## 24.03.2025 Week 6 exercises: MapReduce

Exercise 1:

You are given a symmetric social network (like Facebook) where $a$ is a friend of $b$ implies that $b$ is also a friend of $a$. The input is a dataset $D$ (sharded) containing such pairs of identifiers $(a, b)$—no assumption is made regarding the order between $a$ and $b$. Pairs appear exactly once and are not repeated. Find the last names of those users whose first name is "Kanye" and who have at least 300 friends. You can chain MapReduces if you want (but only if you must, and even then, the shortest possible chain).

You don't need to write code—pseudocode is fine as long as it is understandable. Your pseudocode may assume the presence of appropriate primitives (e.g., "`firstname(user_id)`", etc.). The Map function takes as input a tuple (`key = a, value = b`).

```
1: procedure MAP(a: User, b: User)
2:     if firstname(a) is "Kanye" then
3:         emit(a, b)
4:     if firstname(b) is "Kanye" then
5:         emit(b, a)
6: procedure REDUCE(u: User, friends: List[User])
7:     if length(friends) ≥ 300 then
8:         emit(lastname(u))
```

Exercise 2:

For an asymmetrical social network, you are given a dataset $D$ where lines consist of $(a, b)$ which means user $a$ follows user $b$. Write a MapReduce program (Map and Reduce separately) that outputs the list of all users U who satisfy the following three conditions simultaneously: i) user $U$ has at least 2 million followers, and ii) $U$ follows fewer than 20 other users, and iii) all the users that $U$ follows, also follow $U$ back.

```
1: procedure MAP(a: User, b: User)
2:     emit(a, ⟨b, 1⟩)                              ▷ a follows b
3:     emit(b, ⟨a, 0⟩)                              ▷ b is followed by a
4: procedure REDUCE(u: User, list: List[⟨User, Int⟩])
5:     follows ← ∅
6:     count₀ ← 0
7:     count₁ ← 0
8:     for all pair in list do                      ▷ we assume pair has key and value attributes
9:         if pair.value = 0 then
10:            count₀ ← count₀ + 1                   ▷ count of followers
11:        if pair.value = 1 then
12:            count₁ ← count₁ + 1                   ▷ count of follows
13:            follows ← follows ∪ pair.key          ▷ set of follows
14:    if count₀ ≥ 2M and count₁ < 20 then
15:        for all user in follows do
16:            if ⟨user, 0⟩ ∉ list then
17:                return
18:        emit(u)
```

**Exercise 3** (Final Exam 2021)**:**

Matrix multiplication is a fundamental operation in machine learning. Assume we now want to perform matrix multiplication over two extremely large matrices $A$ and $B$.

Design an $O(1)$-round Map–Reduce program for computing a matrix–matrix product $M = AB$, where $A \in \mathbb{R}^{m \times n}$ is an m x n matrix and $B \in \mathbb{R}^{n \times m}$.

Matrices $A$ and $B$ are represented through $mn$ pairs, and each pair corresponds to $(A, i, j, A[i, j])$ or $(B, i, j, B[i, j])$ (the first three entries are string while the last one is a floating point). The Map function reads each available pair as the input.

Recall how to perform matrix multiplication:

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix} \times \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix} = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} & a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} & a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32} \end{pmatrix}$$

1: **procedure** MAP($\langle X, i, j, X_{i,j} \rangle$: $\langle$String, String, String, Float$\rangle$)                    ▷ Matrix element
2:     **if** $X$ = "A" **then**
3:         **for** $k$ in $\{0, 1, 2, ..., m\text{-}1\}$ **do**
4:             emit($\langle i, k \rangle, \langle X, j, X_{ij} \rangle$)                    ▷ Computation of $M[i, k]$ needs $A[i, j]$
5:     **if** $X$ = "B" **then**
6:         **for** $k$ in $\{0, 1, 2, ..., m\text{-}1\}$ **do**
7:             emit($\langle k, j \rangle, \langle X, i, X_{ij} \rangle$)                    ▷ Computation of $M[k, j]$ needs $B[i, j]$
8: **procedure** REDUCE($\langle i, j \rangle$: $\langle$String, String$\rangle$, list: List[$\langle$String, String, Float$\rangle$])
9:     Avec $\leftarrow \emptyset$
10:     Bvec $\leftarrow \emptyset$
11:     $M_{i,j} \leftarrow 0$
12:     **for all** $\langle X, k, X_k \rangle$ in list **do**                    ▷ Separate out matrix elements
13:         **if** $X$ = "A" **then**
14:             Avec $\leftarrow$ Avec $\cup \langle k, X_k \rangle$
15:         **if** $X$ = "B" **then**
16:             Bvec $\leftarrow$ Bvec $\cup \langle k, X_k \rangle$
17:     sort(Avec)
18:     sort(Bvec)
19:     **for** $k$ in $\{0, 1, 2, ..., n\text{-}1\}$ **do**
20:         $M_{i,j} \leftarrow M_{i,j} + \text{Avec}[k] * \text{Bvec}[k]$                    ▷ Multiply corresponding elements
21:     emit($\langle M, i, j, M_{i,j} \rangle$)

