

Lecture 1: Sorting and General Computational Problems

Harvard SEAS - Fall 2021

Sept. 2, 2021

1 Announcements

- Put up a name tent with your name and an emoji
- Watch course overview video if you haven't already done so
- Staff introductions
- My OH today: 1-2pm in the SEC 3.327; zoom option available from my webpage
- TF sections and OH 0
- Revised problem set 0 posted
- Join/follow Ed even during class. We'll use it for quizzes and continuous chat/Q&A.
- Collaboration policies
-

2 Recommended Reading

- CS50 Week 3: <https://cs50.harvard.edu/x/2021/weeks/3/>
- Cormen-Leiserson-Rivest-Stein Section 1.3
- Roughgarden I, Sections 1.4 and 1.5
- We've ordered all of these at the Coop and for reserve through Harvard libraries (check on status through HOLLIS); apparently CLRS is already available to read as an e-book.

3 Motivating Problem: Web Search

Simplified and outdated description Google's original search algorithm:

1. (Calculate Pageranks) For every URL url on the entire world-wide web WWW, calculate its *pagerank*, $PR_{url} \in [0, 1]$.
2. (Keyword Search) Given a search keyword w , let S_w be the set of all webpages containing w . That is, $S_k = \{url \in WWW : w \text{ is contained on the webpage at } url\}$.
3. (Sort Results) Return the list of URLs in S_w , sorted in decreasing order of their pagerank.

The definition and calculation of pageranks (Step 1) was the biggest innovation in Google's search, and is the most computationally intensive of these steps. However, it can be done offline, with periodic updates, rather than needing to be done in real-time response to an individual search query. Pageranks are outside the scope of CS 120, but you can learn more about them in courses like CS 222 (Algorithms at the End of the Wire) and CS 229r (Spectral Graph Theory in Computer Science).

The keyword search (Step 2) can be done by creating a *trie* data structure for each webpage, also offline. Covered in CS50.

Our focus here is Sorting (Step 3), which needs to be extremely fast (unlike `my.harvard!`) in response to real-time queries, and operates on a massive scale (e.g. millions of pages).

4 The Sorting Problem

Input : An array A of item-key pairs $((I_0, K_0), \dots, (I_{n-1}, K_{n-1}))$, where each key $K_i \in \mathbb{R}$

Output : An array A' of item-key pairs $((I'_0, K'_0), \dots, (I'_{n-1}, K'_{n-1}))$ that is a *valid sorting* of A . That is, A' should be:

1. sorted by key values, i.e. $K'_0 \leq K'_1 \leq \dots \leq K'_{n-1}$. and
2. a permutation of A , i.e. \exists a permutation $\pi : [n] \rightarrow [n]$ such that $(I'_i, K'_i) = (I_{\pi(i)}, K_{\pi(i)})$ for $i = 0, \dots, n - 1$.

Computational Problem Sorting

Above and throughout the course, $[n]$ denotes the set of numbers $\{0, \dots, n - 1\}$. In combinatorics, it is more standard for $[n]$ to be the set $\{1, \dots, n\}$, but being computer scientists, we like to index starting at 0. Similarly, for us, the natural numbers are $\mathbb{N} = \{0, 1, 2, \dots\}$

- Application to web search:
 - Items = urls
 - Keys = $1 - \text{PR}$ (Note that we flip the pageranks so higher PR appears towards the start of the list)
- Many other applications! Database systems (both Relational and NoSQL), Machine learning systems, Ranking professional surfers by points accumulated,...
- Is the output uniquely defined?

In the subsequent sections, we will see pseudocode for three different sorting algorithms, and compare those algorithms to each other.

5 Exhaustive-Search Sort

Input : An array $A = ((I_0, K_0), \dots, (I_{n-1}, K_{n-1}))$, where each $K_i \in \mathbb{R}$
Output : A valid sorting of A

```

1 foreach permutation  $\pi : [n] \rightarrow [n]$  do
2   | if  $K_{\pi(0)} \leq K_{\pi(2)} \leq \dots \leq K_{\pi(n-1)}$  then
3   |   | return  $(I_{\pi(0)}, K_{\pi(0)}), (I_{\pi(1)}, K_{\pi(1)}), \dots, (I_{\pi(n-1)}, K_{\pi(n-1)})$ 
4   |
```

Algorithm 1: Exhaustive-Search Sort

Example: $A = ((a, 6), (b, 1), (c, 6), (d, 9))$.

As our algorithm runs, we try the following permutations and see if they produce a sorted array.

