

Data Structure And Analysis

(CSC263)

Lecture Notes

Tianquan Di

February 11, 2019

How many shortest-length paths are there to get from your house to the doughnut shop?

4 up's
7 right's

$\binom{11}{7} = \binom{11}{4} = 330$ paths

$\binom{n}{k} = \frac{n!}{(n-k)!k!}$

$e^{i\pi} + 1 = 0$

Find $7 + 12 + 17 + 22 + \dots + 342$

$S_n = 7 + 12 + 17 + 22 + \dots + 342$
 $\uparrow S_n = 342 + 337 + 332 + 327 + \dots + 7$
 $2S_n = 349 + 349 + 349 + 349 + \dots + 349$
 $2S_n = 349 \times 68$
 $S_n = \frac{349 \times 68}{2}$
 $S_n = 11866$

Original:
 $\exists x \forall y (x \geq 2y \rightarrow x > y + 1)$
 Converse:
 $\exists x \forall y (x > y + 1 \rightarrow x \geq 2y)$
 Negation:
 $\neg [\exists x \forall y (\neg (x \geq 2y) \vee x > y + 1)]$
 $\forall x \exists y (x \geq 2y \wedge x \leq y + 1)$
 Contrapositive:
 $\exists x \forall y (x \leq y + 1 \rightarrow x < 2y)$

$v - e + f = 2$

P.I.E. Example:
 $6! - \left[\binom{6}{2}5! - \binom{6}{2}4! + \binom{6}{3}3! - \binom{6}{4}2! + \binom{6}{5}1! \right]$

There are six dogs to give 13 tacos. Use a 'stars and bars' diagram to illustrate the first and sixth dog get 3 tacos, the second dog gets none, the third dog gets 5 and the fourth dog gets one.

||**|*||***|

$A = \{2, 4, 10, \text{dog}\}$

Onto

One-to-One

$(A \cup B \cup C) \cup (A \cap B \cap C)$

$K_{3,3}$

Abstract

Why do we need to learn data structures? This course helps you to explore a new perspective on the algorithm analysis, in the field of Computer Science. CSC263 will also provide you the foundation of technical interviews from your dream tech companies, such as Microsoft, Google, etc. This lecture notes will primarily follow Prof. Sam Toueg's lecture, Prof. David Liu's lecture notes, as well as CLRS. According to Prof. Sam Toueg's description of the course:

Data structures are ways of organizing the data involved in computation, suitable for representation in and manipulation by computers. Algorithms are precisely stated, general problem solving methods. Data structures and algorithms are central to computer science. They are also integrally related: neither can be studied fruitfully without knowledge of the other. This course has two goals: First, to learn several important data structures and algorithms, including how to choose and/or modify a data structure to solve various problems; and second, to introduce the basic tools and techniques for the analysis of algorithms and data structures.

Course Website:

<http://www.cs.toronto.edu/~sam/teaching/263/>

Textbook: Introduction to Algorithms (i.e. CLRS)

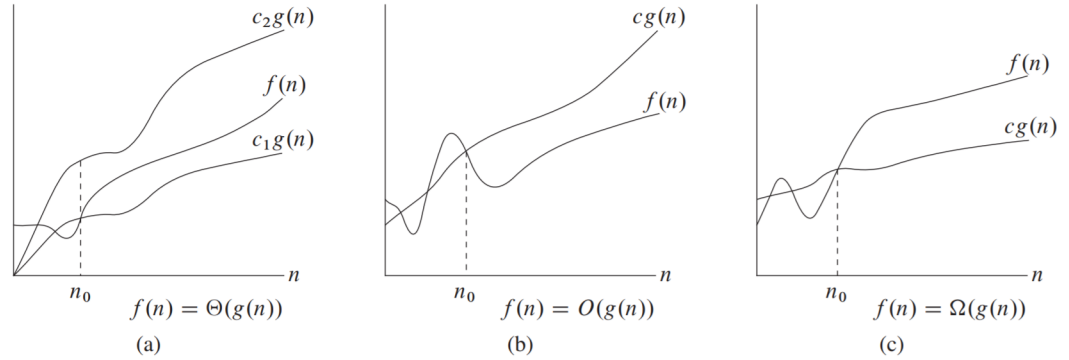
This document assumes that you have the knowledge provided from CSC165 and CSC236.

