| = |V^{(i,s)}| = \sqrt{n} \quad \forall 0 \leq s < \sqrt{n}$$

$$\text{Define } W_s[i] = A^{(i,s)} \circ V^{(i,s)}$$

$$\Rightarrow W[i] = \sum_{s=0}^{\sqrt{n}-1} W_s[i]$$

Round 1

Map Phase

- * Replicate each entry of V m times (as before)
- * Assign $\text{key}(i, s)$ to each entry of segment s of $A^{(i)}$ (i.e., $A^{(i,s)}$) and to each entry of segments of the i -th copy of V (i.e., $V^{(i,s)}$)

Reduce Phase $\forall 0 \leq i < m$ and $\forall 0 \leq s < \sqrt{n}$

Compute $W_s[i] = A^{(i,s)} \circ V^{(i,s)}$
returning $(i, W_s[i])$

Round 2

Map Phase : empty

Reduce Phase :

$$\text{Compute } (i, W[i]) = \sum_{s=0}^{\sqrt{n}-1} W_s[i]$$

These are the values of
the pairs produced by R1
with key i

Exercise: fill in the details

Analyze s_{rs}

$$M_A = O(m \cdot n) \quad \text{as before}$$

M_L :

* Map Phase of R_1 : $O(m)$

To m copies
of V

* Reduce Phase of R_1 : $O(\sqrt{n})$ ← to compute the
inner products of

* Map Phase of R_2 : -

segments of length
 \sqrt{n}

* Reduce Phase of R_2 : $O(\sqrt{n})$

$$\Rightarrow M_L = O(\max\{m, \sqrt{n}\}) = O(\sqrt{n}) \quad \text{by the hypothesis}$$

$$m \leq \sqrt{n}$$