$\pi(0)$	$\pi(1)$	$\pi(2)$	$\pi(3)$
0	1	2	3
0	1	3	2
0	2	1	3
0	2	3	1
0	3	1	2
0	3	2	1
1	0	2	3
1	0	3	2

Permutation is valid if in increasing order (i.e., keys in order). So the first valid permutation is 1 0 2 3, hence the output is $(b, 1), (a, 6), (c, 6), (d, 9)$. Note that the valid sorting is not unique. E.g., in this case, $(a, 6)$ and $(c, 6)$ can be permuted.

For correctness, we check two things:

1. On the iteration where we consider permutation π , we produce an output if and only if π defines a valid sorting of the input array.
2. For every input array, there is at least one permutation π that defines a valid sorting.

Together these imply that Exhaustive-Search Sort will always produce an output, and it will always be a valid sorting of the input array.

6 Insertion Sort

Input : An array $A = ((I_0, K_0), \dots, (I_{n-1}, K_{n-1}))$, where each $K_i \in \mathbb{R}$
Output : A valid sorting of A

```

1 /* "in-place" sorting algorithm that modifies A until it is sorted */
2 foreach  $i = 1, \dots, n - 1$  do
3   | /* loop invariant:  $(A[0], A[1], \dots, A[i - 1])$  is a valid sorting of the first
4   |   |  $i$  elements of the original input array */
4   |   | Insert  $A[i]$  into the correct place in  $(A[0], \dots, A[i - 1])$ 
5 return  $A$ 
```

Algorithm 2: Insertion Sort

Example: $A = ((a, 6), (b, 2), (c, 1), (d, 4))$.

As our algorithm runs, we produce the following sorted sub-arrays:

Iteration	Sorted Sub-Array
i=0	((a,6))
i=1	((b,2),(a,6))
i=2	((c,1),(b,2),(a,6))
i=3	((c,1),(b,2),(d,4),(a,6))

Correctness: proof by induction on i that prior to the i -th loop iteration,

$$A[0], A[1], \dots, A[i-1]$$

is a valid sorting of the first i elements of the original array.

7 Merge Sort

```

1 MergeSort(A)
  Input   : An array  $A = ((I_0, K_0), \dots, (I_{n-1}, K_{n-1}))$ , where each  $K_i \in \mathbb{R}$ 
  Output  : A valid sorting of  $A$ 
2 if  $n \leq 1$  then return  $A$ ;
3 else if  $n = 2$  and  $K_0 \leq K_1$  then return  $A$ ;
4 else if  $n = 2$  and  $K_0 > K_1$  then return  $((I_1, K_1), (I_0, K_0))$ ;
5 else
6    $i = \lceil n/2 \rceil$ 
7    $A_1 = \text{MergeSort}(((I_0, K_0), \dots, (I_i, K_i)))$ 
8    $A_2 = \text{MergeSort}(((I_{i+1}, K_{i+1}), \dots, (I_{n-1}, K_{n-1})))$ 
9   return Merge ( $L_1, L_2$ )
10
```

Algorithm 3: Merge Sort

We omit the implementation of `Merge`, which you can find in the readings.

Example: $A = (7, 4, 6, 9, 7, 1, 2, 4)$.

We sort $(7, 4, 6, 9)$ and $(7, 1, 2, 4)$ independently and obtain $(4, 6, 7, 9)$ and $(1, 2, 4, 7)$. We then merge the two sorted halves and obtain $(1, 2, 4, 4, 6, 7, 9)$.

For the proof of correctness, we use strong induction. We show the base case (that we sort length 1, 2 arrays correctly), and that if we sort arrays of size up to $n - 1$ correctly, we also sort arrays of size up to n correctly.

8 Computational Problems

Definition 8.1. A *computational problem* Π is a pair (\mathcal{I}, f) where:

- \mathcal{I} is a (typically infinite) set of possible inputs x .
- For every input $x \in \mathcal{I}$, a set $f(x)$ of valid solutions.

Example: sorting

- $\mathcal{I} = \{\text{All arrays of item-key pairs with keys in } \mathbb{R}\}$
- $f(x) = \{\text{All valid sorts of } x\}$

(Note that there are multiple valid solutions, which is why $f(x)$ is a set)

Informal Definition 8.2. An *algorithm* is a well-defined “procedure” A for “transforming” any input x into an output $A(x)$.

Note: this is an informal definition. We will be more formal in a couple of weeks.

Definition 8.3. Algorithm A *solves* computational problem $\Pi = (\mathcal{I}, f)$ if for every input $x \in \mathcal{I}$, we have $A(x) \in f(x)$.

Note that we want a *single* algorithm A (with a fixed, finite description) that is going to correctly solve the problem Π for *all of the* (infinitely many) inputs in the set \mathcal{I} .

Important point: we distinguish between computational problems and algorithms that solve them. A single computational problem may have many different algorithms that solve it (or even no algorithm that solves it!), and our focus will be on trying to identify the most efficient among these.