

# Theo's DSA Review

Theo Park

Compiled: July 2, 2025

# Contents

<b>1</b>	<b>Runtime Analysis</b>	<b>4</b>
1.1	Power Law . . . . .	4
1.2	Runtime Expressions . . . . .	4
1.3	Asymptotic Runtime Analysis . . . . .	5
<b>2</b>	<b>Intro to Data Structures</b>	<b>6</b>
2.1	Array . . . . .	6
2.2	Linked List . . . . .	6
2.3	Stack . . . . .	6
2.4	Queue . . . . .	7
2.4.1	Implementing Queue Using a Circular Array . . . . .	7
2.5	Tree . . . . .	7
2.6	Binary Heap . . . . .	9
2.6.1	Heapify . . . . .	9
<b>3</b>	<b>Sorting Algorithms</b>	<b>10</b>
3.1	Bubble, Selection, Insertion Sort . . . . .	10
3.2	Shell Sort . . . . .	12
3.3	Heap Sort . . . . .	12
3.3.1	Sort Down Algorithm . . . . .	12
3.4	Merge Sort . . . . .	13
3.5	Quick Sort . . . . .	14
3.5.1	Partition and Pivot . . . . .	15
3.6	Comparison Sorting Algo Lower Bound . . . . .	15
3.7	Bucket Sort . . . . .	16
3.8	Counting Sort . . . . .	16
3.9	Radix Sort . . . . .	17
3.9.1	Lexicographic Order . . . . .	17
<b>4</b>	<b>Hash Tables</b>	<b>18</b>
4.1	Division Method . . . . .	18
4.2	Multiplication Method . . . . .	18
4.3	Collision . . . . .	18
4.3.1	Chaining . . . . .	18
4.3.2	Open Addressing . . . . .	18
<b>5</b>	<b>Search Tree</b>	<b>19</b>
5.1	Binary Search Tree and Its Limit . . . . .	19
5.2	2-3 Tree . . . . .	19
5.3	Red-Black Tree . . . . .	19
5.4	Left-Leaning Red-Black Tree . . . . .	19
5.4.1	Deletion in LLRBT . . . . .	19

<b>6</b>	<b>Undirected Graph</b>	<b>20</b>
6.1	Adjacency Matrix and List . . . . .	20
6.2	BFS . . . . .	21
6.3	DFS . . . . .	21
<b>7</b>	<b>Directed Graphs</b>	<b>23</b>
7.1	Strong Connectivity . . . . .	23
7.1.1	Strongly Connected Components and Kosaraju's Algorithm . . . . .	24
7.2	Articulation Points . . . . .	24
7.3	Directed Acyclic Graphs . . . . .	25
7.3.1	Topological Sort . . . . .	25
<b>8</b>	<b>Weighted Graphs</b>	<b>27</b>
8.1	Shortest Path . . . . .	27
8.1.1	Initialization and Edge Relaxation for Single-source Shortest-path Algorithms . . . . .	27
8.1.2	Bellman-Ford Algorithm . . . . .	28
8.1.3	Dijkstra's Algorithm . . . . .	28
8.1.4	Floyd-Warshall All Pairs Shortest Path Algorithm . . . . .	30
8.2	Minimum Spanning Tree . . . . .	30
8.2.1	Cycle and Cut Properties . . . . .	31
8.2.2	Prim's Algorithm . . . . .	31
8.2.3	Union-Find . . . . .	33
8.2.4	Kruskal MST Algorithm . . . . .	36
<b>9</b>	<b>Strings</b>	<b>38</b>
9.1	Brute-force String Pattern Matching . . . . .	38
9.2	KMP Algorithm . . . . .	38
9.3	Trie . . . . .	40
9.4	PATRICIA . . . . .	43
9.5	Huffman Coding . . . . .	43

# Preface

This was a small project of mine while I was taking a gap year, and I had great pleasure writing this. I am still hoping to further refine this review with more examples and explanations, but it is time to return to school and learn new things.