Contents

1	Review of Growth of Functions	4
1.1	Asymptotic bounds on worst-case time complexity	4
2	Priority Queues & Heaps	5
2.1	Abstract Data Type (ADT) vs. Data Structure	5
2.2	The Priority Queue ADT	5
2.3	Heaps	5
2.4	Mergeable Priority Queue ADT and Binomial Heaps	8
2.5	Implementation of operations on Binomial Heaps	9
3	Dictionary I: AVL Trees	10
3.1	Abstract Data Type: Dictionary	10
3.2	AVL Trees: Intro	10
3.3	AVL Trees: Insert Operation	11

1 Review of Growth of Functions

1.1 Asymptotic bounds on worst-case time complexity



Let $t(x)$ be the number of steps taken by algorithm \mathcal{A} on the input x .

Let $T(n)$ be the *worst-case* time complexity of algorithm \mathcal{A} . Then:

$$T(n) = \max_{\text{all inputs } x \text{ of size } n} t(x) = \max\{t(x) : x \text{ is an input of size } n\}$$

To prove that $T(n)$ is $O(g(n))$: we must show that:

$$\exists c > 0, \exists n_0 > 0, \forall n > n_0 : T(n) \leq c \cdot g(n)$$

The above statement is equivalent to:

$$\Leftrightarrow \max\{t(x) : x \text{ is an input of size } n\} \leq c \cdot g(n)$$

$$\Leftrightarrow \text{For **every** input } x \text{ of size } n, t(x) \leq c \cdot g(n)$$

$$\Leftrightarrow \text{For **every** input of size } n, \mathcal{A} \text{ takes **at most** } c \cdot g(n) \text{ steps.}$$

To prove that $T(n)$ is $\Omega(g(n))$: we must show that:

$$\exists c > 0, \exists n_0 > 0, \forall n > n_0 : T(n) \geq c \cdot g(n)$$

The above statement is equivalent to:

$$\Leftrightarrow \max\{t(x) : x \text{ is an input of size } n\} \geq c \cdot g(n)$$

$$\Leftrightarrow \text{For **some** input } x \text{ of size } n, t(x) \geq c \cdot g(n)$$

$$\Leftrightarrow \text{For **some** input of size } n, \mathcal{A} \text{ takes **at least** } c \cdot g(n) \text{ steps.}$$

To prove that $T(n)$ is $\Theta(g(n))$: we must show that:

$$T(n) \text{ is } O(g(n)) \text{ and } T(n) \text{ is } \Omega(g(n))$$

For example, the worst-case time complexity for bubble sort is $\Theta(n^2)$. This means that bubble sort is both $O(n^2)$ and $\Omega(n^2)$. Since bubble sort is $O(n^2)$, for **every** input of size n , the algorithm takes **at most** $c_1 \cdot n^2$ steps. Since bubble sort is $\Omega(n^2)$, this means that for **some** input of size n , the algorithm takes **at least** $c_2 \cdot n^2$ steps.

2 Priority Queues & Heaps

2.1 Abstract Data Type (ADT) vs. Data Structure

Before we start talking about any else, let us discuss about the definition of **abstract data type** and **data structure**.

An **abstract data type (ADT)** is a theoretical model of an entity and the set of operations that can be performed on that entity. In other words, it describes an object and its operations.

A **data structure** is some specific implementation of Abstract Data Type.

2.2 The Priority Queue ADT

Consider the following example, which is the definition for the Priority Queue ADT.

The Priority Queue ADT

Object: Set S of elements with “keys” (priority) that can be compared

Operations:

- **Insert(S, x):** Insert element x into S .
- **Max(S):** Returns the element with the highest priority (i.e. largest key) in S .
- **Extract_Max(S, x):** Returns $\text{Max}(S)$ and removes it from S .

A classic example of priority queues in practice is a hospital waiting room: more severe injuries and illnesses are generally treated before minor ones, regardless of when the patients arrived at the hospital.¹

2.3 Heaps

Our goal is to do all the operations above in $\Theta(\log(n))$ time. The solution is to use a **data structure** called **Max-Heap**.

Worst Case Time For	Insert	Extract_Max
Unordered Linked List	$\Theta(1)$	$\Theta(n)$
Ordered Linked List	$\Theta(n)$	$\Theta(1)$
Max Heap	$\Theta(\log(n))$	$\Theta(\log(n))$

Before introducing the max-heap, let us talk about the complete binary tree. Definition of a complete binary tree **complete binary tree (CBT)** is a binary tree such that, for every level L (except the bottom level), the binary tree has 2^L nodes. Also, all the nodes at the bottom level are as far left as possible.

