

Homework #4

(Due: May 6)

Task 1. [160 Points] Sequence Alignment with General Gap Penalty.

The problem of *sequence alignment with general gap penalty (gap problem)* is a generalization of the edit distance problem that arises in molecular biology, geology, and speech recognition. When transforming a string $X = x_1x_2 \dots x_m$ into another string $Y = y_1y_2 \dots y_n$, a sequence of consecutive deletes corresponds to a gap in X , and a sequence of consecutive inserts corresponds to a gap in Y . An affine gap penalty function is predominantly used in bioinformatics. However, in many applications the cost of such a gap is not necessarily equal to the sum of the costs of each individual deletion (or insertion) in that gap. In order to handle this general case, two new cost functions w_{ins} and w_{del} are defined, where $w_{del}(p, q)$ ($0 \leq p < q \leq m$) is the cost of deleting $x_{p+1} \dots x_q$ from X , and $w_{ins}(p, q)$ ($0 \leq p < q \leq n$) is the cost of inserting $y_{p+1} \dots y_q$ into X . The substitution function $s(x_i, y_j)$ is the same as that of the standard edit distance problem. Let $G[i, j]$ denote the minimum cost of transforming $X_i = x_1x_2 \dots x_i$ into $Y_j = y_1y_2 \dots y_j$ (where $0 \leq i \leq m$ and $0 \leq j \leq n$) under this general setting. Then $G[0, 0] = 0$, $G[0, j] = w_{del}(0, j)$ for $1 \leq j \leq n$, and $G[i, 0] = w_{ins}(0, i)$ for $1 \leq i \leq m$. Otherwise,

$$G[i, j] = \min \left\{ \begin{array}{l} G[i-1, j-1] + s(x_i, y_j), \\ \min_{0 \leq q < j} \{ G[i, q] + w_{del}(q, j) \}, \\ \min_{0 \leq p < i} \{ G[p, j] + w_{ins}(p, i) \} \end{array} \right\}. \quad (1)$$

Figures 2 and 3 show two different approaches for evaluating the entries of $G[1 : m, 1 : n]$. While the pull-based approach shown in Figure 2 is the standard way of evaluating Recurrence 1, the non-standard push-based approach shown in Figure 3 seems more suitable for a MapReduce implementation.

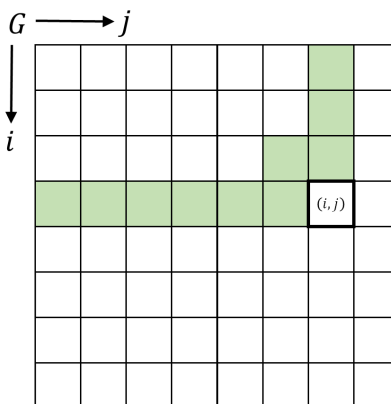


Figure 1: Dependency structure of the gap DP: cell $G[i, j]$ depends on the green cells.

LOOP-GAP-PULL(G, x, m, y, n) {Inputs are sequences $X = x_1x_2 \dots x_m$ and $Y = y_1y_2 \dots y_n$, and an $(m+1) \times (n+1)$ matrix $G[0 : m, 0 : n]$. Row 0 and column 0 of G are assumed to be appropriately initialized.}

1. **for** $t \leftarrow 1$ **to** $m + n - 1$ **do**
2. **for** $i \leftarrow \max\{1, t + 1 - n\}$ **to** $\min\{t, m\}$ **do**
3. $j \leftarrow t + 1 - i$
4. $G[i, j] \leftarrow G[i - 1, j - 1] + s(x_i, y_j)$
5. **for** $q \leftarrow 0$ **to** $j - 1$ **do**
6. $G[i, j] \leftarrow \min\{G[i, j], G[i, q] + w_{ins}(q, j)\}$
7. **for** $p \leftarrow 0$ **to** $i - 1$ **do**
8. $G[i, j] \leftarrow \min\{G[i, j], G[p, j] + w_{ins}(p, i)\}$
9. **return** $G[m, n]$

Figure 2: [PULL-BASED APPROACH] Looping code implementing Recurrence 1 in which each cell $G[i, j]$ pulls values from the cells it depends on.

LOOP-GAP-PUSH(G, x, m, y, n) {Inputs are sequences $X = x_1x_2 \dots x_m$ and $Y = y_1y_2 \dots y_n$, and an $(m+1) \times (n+1)$ matrix $G[0 : m, 0 : n]$. Row 0 and column 0 of G are assumed to be appropriately initialized.}