The materials in this review are loosely based on:

- Purdue University CS251: Data Structures & Algorithms, taught in the fall of 2022 by Prof. Andres Bejarano – Thanks, Prof. Bejarano
- Introduction to Algorithms, Third Edition by "CLRS"
- My notes for various Leetcode questions

Most of the algorithms in the review are written in Python to provide real-life implementations. You can find the sources in the `CODE` directory of the GitHub repository. Some test cases are generated using generative AI models, such as the Windsurf model with the Neocodeium Neovim plugin or OpenAI ChatGPT. However, the use of LLM was restricted to the test cases, and all codes are either a direct Python translation of the pseudocodes in CLRS or my own creation.

For more details regarding compilation of the  $\text{\TeX}$  source, see `README.MD` in the GitHub repository.

Constructive criticism is always welcome, but I do not take any responsibility for any errors in the document.

**Every file in the repository, including but not limited to the  $\text{\TeX}$  source files, is licensed under the following.**

MIT License

Copyright (c) 2022–present Theo Park

All rights reserved.

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

# Chapter 1

## Runtime Analysis

*Algorithms* are any well-defined computational procedures that take some value(s) as input and produce more value(s) as output. They are **effective**, **precise**, and **finite**. There are several ways to analyze the runtime of an algorithm.

### 1.1 Power Law

Construct a hypothesis using  $T(n) = an^b$ , where  $b = \log(\text{ratio})$ . This section is not terribly important.

### 1.2 Runtime Expressions

<hr/> <b>for</b> $i = 1$ <b>to</b> $n$ <b>do</b> $\triangleright$ <i>for</i> ( $i = 1$ ; $i \leq n$ ; $i++$ ) └ A constant time operation <hr/> $T(n) = \sum_{i=1}^n c = cn$	<hr/> <b>for</b> $i = 0$ <b>to</b> $n - 1$ <b>do</b> $\triangleright$ <i>for</i> ( $i = 0$ ; $i < n$ ; $i++$ ) └ A constant time operation <hr/> $T(n) = \sum_{i=0}^{n-1} c = c(n - 1 + 1) = cn$
<hr/> 1: $\triangleright$ <i>for</i> ( <i>int</i> $i = 0$ ; $i \leq n$ ; $i += 2$ ) $\triangleleft$ 2: <b>for</b> $i = 0$ <b>to</b> $n$ , increment by 2 <b>do</b> 3:   └ $n$ time operation <hr/> Note that $i$ increases by 2, so express its sequence in terms of increments by 1. $i = 0, 2, 4, 6, \dots, n = 2 \times 0, 2 \times 1, 2 \times 2, \dots, 2 \times \frac{n}{2}$ $T(n) = \sum_{k=0}^{\frac{n}{2}} n = n \sum_{k=0}^{\frac{n}{2}} 1 = n\left(\frac{n}{2} + 1\right) = \frac{1}{2}n^2 + n$	<hr/> 1: <b>for</b> $i = 0$ <b>to</b> $n$ <b>do</b> 2:   └ <b>for</b> $j = 1$ <b>to</b> $n - 1$ <b>do</b> 3:     └ A constant time operation <hr/> $T(n) = \sum_{i=0}^n \sum_{j=1}^{n-1} c = \sum_{i=0}^n c(n - 1) = c(n - 1)(n + 1) = cn^2 - c$
<hr/> 1: <b>for</b> $i = 0$ <b>to</b> $n$ , increment by 2 <b>do</b> 2:   └ <b>for</b> $j = 1$ <b>to</b> $n - 1$ <b>do</b> 3:     └ $n$ time operation <hr/> $T(n) = \sum_{k=0}^{\frac{n}{2}} \sum_{j=1}^{n-1} n = \sum_{k=0}^{\frac{n}{2}} n(n - 1) = n(n - 1)\left(\frac{n}{2} + 1\right) = \frac{1}{2}n^3 + \frac{1}{2}n^2 - n$	<hr/> 1: <b>for</b> $i = 0$ <b>to</b> $n$ <b>do</b> 2:   └ <b>for</b> $j = i$ <b>to</b> $n - 1$ <b>do</b> 3:     └ $n$ time operation <hr/> $T(n) = \sum_{i=0}^n \sum_{j=i}^{n-1} 1 = \sum_{i=0}^n (\sum_{j=0}^{n-1} n - \sum_{j=0}^{i-1} n) = \dots = n^2(n + 1) - \frac{1}{2}n^2(n + 1) = \frac{1}{2}n^2(n + 1)$