**Exercise 4** ([10 points] Final Exam 2022 – multiple choice: Only 1 correct answer. No negative marking.)**:**

1) For each of the following questions, compute the total communication cost between the mappers and the reducers, i.e., the total number of (key, value) pairs outputted by the mappers. Use the optimal algorithm that runs in one MapReduce step for each of the tasks. Also, assume that there is no combiner.

(a) [2 points] Word count for a dataset comprising $W$ total words with $d$ distinct words. The mappers receive a single word as input.

( ) $d$

( ) $W$

    The algorithm will emit one key-value pair per word.

( ) $\frac{W}{d}$

( ) $dW$

( ) $2W$

```
1: procedure MAP(word: String)
2:     emit(word, 1)
3: procedure REDUCE(word: String, counts: List[Integer])
4:     total ← 0
5:     for all count in counts do
6:         total ← total + count
7:     emit(word, total)
```

(b) [3 points] Matrix multiplication of two matrices of size $m \times n$ and $n \times p$. The mappers read input tuples in the form `<Matrix identifier, row index, column index, value>`.

( ) $mp$

( ) $n(m + p)$

( ) $2n(m + p)$

( ) $mnp$

( ) $2mnp$

    The final matrix will be $m \times p$. Each element $(i, j)$ will require $n$ elements of the first matrix and $n$ elements of the second matrix.

```
1: procedure MAP(MatrixID: String, i: Integer, j: Integer, value: Float)
2:     if MatrixID = "A" then
3:         for k ← 1 to p do
4:             emit((i, k), (A, j, value))
5:     else
6:         for k ← 1 to m do
7:             emit((k, j), (B, i, value))
```

(c1) [2 points] Computing the INNER JOIN of two relations defined as follows: relation `R1(X, Y)` with four tuples $\{(5,21), (7,16), (15,3), (3,21)\}$ and relation `R2(Y, Z)` with $\{(3,1), (4,8), (21,28)\}$. The mappers read input tuples in the form `<Relation identifier, column-1 value, column-2 value>`.

*Recall: The database INNER JOIN operation of two relations outputs concatenated tuples from R1 and R2 having equal values on the common column.*

( ) 12

( ) 7

The comparison has to be done in the reducer. The mapper will emit all the rows of both the relations with a tag.

( ) 4

( ) 3

( ) 2

(c2) [1 point] How many output tuples are produced by the reducer(s) in question (c1) ?

( ) 12

( ) 7

( ) 4

( ) 3

(5, 21, 28), (15, 3, 1), (3, 21, 28)

( ) 2

```
1: procedure MAP(RelationID: String, X: Integer, Y: Integer)
2:     if RelationID = "R1" then
3:         emit(Y, (RelationID, X))
4:     else
5:         emit(X, (RelationID, Y))
6: procedure REDUCE(Key: Integer, Values: List[Tuple])
7:     R1_list ← []
8:     R2_list ← []
9:     for all (RelationID, Value) in Values do
10:        if RelationID = "R1" then
11:            append R1_list with Value
12:        else
13:            append R2_list with Value
14:    for all v1 in R1_list do
15:        for all v2 in R2_list do
16:            emit((v1,Key,v2))
```

(d) [2 points] Difference of two sets $X$ and $Y$ containing a total of $x$ and $y$ elements respectively. The mappers read input tuples in the form <Set identifier, Value>.

*Recall: The difference of two sets X and Y is a set that contains those elements of X that are NOT in Y.*

( ) $x + y$

    All the values have to be emitted once from both the sets.

( ) $x - y$

( ) $x$

( ) $y$

( ) $xy$

```
1: procedure MAP(SetId: String, Value: Integer)
2:     emit(Value, SetId)
3: procedure REDUCE(Value: Integer, SetIds: List[String])
4:     if "X" in SetIds and "Y" not in SetIds then
5:         emit(Value)
```