1. **for** $i \leftarrow 1$ **to** m **do**
2. **for** $j \leftarrow 1$ **to** n **do**
3. $G[i, j] \leftarrow \min\{G[i, 0] + w_{del}(0, j), G[0, j] + w_{ins}(0, i)\}$
4. **for** $i \leftarrow 1$ **to** m **do**
5. $G[i, 1] \leftarrow \min\{G[i, 1], G[i - 1, 0] + s(x_i, y_1)\}$
6. **for** $j \leftarrow 1$ **to** n **do**
7. $G[1, j] \leftarrow \min\{G[1, j], G[0, j - 1] + s(x_1, y_j)\}$
8. **for** $t \leftarrow 1$ **to** $m + n - 1$ **do**
9. **for** $i \leftarrow \max\{1, t + 1 - n\}$ **to** $\min\{t, m\}$ **do**
10. $j \leftarrow t + 1 - i$
11. **if** $i < m$ **and** $j < n$ **then** $G[i + 1, j + 1] \leftarrow \min\{G[i + 1, j + 1], G[i, j] + s(x_{i+1}, y_{j+1})\}$
12. **for** $q \leftarrow j + 1$ **to** n **do**
13. $G[i, q] \leftarrow \min\{G[i, q], G[i, j] + w_{ins}(j, q)\}$
14. **for** $p \leftarrow i + 1$ **to** m **do**
15. $G[p, j] \leftarrow \min\{G[p, j], G[i, j] + w_{ins}(i, p)\}$
16. **return** $G[m, n]$

Figure 3: [PUSH-BASED APPROACH] Looping code implementing Recurrence 1 in which each cell $G[i, j]$ pushes its finalized value to all cells that depend on it. This approach is better suited for MapReduce.

In the rest of the task we will assume $m = n$ for simplicity. We will also assume $w_{ins}(i, j) = j^2 - i^2$, $w_{del}(i, j) = j^3 - i^3$, and $s(x_i, y_j) = 0$ when $x_i = y_j$ and 1 otherwise.

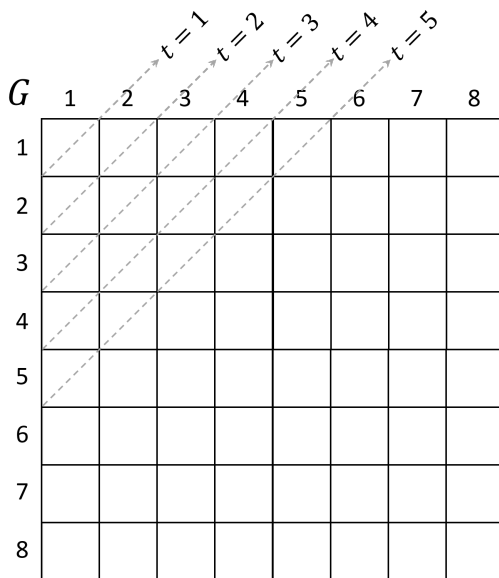


Figure 4: The dashed arrows show the cells that push their finalized values to other cells for various values of t in the approach given in Figure 3. For each value of t all cells to the right and below of the corresponding dashed arrow get updated due to the pushes.

- (a) [**80 Points**] Design a MapReduce algorithm for solving Recurrence 1 based on the push-based approach given in Figure 3. Write pseudocode, and analyze its time and space complexities. Modify your algorithm to reduce space complexity.
- (b) [**80 Points**] Implement your algorithm from part (a) on Hadoop. As a practical performance optimization you may want to avoid executing the last few values of t as MapReduce rounds because the number of emits will be very low. Similarly for the first few values of t . Find the smallest value of n for which you are unable to allocate the entire matrix G on a single machine. Let n' be that value. Now use your Hadoop implementation for solving the problem for $n = n'$ and $n = 2n'$. Report running times using 10 data nodes.

Task 2. [90 Points] Floyd-Warshall's All-Pairs Shortest Path (APSP).

This task is an extension of Task 1 of Homework 1.

Consider a directed graph $\mathcal{G} = (V, E)$, where $V = \{v_1, v_2, \dots, v_n\}$, and each edge (v_i, v_j) is labeled by a real-valued length $l(v_i, v_j)$. If $(v_i, v_j) \notin E$, $l(v_i, v_j)$ is assumed to have a value $+\infty$. The length of a path is defined as the sum of the lengths of the edges in the path, taken in order. For each pair $v_i, v_j \in V$, $\delta[i, j]$ is defined to be the smallest of the lengths of all paths going from v_i to v_j . We assume that G does not have any negative length cycle.

Floyd-Warshall's algorithm for computing the length of the shortest path between every pair of vertices of G has a standard iterative implementation as shown in Figure 5. The input is an $n \times n$ matrix $D[1 \dots n, 1 \dots n]$ in which for every pair $i, j \in [1, n]$, $D[i, j]$ is initialized with $l(v_i, v_j)$. After the algorithm terminates $D[i, j] = \delta[i, j]$ for all $i, j \in [1, n]$.

```

LOOP-FLOYD-WARSHALL-APSP-PULL( $D, n$ )
1. for  $k \leftarrow 1$  to  $n$  do
2.   for  $i \leftarrow 1$  to  $n$  do
3.     for  $j \leftarrow 1$  to  $n$  do
4.        $D[i, j] \leftarrow \min(D[i, j], D[i, k] + D[k, j])$ 

```

Figure 5: [PULL-BASED APPROACH] Looping code implementing Floyd-Warshall's APSP algorithm.

- (a) [**90 Points**] Design a MapReduce algorithm for Floyd-Warshall's APSP shown in Figure 5. Write pseudocode, and analyze its time and space complexities. Modify your algorithm to reduce space complexity.

APPENDIX 1: What to Turn in

One compressed archive file (e.g., zip, tar.gz) containing the following items.

- Source code, makefiles and job scripts.
- A PDF document containing all answers and plots.