---

```

1: for  $i = 0$  to  $n$ , increment by 2 do
2:   for  $j = i$  to  $n - 1$  do
3:      $n$  time operation

```

---

$$T(n) = \sum_{k=0}^{\frac{n}{2}} \sum j = 2k^{n-1}n = \sum_{k=0}^{\frac{n}{2}} (\sum_{j=0}^{n-1} n - \sum_{j=0}^{2k-1} n) = \dots = n^2(\frac{n}{2} + 1) - 2n\frac{\frac{n}{2}(\frac{n}{2}+1)}{2} = \frac{n^3}{2} + n^2 - \frac{n^3}{4} - \frac{n^2}{2} = \frac{1}{4}n^3 + \frac{1}{2}n^2$$


---

```

1: for  $i = 1$  to  $n, i = i * 3$  do
2:    $\log(n)$  time operation

```

---

$$T(n) = \sum_{k=0}^{\log_3(n)} \log(n) = \log(n)(\log_3(n) + 1) \approx \log(n)\log(n) + \log(n)$$


---

---

```

1: for  $i = 0$  to  $n, i = i * 2$  do
2:   A constant time operation

```

---

Express the sequence in terms of increments by 1

- $i = 1 = 2^0$
- $i = 2 = 2^1$
- $i = 4 = 2^2$
- $i = 2^k \leq n \rightarrow k \leq \log_2(n)$

$$T(n) = \sum_{k=0}^{\log_2(n)} c = c(\log_2(n) + 1) = c\log_2(n) + c$$


---

```

1: for  $i = 1$  to  $n$  do
2:    $n$  time operation

```

---

$$T(n) = \sum_{k=0}^{\log_2(n)} n = n(\log_2(n) + 1) = n\log_2(n) + n \approx n\log(n) + n$$


---

### 1.3 Asymptotic Runtime Analysis

- Big-O:  $f(n) \in O(g(n))$  if there exist positive constants  $c$  and  $n_0$  such that  $0 \leq f(n) \leq cg(n)$  for all  $n \geq n_0$
- Big-Omega:  $f(n) \in \Omega(g(n))$  if there exist positive constants  $c$  and  $n_0$  such that  $0 \leq c(g)n \leq f(n)$  for all  $n \geq n_0$
- Big-Theta: One number is both Big-O and Big-Omega

## Chapter 2

# Intro to Data Structures

*Data structures* are collections of data values, the relationships among them, and the functions or operations that can be applied to the data. All three characteristics need to be present.

## 2.1 Array

*Array* is a linear container of items.

Array length 6	0	27	12	39	24	9
	0	1	2	3	4	5

- Access time:  $\Theta(1)$
- Inserting  $n$  items in the *tail* for array size  $n$ :  $\Theta(1)$  per item,  $n \times \Theta(1) \in \Theta(1)$
- Inserting  $n$  items in the *tail* for array size *unknown*:  $\Theta(n)$  per item,  $n \times \Theta(n) \in \Theta(n)$
- Resizing:  $\Theta(n)$

Lesson? **Keep track of the size!**

## 2.2 Linked List



- Access time:  $\Theta(n)$
- Insertion:  $\Theta(1)$
- Deletion:
  - Singly LL:  $O(n)$
  - Doubly LL:  $\Theta(1)$

## 2.3 Stack

-	-	-	7	-	-	-	
-	-	5	5	5	-	-	Error
-	3	3	3	3	3	-	
init()	push(3)	push(5)	push(7)	pop()	pop()	pop()	pop()

Stack is a "last in, first out" (LIFO) data structure.

