# **CSC263 Notes**

Data Structures and Analysis

https://github.com/ICPRplshelp/

Last updated January 12, 2023

# 1 Data Types, Data Structures

### **ADTs**

- Specification
  - Objects we're working with
  - The operations (WHAT but not how)

#### **Data structures**

- Implementation (how)
  - Data
  - Algorithms

## **Analysis (runtime or complexity)**

- Worst case
- Best case
- Upper bounds
  - 0
- Lower bounds
  - **-** Ω
- Tight bounds
  - **-** Θ

If we have algorithm A and input x, the runtime  $t_A(x) =$  number of constant time operations independent of x.

Ultimately, we want a measure of running time that is a function of the input size. We have lots of inputs for each input size. So if we want to prove an upper bound when looking for the worst case running time:

- I need two functions to prove an upper bound
  - A simple algebraic expression
  - The running time
    - \* However, the pure runtime function's codomain isn't  $\mathbb{R}^{\geq 0}$  but rather, a list of running times. To turn it into a raw function that outputs  $\mathbb{R}^{\geq 0}$ , we can take the largest of the list I just described.

Worst case is just us narrowing down a bunch of possible runtimes to the worst one.

My upper bound will always be some value that is larger or equal to the worst case, and the lower bound must be below the worst case but not all the worst cases.

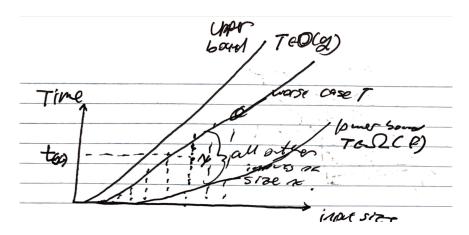


Figure 1: Upper and lower bounds of the worst case

# 1.1 Average-case Running Time

For a particular input size n, I have all these inputs x. We need to precisely define  $S_n = \{\text{set of all inputs of size } n\}$ . We need a probability distribution over our set of inputs.

- For each input  $x \in S_n$ , t(x) is a random variable.
  - It's something that assigns a number to each element of our sample space. Now,  $S_n$  (a finite set) becomes a sample space.

Given a **discrete** probability distribution over  $S_n$ , t(x) = the running time for input x. But what we'll get is that if I only know the input size,  $t_n =$  **average-case running time for size** n **is:** 

$$t_n := E[t(x)] \text{ over } S_n$$
  
=  $\sum_{x \in S_n} t(x) \Pr(x)$ 

This summation is not as easy to figure out. There's a way we can get a particular value, but there is another way:

## **EX:** The linear search algorithm

```
def LinSearch(L: LinkedList, x: T) -> Node | None:
    """Pre: L is a linked list, x is a value
    Post: return the node that contains x
    or none otherwise"""

z = L.head
    while z != None and z.key != x:
    z = z.next
    return z
```

Standard linked list search. Here's what we need to do:

- 1. Define our sample space (a family of sample space, one for each sample size):
  - a. Let n be arbitrary.
  - b.  $S_n = \{\text{every input of size } n\}$
- 2. What should the sample space be? It shouldn't necessarily have infinitely many inputs with a given running time.
  - a. The number of different possible running times I have is finite for this algorithm.
- 3. **INSIGHT.** One representative input for each possible behavior. Behavior = running time.

a. 
$$S_3 = \{([1, 2, 3], 1), ([2, 1, 3], 1), ([2, 3, 1], 1), ([2, 3, 4], 1)\}$$

- b. Alternatively,  $S_n = \{([1, 2, ..., n], 1), ([1, 2, ..., n], 2), ...\}$
- c. Which is  $\{([1, 2, ..., n], x) : x \in [0, n] \land x \in \mathbb{N}\}$
- 4. The probability distribution becomes important. How do we decide how likely we want each input to be? How can we tell? That is a tricky question. What are we trying to do, and there's no obvious way to choose. In practice, in any kind of real-life scenario, if you want to judge how well an algorithm performs on average if you have some idea of what your real-life inputs are going to look like. If you have no information at all, where it is all abstract, then we'll just uniformly distribute.

a. 
$$\Pr([1, 2, ..., n], i) = \frac{1}{n+1}$$
 for  $i = 0, 1, ..., n$ 

b. Now we have this, we can calculate the expected value:

c. 
$$E[t(x)] = \sum_{(L, i) \in S_n} t(L, i) \cdot \Pr(L, i)$$

d. = 
$$\sum_{i=0}^{n} t([1, 2, ..., n], i) \cdot \frac{1}{n+1}$$

- e. When we're doing an average case, we cannot calculate an expected value with  $\mathscr O$  expressions in there. We need a precise expression we can add up and average out. We need an **exact** expression
  - i. Not in the sense that there's one right answer, but we need to fix a particular way of counting and count the same way for every input.
- f. There is one trick: pick some representative operation that we know if we count that, the number of representative operations is  $\Theta$  (runtime).
  - i. In the example, the number of times  $z \cdot key == x$  is run, which will be the **representative operation**
- g. Ignore the constant time operations. The thing that matters is the loop. The loop does a constant amount of work each operation.

h. = 
$$t([1, 2, ..., n], 0) \cdot \frac{1}{n+1} + \sum_{i=1}^{n} i \cdot \frac{1}{n+1}$$

i. It's all algebra by this point. Do all of it, and you should end up with  $\frac{n}{2}+\frac{n}{n+1}\in\Theta(n)$ .