¹Liu, n.d., 23

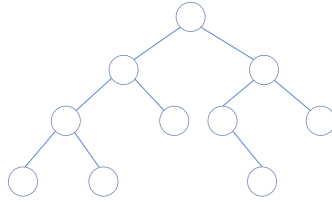


Figure 1: An example of **non-complete** binary tree

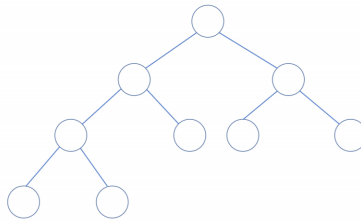


Figure 2: An example of **complete** binary tree

Fact: Height of a CBT with n nodes is $\lfloor \log_2(n) \rfloor$

A **max-heap** is a complete binary tree that satisfies **max-heap property**, which is that the value in each internal node is greater than or equal to the values in the children of that node. We typically use an **array** to store a max-heap into the memory.

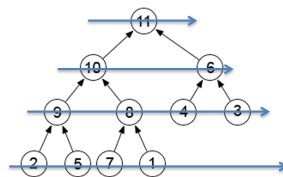


Figure 3: The order of elements in an array implementation of a max-heap.

Notice from the figure that the array is just the **level-order traversal** of the max-heap, from **top to bottom**. Now let us introduce the operations for the max-heap.

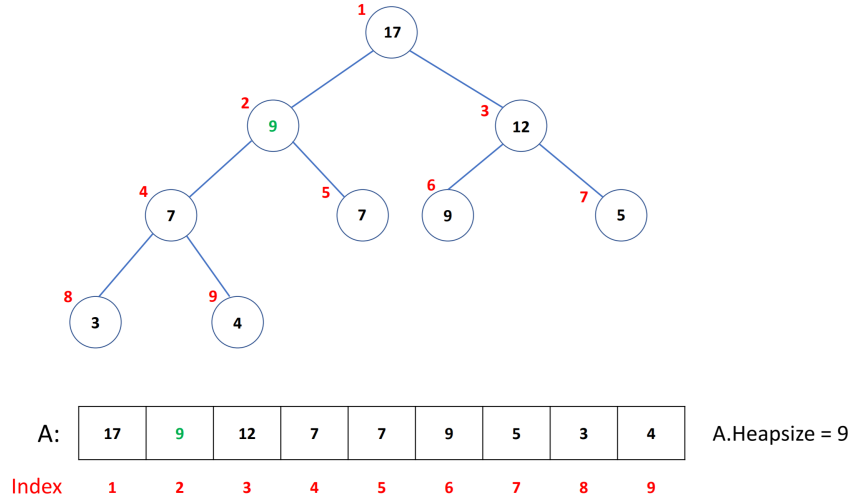


Figure 4: A typical example of array representation of a Max-Heap.

High level idea for operations on Max-Heap

- Maintain the Complete Binary Tree shape.
- Maintain the Max-Heap property.

Insert(S, x): Insert places a key in a new node that is the last node in a level-order-traversal of the heap. The inserted key is then “bubbled” upwards until the heap property is satisfied.

Algorithm 1 The algorithm for insert operation of max-heap.

procedure INSERT(A, x)

 Put x at the bottom left of the tree:

 Increment A.heapsize and set $A[A.heapsize] = x$

 Bubble x up to the top of the tree:

while x is not root **AND** priority of $x >$ priority of its parent **do**

 Swap x with its parent

Putting x at the bottom left of the tree **maintains Complete Binary Tree shape**, and bubbling x up to the top of the tree **maintains the Max-Heap property**. Since the max-heap is a CBT with height $\lfloor \log_2(n) \rfloor$, bubbling x up the tree takes $O(\log(n))$. Thus, the worst case running time of the algorithm is $O(\log(n))$.

Max(S): We simply get the root of the tree. Thus, the Max algorithm is constant runtime.

Extract_Max(S, x):

Algorithm 2 The algorithm for extract-max operation of max-heap.

procedure EXTRACT_MAX(A, x)

 Return the root $A[1]$.

