

# DSC 190

DATA STRUCTURES & ALGORITHMS

Lecture 1 | Part 1

**Welcome!**

# Advanced Data Structures and Algorithms

(for data science)

- ▶ Third time being taught.
- ▶ Modeled (partly) after CSE 100/101.
- ▶ But with more data science flavor.

# Roadmap

- ▶ Advanced Data Structures
  - ▶ Dynamic Arrays
  - ▶ AVL Trees
  - ▶ Heaps
  - ▶ Disjoint Set Forests
  
- ▶ Nearest Neighbor Queries
  - ▶ KD-Trees
  - ▶ Locality Sensitive Hashing

# Roadmap

- ▶ Strings
  - ▶ Tries and Suffix Trees
  - ▶ Knuth-Morris-Pratt and Rabin-Karp string search
  
- ▶ Algorithm Design
  - ▶ Divide and Conquer
  - ▶ Greedy Algorithms
  - ▶ Dynamic Programming (Viterbi Algorithm)
  - ▶ Backtracking, Branch and Bound
  - ▶ Linear Time Sorting; Sort with Noisy Comparator

# Roadmap

- ▶ Sketching and Streaming
  - ▶ Count-min-sketch
  - ▶ Bloom filters
  - ▶ Reservoir Sampling?
- ▶ Theory of Computation
  - ▶ NP-Completeness and NP-Hardness
  - ▶ Computationally-hard problems in ML/DS

# Prerequisite Knowledge

- ▶ Python
- ▶ Basic Data Structures and Algorithms
  - ▶ DSC 30, DSC 40B

# Syllabus

- ▶ All course materials can be found at [dsc190.com](https://dsc190.com).

# DSC 190

DATA STRUCTURES & ALGORITHMS

Lecture 1 | Part 2

**Review of Time Complexity Analysis**



# Time Complexity Analysis

- ▶ Determine efficiency of code **without** running it.
- ▶ Idea: find a formula for time taken as a function of input size.

# Advantages of Time Complexity

1. Doesn't depend on the computer.
2. Reveals which inputs are slow, which are fast.
3. Tells us how algorithm scales.

# Counting Operations

- ▶ Abstraction: certain basic operations take **constant time**, no matter how large the input data set is.
- ▶ Example: addition of two integers, assigning a variable, etc.
- ▶ Idea: count basic operations

## Example

```
def mean(numbers):  
    total = 0  
    n = len(numbers)  
    for x in numbers:  
        total += x  
    return total / n
```

time to exec once	# of execs
$C_1$	1
$C_2$	1
$C_3$	$n+1$
$C_4$	$n$
$C_5$	1

$O(n)$

# Theta Notation, Informally

- ▶  $\Theta(\cdot)$  forgets constant factors, lower-order terms.

$$5n^3 + 3n^2 + 42 = \Theta(n^3)$$

## Theta Notation, Informally

- ▶  $f(n) = \Theta(g(n))$  if  $f(n)$  “grows like”  $g(n)$ .

$$5n^3 + 3n^2 + 42 = \Theta(n^3)$$

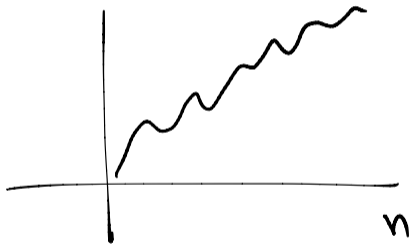
$$2^n \neq \Theta(3^n)$$

## Theta Notation Examples

▶  $4n^2 + 3n - 20 = \Theta(n^2)$

▶  $3n + \sin(4\pi n) = \Theta(n)$

▶  $2^n + 100n = \Theta(2^n)$



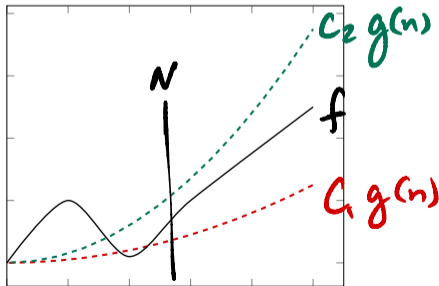
$$2^{n+1} = \Theta(2^n)$$

$$2^{n+1} = \underbrace{2}_{\sim} \cdot 2^n$$

## Definition

We write  $f(n) = \Theta(g(n))$  if there are positive constants  $N$ ,  $c_1$  and  $c_2$  such that for all  $n \geq N$ :

$$c_1 \cdot g(n) \leq f(n) \leq c_2 \cdot g(n)$$





## Main Idea

If  $f(n) = \Theta(g(n))$ , then  $f$  can be “sandwiched” between copies of  $g$  when  $n$  is large.