- `push(item)`: Inserting an item on the top –  $\Theta(1)$
- `pop()`: Removes and returns the item on the top –  $\Theta(1)$
- `peek()`: Returns the item on the top
- `size()`: Returns the number of item in the stack
- `isEmpty()`: Checks if `size == 0`

## 2.4 Queue

```
- - - 3 - - 3 5 - 3 5 7 - 5 7 - - 7 - - - Error
init() e(3)   e(5)   e(7)   d()   d()   d()   d()
```

Queue is a "first in, first out" (FIFO) data structure.

- `enqueue(item)`: Inserts an item at the end -  **$\Theta(1)$**
- `dequeue()`: Removes the item at the front of the queue and return it -  **$\Theta(1)$**
- `peek()`: Returns the item at the front of the queue w/o removing it
- `size()`: Returns the number of the items in the queue
- `isEmpty()`: Checks if `size == 0`

### 2.4.1 Implementing Queue Using a Circular Array

---

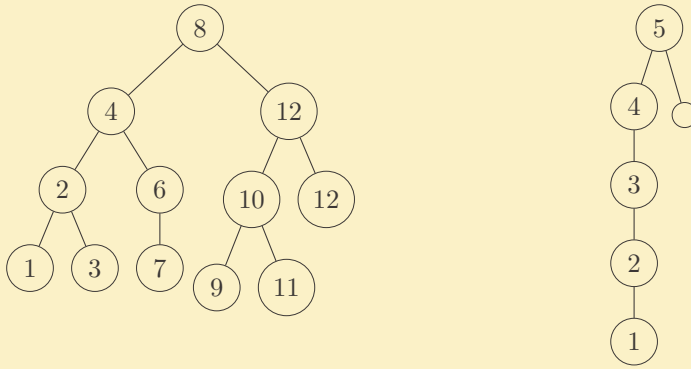
```
1: function CONSTRUCTOR( $k$ )                                ▷  $k$  is the initial size of the array
2:    $Q \leftarrow$  an array of size  $k$ 
3:    $size \leftarrow 0$ 
4:    $head \leftarrow 0$ 
5:    $tail \leftarrow k - 1$ 
1: function ENQUEUE( $n$ )                                    ▷  $n$  is the new item to add
2:   if  $size = k$  then return False
3:    $tail \leftarrow (tail + 1) \bmod k$ 
4:    $Q[tail] \leftarrow n$ 
5:    $size \leftarrow size + 1$ 
1: function DEQUEUE
2:   if  $size = 0$  then return False
3:    $tmp \leftarrow Q[head]$ 
4:    $Q[head] \leftarrow nil$ 
5:    $head \leftarrow (head + 1) \bmod k$ 
6:    $size \leftarrow size - 1$ 
   return tmp
```

---

## 2.5 Tree

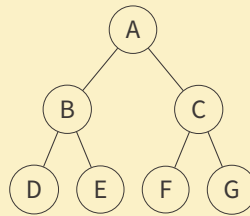
- Search:  $O(\log(n))$  for a balanced tree (right),  $O(n)$  for unbalanced tree (left)





- **Maximum number of leaves:**  $2^h$ , where  $h$  is the height of the tree
- **Maximum number of nodes:**  $\sum_{i=0}^h 2^i = 2^{h+1} - 1$ , where  $h$  is the height of the tree

Following calculations and traversals will be explained using the following **full binary tree**.