 Remove the returned element from the heap:

 Set $A[1] = A[A.\text{heapsize}]$ and decrement $A.\text{heapsize}$

 Drip the element in $A[1]$ down the tree:

 Let x be the element in $A[1]$

while x is not a leaf **AND** priority of some child of x > priority of x **do**

 Swap x with the highest-priority of child of x

2.4 Mergeable Priority Queue ADT and Binomial Heaps

Abstract Data Types	Data Structures	Insert	Min	Extract_Min	Union
Priority Queues	Min Heap	Yes	Yes	Yes	No
Mergable Priority Queues	Min Binomial Heaps	Yes	Yes	Yes	Yes

Elements of binomial heaps are stored in a **sequence** of **binomial trees**.
(i.e. Binomial Forests)

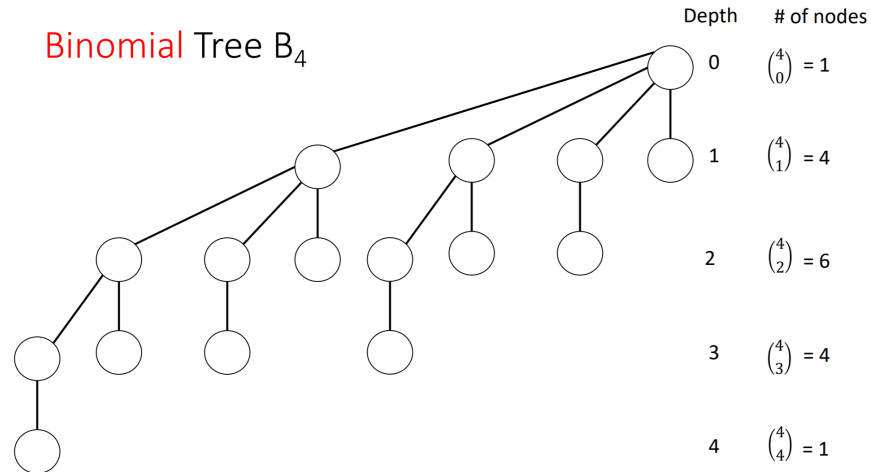


Figure 5: A visualization of binomial tree, B_4 .

Properties of Binomial Tree, B_k :

- Height k
- 2^k nodes

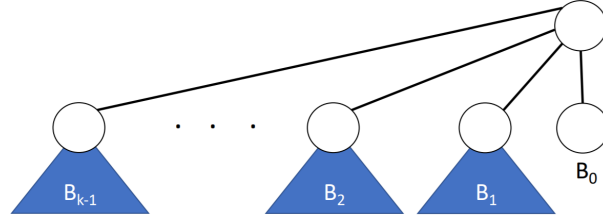


Figure 6: A more general visualization of binomial tree, B_k .

- $\binom{k}{d}$ nodes at depth d

A **Binomial Forest** is a sequence of B_k trees with strictly decreasing k 's and a total of n nodes.

Binomial Forest F_n :

- F_n has n nodes.
- $n = \langle b_t, b_{t-1}, \dots, b_0 \rangle_2$. Note that $t = \lfloor \log_2(n) \rfloor$
- $F_n = \langle \text{all trees } B_i \text{ such that bit } b_i = 1 \rangle$

2.5 Implementation of operations on Binomial Heaps

Union(T, Q): $S \leftarrow \text{Union}(T, Q)$: Think about binary addition. Additional info about this to be appeared later.

To analyze the worst-case time complexity of Union: Let $|T| \leq n$ and $|Q| \leq n$ (i.e. each contains at most n elements). This implies that each of T, Q have $O(\log(n))$ B_k trees. Then $\text{Union}(T, Q)$ takes at most $O(\log(n))$ key-comparisons. Hence the worst-case time complexity of the algorithm is $O(\log(n))$.

Insert(T, x): Simply apply this statement: $S \leftarrow \text{Union}(T, x)$

Min(T): Scan the roots of the B_k trees of T and return the smallest key.

Extract_Min(T): The idea of Extract_Min can be expressed as follow.

Algorithm 3 The algorithm for Extract_Min for binomial min heap.

```
procedure EXTRACT_MIN( $T$ )
  // Do Min( $T$ ) to locate the smallest element
   $U \leftarrow T - B_i$ 
  // Delete root of  $B_i$ . By lemma 2, we get a binomial heap  $S$ , where:
   $S \leftarrow B_i - \text{root of } B_i$ 
   $T \leftarrow \text{Union}(U, S)$ 
```

3 Dictionary I: AVL Trees

3.1 Abstract Data Type: Dictionary

Object: Set S of (element with) keys

Operations:

- Search(S, x): Returns element with key x , if x is in S ; else return “Not Found”
- Insert(S, x): Inserts x into S .
- Delete(S, x): Deletes x from S .

We want all the operations to be done in log time complexity. That is $\Theta(\log n)$. Possible solution: Binary Search Tree(BST).

3.2 AVL Trees: Intro

Use BST as the data structure for the dictionary ADT can be slow. For example, if we have a balanced BST and we keep inserting node to the rightmost corner. Then the search would take $\Theta(n)$ time.