# Other Bounds

- ▶  $f = \Theta(g)$  means that  $f$  is both **upper** and **lower** bounded by factors of  $g$ .
- ▶ Sometimes we only have (or care about) upper bound or lower bound.
- ▶ We have notation for that, too.

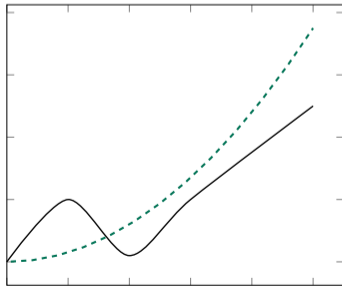
# Big-O Notation, Informally

- ▶ Sometimes we only care about upper bound.
- ▶  $f(n) = O(g(n))$  if  $f(n)$  “grows at most as fast” as  $g(n)$ .
- ▶ Examples:
  - ▶  $4n^2 = O(n^{100})$
  - ▶  $4n^2 = O(n^3)$
  - ▶  $4n^2 = O(n^2)$  and  $4n^2 = \Theta(n^2)$

## Definition

We write  $f(n) = O(g(n))$  if there are positive constants  $N$  and  $c$  such that for all  $n \geq N$ :

$$f(n) \leq c \cdot g(n)$$



$\omega(n)$ 

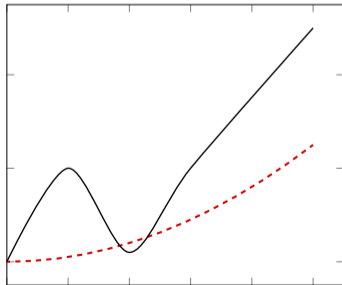
## Big-Omega Notation

- ▶ Sometimes we only care about lower bound.
- ▶ Intuitively:  $f(n) = \Omega(g(n))$  if  $f(n)$  “grows at least as fast” as  $g(n)$ .
- ▶ Examples:
  - ▶  $4n^{100} = \Omega(n^5)$
  - ▶  $4n^2 = \Omega(n)$
  - ▶  $4n^2 = \Omega(n^2)$  and  $4n^2 = \Theta(n^2)$

## Definition

We write  $f(n) = \Omega(g(n))$  if there are positive constants  $N$  and  $c$  such that for all  $n \geq N$ :

$$c_1 \cdot g(n) \leq f(n)$$



$$n^2 + n^3 = \Theta(n^3)$$

## Sums of Theta

- ▶ If  $f_1(n) = \Theta(\overset{n^2}{g_1(n)})$  and  $f_2(n) = \Theta(\overset{n^3}{g_2(n)})$ , then

$$\begin{aligned} f_1(n) + f_2(n) &= \Theta(g_1(n) + g_2(n)) \\ &= \Theta(\max(g_1(n), g_2(n))) \end{aligned}$$

- ▶ Useful for sequential code.

## Products of Theta

- If  $f_1(n) = \Theta(n^2)$  and  $f_2(n) = \Theta(n^3)$ , then

$$f_1(n) \cdot f_2(n) = \Theta(n^2 \cdot n^3)$$

$$n^2 \times n^3 = \Theta(n^2 \times n^3) = \Theta(n^5)$$



$$f_1 \times f_2 = \Theta(n^5)$$



$$\Theta(n^2) \quad \Theta(n^3)$$

## Example

```
def foo(n):  
    for i in range(3*n + 4, 5n**2 - 2*n + 5):  
        for j in range(500*n, n**3):  
            print(i, j)
```

# Linear Search

- ▶ **Given:** an array `arr` of numbers and a target `t`.
- ▶ **Find:** the index of `t` in `arr`, or **None** if it is missing.

```
def linear_search(arr, t):  
    for i, x in enumerate(arr):  
        if x == t:  
            return i  
    return None
```

## Exercise

What is the time complexity of `linear_search`?

## The **Best** Case

- ▶ When  $t$  is the very first element.
- ▶ The loop exits after one iteration.
- ▶  $\Theta(1)$  time?

## The **Worst** Case

- ▶ When  $t$  is not in the array at all.
- ▶ The loop exits after  $n$  iterations.
- ▶  $\Theta(n)$  time?

# Time Complexity

- ▶ `linear_search` can take vastly different amounts of time on two inputs of the **same size**.
  - ▶ Depends on **actual elements** as well as size.
- ▶ There is no single, overall time complexity here.
- ▶ Instead we'll report **best** and **worst** case time complexities.