Definitions:

- $I$  (number of internal nodes) = 3 (A, B, C – root is an internal node unless it is a leaf)
- $N$  (number of nodes) = 7
- $L$  (number of leaves) = 4

Full binary tree properties:

- $L = I + 1 = \frac{N+1}{2}$
- $N = 2I + 1 = 2L - 1$
- $I = \frac{N-1}{2} = L - 1$

Traversal (applicable for all trees, including non-full-binary trees):

- **Preorder – [N]ode [L]eft [R]ight:** A B D E C F G
- **Inorder – [L][N][R]:** D B E A F C
- **Postorder – [L][R][N]:** D E B F C A
- Level – by height: A B C D E F

## 2.6 Binary Heap

- Max heap: Key in each node is larger than or equal to the keys in node's two children
- Min heap: Key in each node is less than or equal to the keys in node's two children

Array implementation for the above max-heap is as follow:

- $\text{LEFT-CHILD}(i) = 2i + 1$
- $\text{RIGHT-CHILD}(i) = 2i + 2$
- $\text{PARENT}(i) = \lfloor \frac{i-1}{2} \rfloor, i > 0$

### 2.6.1 Heapify

1. Define a function MAX-HEAPIFY that given an index  $i$ , it "sinks down"  $A[i]$  to find its correct heap position.
2. Starting at the last non-leaf node (i.e., parent of the last node), run MAX-HEAPIFY to heapify the subtree.

$O(n)$

```
def max_heapify(arr, i):
    n = len(arr)

    l = 2 * i + 1
    r = 2 * i + 2

    largest = i
    if l < n and arr[l] > arr[i]:
        largest = l

    if r < n and arr[r] > arr[largest]:
        largest = r

    if largest != i:
        arr[i], arr[largest] = arr[largest], arr[i]
        max_heapify(arr, largest)

def build_heap(arr):
    n = len(arr)

    lastLeftNode = (n // 2) - 1 # floor(n / 2) - 1
    for i in range(lastLeftNode, -1, -1):
        max_heapify(arr, i)
```

## Chapter 3

# Sorting Algorithms

Once you store all the items in a data structure, you might want to organize them for the future use (such as selecting nth largest element). For this, you have to *sort* the data structure (in this book, array will be assumed). *Sorting* is deciding how to permute the array elements until they are sorted.

There are couple aspects of sorting algorithms you need to consider:

- Runtime: When analyzing a runtime of a sorting algorithm, both number of compares and number of swaps are considered. **Most sorting algorithms make more comparisons than swaps**, but if a sorting algorithm makes more swaps, it must be used for the asymptotic runtime analysis
- Stability: An algorithm is stable if it preserves the input ordering of equal items.

unsorted list	<table><tr><td>5</td><td>3</td><td>5'</td><td>2</td><td>7</td></tr></table>	5	3	5'	2	7
5	3	5'	2	7		
sorted with stable sorting algorithm	<table><tr><td>2</td><td>3</td><td>5</td><td>5'</td><td>7</td></tr></table>	2	3	5	5'	7
2	3	5	5'	7		
sorted with unstable sorting algorithm	<table><tr><td>2</td><td>3</td><td>5'</td><td>5</td><td>7</td></tr></table>	2	3	5'	5	7
2	3	5'	5	7		

- In-place: An algorithm is in-place if it can directly sorts the items without making a copy or extra array(s)

### 3.1 Bubble, Selection, Insertion Sort

The first three sorting algorithms we will discuss are BUBBLE-SORT, SELECTION-SORT, and INSERTION-SORT. They are simple to understand and implement, however, they all have the worst-case runtime of  $O(n^2)$ , which limits their usages.

