

1. Compute the edit distance between “Paris” and “alice”

EDITDISTANCE( $s_1, s_2$ )

```

1  int  $m[i, j] = 0$ 
2  for  $i \leftarrow 1$  to  $|s_1|$ 
3  do  $m[i, 0] = i$ 
4  for  $j \leftarrow 1$  to  $|s_2|$ 
5  do  $m[0, j] = j$ 
6  for  $i \leftarrow 1$  to  $|s_1|$ 
7  do for  $j \leftarrow 1$  to  $|s_2|$ 
8      do  $m[i, j] = \min\{m[i-1, j-1] + \text{if } (s_1[i] = s_2[j]) \text{ then } 0 \text{ else } 1, \text{ fi},$ 
9           $m[i-1, j] + 1,$ 
10          $m[i, j-1] + 1\}$ 
11  return  $m[|s_1|, |s_2|]$ 

```

**SOLUTION.**

		a	l	i	c	e
	0	1	2	3	4	5
p	1	2	3	4	5	6
a	2	3	2	3	4	5
r	3	4	3	4	5	6
i	4	5	4	3	4	5
s	5	6	5	4	3	4

2. Write down the entries in permutation index dictionary that are generated by the term “mama”

**SOLUTION.**

mama\$, ama\$m, ma\$ma, a\$mam, \$mama.

3. Apply map reduce to problem of counting how often each term occurs in a set of files. Specify map and reduce operations for this task. Write down an example along the lines

4.

If we need  $n \log_2 n$  comparisons (where  $n$  is the number of termID-docID pairs) and 2 disk seeks for each comparison, how much time would index construction for Reuters-RCV1 take if we used disk instead of memory for storage and an unoptimized sorting algorithm (i.e., not an external sorting algorithm)? Use the system parameters in Table 4.1.

### SOLUTION.

4.1 An unoptimized sorting algorithm would be to have all postings stored on disk and transfer them back and forth from disk to memory to make comparisons.

A trivial Index construction task would consist of these 2 steps:

1. Parsing the documents and creating the postings(n)
2. Sorting the postings

Step 1 would take  $O(n)$  time

Step 2 takes  $2 \cdot (n \cdot \log_2 n) \cdot \text{disk seek time}$ . For RCV-1 corpus,  $n = 100 \text{ million} = 10^8$

In this case, Step 2 dominates the total time by a large factor.

Total time  $= 2 \cdot (10^8 \cdot \log_2 10^8) \cdot 5 \cdot 10^{-5} \text{ s} =$

5. Can the tf-idf weight of a term in a document exceed 1?

Yes  $\text{idf}_t = \log \frac{N}{\text{df}_t}$ .

6.

One measure of the similarity of two vectors is the *Euclidean distance* (or  $L_2$  distance) between them:

$$|\vec{x} - \vec{y}| = \sqrt{\sum_{i=1}^M (x_i - y_i)^2}$$

Given a query  $q$  and documents  $d_1, d_2, \dots$ , we may rank the documents  $d_i$  in order of increasing Euclidean distance from  $q$ . Show that if  $q$  and the  $d_i$  are all normalized to unit vectors, then the rank ordering produced by Euclidean distance is identical to that produced by cosine similarities.

word	query					document			
	tf	wf	idf	$q_i = \text{wf-idf}$		tf	wf	$d_i = \text{normalized wf}$	$q_i \cdot d_i$
digital			10,000						
video			100,000						
cameras			50,000						

► **Table 4** Cosine computation for Exercise 6.19.

**SOLUTION.**

$$\sum (q_i - w_i)^2 = \sum q_i^2 - 2\sum q_i w_i + \sum w_i^2 = 2(1 - \sum q_i w_i)$$

$$\text{Thus: } \sum (q_i - v_i)^2 < \sum (q_i - w_i)^2 \Leftrightarrow 2(1 - \sum q_i v_i) < 2(1 - \sum q_i w_i) \Leftrightarrow \sum q_i v_i > \sum q_i w_i$$

**Exercise 0 91**