# Best Case Time Complexity

- ▶ How does the time taken in the **best case** grow as the input gets larger?



## Definition

Define  $T_{\text{best}}(n)$  to be the **least** time taken by the algorithm on any input of size  $n$ .

The asymptotic growth of  $T_{\text{best}}(n)$  is the algorithm's **best case time complexity**.

## Best Case

- ▶ In `linear_search`'s **best case**,  $T_{\text{best}}(n) = c$ , no matter how large the array is.
- ▶ The **best case time complexity** is  $\Theta(1)$ .

# Worst Case Time Complexity

- ▶ How does the time taken in the **worst case** grow as the input gets larger?

## Definition

Define  $T_{\text{worst}}(n)$  to be the **most** time taken by the algorithm on any input of size  $n$ .

The asymptotic growth of  $T_{\text{worst}}(n)$  is the algorithm's **worst case time complexity**.

## Worst Case

- ▶ In the worst case, `linear_search` iterates through the entire array.
- ▶ The **worst case time complexity** is  $\Theta(n)$ .

# Faux Pas

- ▶ Asymptotic time complexity is not a **complete** measure of efficiency.
- ▶  $\Theta(n)$  is not always better than  $\Theta(n^2)$ .
- ▶ Why?

## Faux Pas

- ▶ **Why?** Asymptotic notation “hides the constants”.
- ▶  $T_1(n) = 1,000,000n = \Theta(n)$
- ▶  $T_2(n) = 0.00001n^2 = \Theta(n^2)$
- ▶ But  $T_1(n)$  is **worse** for all but really large  $n$ .

## Main Idea

Asymptotic time complexity is not the **only** way to measure efficiency, and it can be misleading.

Sometimes even a  $\Theta(2^n)$  algorithm is better than a  $\Theta(n)$  algorithm, if the data size is small.



# DSC 190

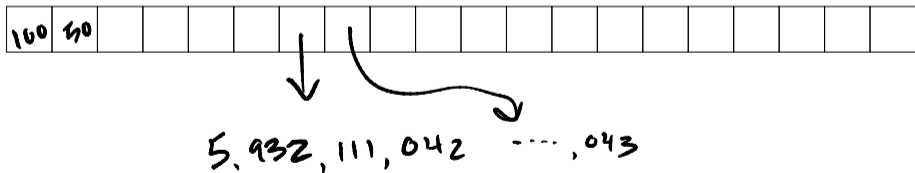
DATA STRUCTURES & ALGORITHMS

Lecture 1 | Part 3

**Arrays and Linked Lists**

# Memory

- ▶ To access a value, we must know its **address**.

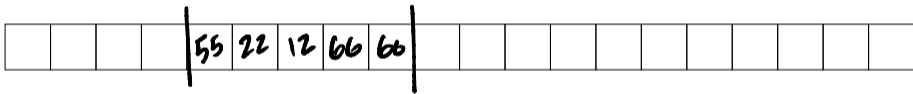


# Sequences

- ▶ How do we store an **ordered sequence**?
  - ▶ e.g.: 55, 22, 12, 66, 60
- ▶ Array? Linked list?

# Arrays

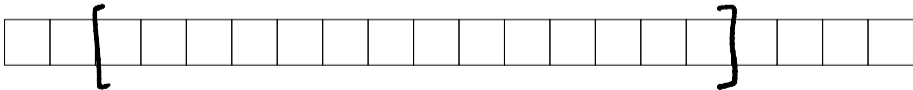
- ▶ Store elements **contiguously**.
  - ▶ e.g.: 55, 22, 12, 66, 60



- ▶ NumPy arrays are... arrays.

# Allocation

- ▶ Memory is shared resource.
- ▶ A chunk of memory of fixed size has to be reserved (**allocated**) for the array.
- ▶ The size has to be known beforehand.



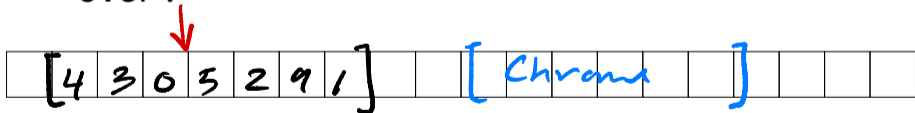
arr[5]

## Arrays

- ▶ To access an element, we need its address.
- ▶ **Key:** Addresses are easily calculated.
  - ▶ For  $k$ th element: address of first + ( $k \times 64$  bits)
- ▶ Therefore, arrays support  $\Theta(1)$ -time access.