Balanced BST: We want BST trees with height $\Theta(\log n)$. So the operations, including search, insert, delete, will have log time complexity. Solution is to use **balanced BSTs**.

AVL Trees: AVL tree is a **type** of self-balancing BST. For AVL Trees, the **balance factor** of each node is always **between -1 and 1**.

Balance Factor(BF) of a node v in AVL Trees:

$$\text{BF}(v) = \text{height}(\text{right subtree of } v) - \text{height}(\text{left subtree of } v)$$

Properties of AVL Trees:

- AVL trees of n nodes has height $\Theta(\log n)$; height $\leq 1.44 \log_2(n + 2)$
- Can do inserts, deletes while maintaining the tree balance in $\Theta(\log n)$ time.

3.3 AVL Trees: Insert Operation

General Idea:

- Insert x into T as in any BST
- Go up from x to the root:
 - For each node:
 - * Adjust the Balance Factor
 - * “Rebalance” if $BF > 1$ or $BF < -1$

Specific Implementation:

- Insert x into T as in any BST:
 - x is now a leaf
 - Set $BF(x)$ to 0 (Since it has no children, balance factor is of course 0.)
- Go up from x to the root and for each node v in this path:
 - Adjust the BF:
 - bif x is in right subtree of v : Increment $BF(v)$ if x is in left subtree of v : Decrement $BF(v)$ if $BF(v) = 0$: Stop
 - * Rebalance if necessary:
 - * if $BF(v) = +2$:
 - if $BF(v.\text{right}) = +1$: Do **Left Rotation**, update BFs of rotated nodes, and **stop**.
 - if $BF(v.\text{right}) = -1$: Do **Right-Left Rotation**, update BFs of rotated nodes, and **stop**.
 - * if $BF(v) = -2$:
 - Symmetric to above case.

Exercise: How to implement Delete and Search?

3.4 Augmentation of AVL Trees

Coming Soon...

4 Randomized Algorithms

You want an algorithm that haves well for **every possible input**.

4.1 Randomized QuickSort Algorithm (RQS)

Input: Set S of n distinct keys

Output: Keys of S in increasing order.
It is a recursive divide-and-conquer algorithm.

High-level Description:

- If S is empty then return.
- If $|S| = 1$ then output the key in S .
- Select a key p (“Pivot”) uniformly at random from S :
 - By comparing p to every other key in S , split S into:
 - * $S_< = \{s \in S | s < p\}$
 - * $S_> = \{s \in S | s > p\}$
- $RQS(S_<)$; output p ; $RQS(S_>)$

Basic Properties of RQS

- Two keys are compared if and only if one of them is selected as a pivot.
- Two keys are compared at most once.
- If two keys are “split apart” in different sets by a pivot (like 2 and 7 are split apart by pivot 4) then they are never compared.

Another Example: Fix an input S . run $RQS(S)$. Let C = number of key comparisons done by $RQS(S)$. In this case, C is a **random variable**. Thus, it has an expected value. **What is the expected value of C , $E(C)$, over all possible random pivot selections done by RQS?**

- Let $z_1 < z_2 < \dots < z_i < \dots < z_j < \dots < z_n$, where $i < j$, each z is a key.
- Let's say there is a random variable c_{ij} . $c_{ij} = \begin{cases} 1, & \text{iff } z_i \text{ and } z_j \text{ are compared by RQS(S)} \\ 0, & \text{otherwise} \end{cases}$

$$C = \sum_{1 \leq i < j \leq n} c_{ij} \quad (1)$$

$$E(C) = E\left(\sum_{i < j} c_{ij}\right) = \sum_{i < j} E(c_{ij}) \quad (2)$$

$$E(c_{ij}) = 1 \cdot P(c_{ij} = 1) + 0 \cdot P(c_{ij} = 0) \quad (3)$$

$$= P(c_{ij} = 1) \quad (4)$$

$$E(c_{ij}) = P(z_i \text{ and } z_j \text{ are compared by RQS(S)}) = \frac{2}{j-i+1} \text{ (Why?)} \quad (5)$$

$$E(C) = \sum_{1 \leq i < j \leq n} \frac{2}{j-i+1} = \dots = 2n \cdot \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n}\right) \quad (6)$$

H_n (H_n is denoted as the harmonic series) is $O(\log n)$. We have a question remaining: Why $E(cij)$ equals that? because for two keys to be compared, if only if one of them is selected as a pivot. If we chose anything in between, they can NEVER be compared. So the case must be that we chose one of them as the pivot. Look at the specific example from the slide. Smallest key is 1 and the largest key is 8.