# Downsides of Arrays

- ▶ Homogeneous; every element must be same size.
- ▶ To **resize** the array, a totally new chunk of memory has to be found; old values copied over<sup>1</sup>.



---

<sup>1</sup>In worst case: see realloc

*malloc*

# Array Time Complexities

- ▶ Retrieve  $k$ th element:  $\Theta(1)$  (**good**).
- ▶ Append element at end:  $\Theta(n)$  (**bad**)<sup>2</sup>.
- ▶ Insert/remove in middle:  $\Theta(n)$  (**bad**).
- ▶ Allocation:  $\Theta(n)$  if initialized,<sup>3</sup> else  $\Theta(1)$

---

<sup>2</sup>At least on average. See: `realloc`

<sup>3</sup>On Linux this is done lazily, as can be seen by timing `np.zeros`



## Aside: np.append

```
»> arr = np.array([1, 2, 3])
»> np.append(arr, 4) # takes Theta(n) time!
array([1, 2, 3, 4])
```

## Aside: np.append

```
results = np.array([])
for i in np.arange(100):
    result = run_simulation()
    results = np.append(results, result)
```

$$\Theta(n^2)$$

## Aside: np.append

- ▶ This was **bad** code!
- ▶ We allocate/copy a **quadratic** number of elements:

$$\underbrace{1}_{\text{1st iter}} + \underbrace{2}_{\text{2nd iter}} + \underbrace{3}_{\text{3rd iter}} + \dots + \underbrace{100}_{\text{last iter}} = \frac{100 \times 101}{2} = 5050$$

## Aside: np.append

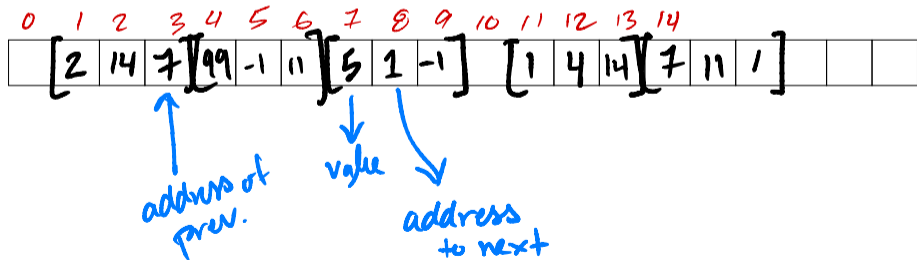
- ▶ Better: pre-allocate.

```
results = np.empty(100)
for i in np.arange(100):
    results[i] = run_simulation()
```

< 5, 2, 7, 1

## (Doubly) Linked Lists

- ▶ Scatter elements throughout memory.
- ▶ For each, store address of next/previous.



# Linked Lists

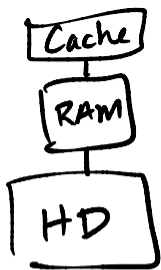
- ▶ Each element has an **address**.
- ▶ Keep track of the address of first/last elements.
- ▶ Have to **find** address of middle elements by looping.

# Linked List Time Complexities

- ▶ Retrieve  $k$ th element:
  - ▶  $\Theta(k)$  if you don't know address (**bad**)<sup>4</sup>
  - ▶  $\Theta(1)$  if you do
- ▶ Append/pop element at start/end:  $\Theta(1)$  (**good**).
- ▶ Insert/remove  $k$ th element:
  - ▶  $\Theta(k)$  if you don't know address (**bad**)
  - ▶  $\Theta(1)$  if you do
- ▶ Allocation not needed! (**good**)

---

<sup>4</sup>assumes search starts from beginning



## Tradeoffs

- ▶ Arrays are better for numerical algorithms.
  - ▶ Arrays have good cache performance.
- ▶ Linked lists are better for stacks and queues.



## Main Idea

Different data structures optimize for different operations.

## Next time...

- ▶ Can we have the best of both?
- ▶ I.e., a data structure with the “growability” of linked lists, but the fast access of arrays.

## Next time...

- ▶ Can we have the best of both?
- ▶ I.e., a data structure with the “growability” of linked lists, but the fast access of arrays.
- ▶ Yes, in a sense: the **dynamic array**.

## Next time...

- ▶ Can we have the best of both?
- ▶ I.e., a data structure with the “growability” of linked lists, but the fast access of arrays.
  
- ▶ Python’s `list` is a dynamic array.