BUBBLE-SORT goes through the array and swap elements that are out of order, and if such element is found, it re-starts from the beginning of the list.

```
def bubble_sort(arr):
    n = len(arr)
    repeat = True
    while repeat:
        repeat = False

        for i in range(n - 1):
            if arr[i] > arr[i + 1]:
                arr[i], arr[i + 1] = arr[i + 1], arr[i]
                repeat = True

    return arr
```

In-place?	Stable?
True	True

-	NumCompares	NumSwaps
Already Sorted	$n - 1$	0
Worst Case	$n^2 - n$	$\frac{1}{2}n^2 - \frac{1}{2}n$

SELECTION-SORT is a sorting algorithm closest to our “natural” thought of sorting an array. It finds the smallest, second smallest, ... elements and place them accordingly. It makes the same number of comparisons in any cases.

```
def selection_sort(arr):
    n = len(arr)
    for i in range(n - 1):
        min_idx = i

        for j in range(i + 1, n):
            if arr[j] < arr[min_idx]:
                min_idx = j

        if i != min_idx:
            arr[i], arr[min_idx] = arr[min_idx], arr[i]

    return arr
```

In-place?	Stable?
True	False

-	NumCompares	NumSwaps
Already Sorted	$\frac{1}{2}n^2 - \frac{1}{2}n$	0
Worst Case	$\frac{1}{2}n^2 - \frac{1}{2}n$	$\lfloor \frac{1}{2}n \rfloor$

INSERTION-SORT takes one element at a time and places it in the correct index of the partially sorted sub-array until the array is sorted.

```
def insertion_sort(arr):
    n = len(arr)

    for i in range(1, n):
        j = i - 1

        while j >= 0 and arr[j] > arr[j + 1]:
            arr[j], arr[j + 1] = arr[j + 1], arr[j]
            j -= 1

    return arr
```

In-place?	Stable?
True	True

-	NumCompares	NumSwaps
Already Sorted	$n - 1$	0
Worst Case	$\frac{1}{2}n^2 - \frac{1}{2}n$	$\frac{1}{2}n^2 - \frac{1}{2}n$

## 3.2 Shell Sort

SHELL-SORT runs insertion sort with “gaps”, the index margin where insertion sort will be performed.

---

```

1: function SHELL-SORT( $A, H$ )                                ▷  $A$  is an array size  $n$ ,  $H$  is the array size  $m$  containing gap values
2:   for  $h = 0$  to  $m - 1$  do
3:      $gap \leftarrow H[h]$ 
4:     for  $i = 1$  to  $n - 1$  do
5:        $j \leftarrow i - 1$ 
6:       while  $j \geq 0$  and  $j + gap < n$  do
7:         if  $A[j] > A[j + gap]$  then
8:           SWAP( $A, j, j + gap$ )
9:            $j \leftarrow j - gap$ 
10:  return  $A$ 

```

---

## 3.3 Heap Sort

HEAP-SORT uses binary max-heap to sort an array. While it’s the first sorting algorithm to utilize a data structure, it’s not preferred in real life due to cache issue.

```

def heap_sort(arr):
    n = len(arr)

    build_heap(arr)

    def sortdown(i):
        # move the first element of the heap (largest) to the back
        arr[i], arr[0] = arr[0], arr[i]
        # rebuild heap with the first i elements
        max_heapify(arr, i, 0)

    for i in range(n - 1, 0, -1):
        sortdown(i)

    return arr

```

The algorithm first builds the heap from the array elements (refer to section 2.6.1 for the MAX\_HEAPIFY and BUILD\_HEAP functions). Then the algorithm calls SORT-DOWN from the last heap elements down to the first.

### 3.3.1 Sort Down Algorithm

```

def sortdown(i):
    # move the first element of the heap (largest) to the back
    arr[i], arr[0] = arr[0], arr[i]
    # rebuild heap with the first i elements
    max_heapify(arr, i, 0)

```

Suppose we have a max-heap  $A = [7, 6, 3, 5, 4, 1]$ . Using the property of max-heap that the largest element is always placed in the first index, HEAP-SORT performs SORT-DOWN  $n - 1$  times, which places the last element to a correct position and rebuild the heap on the last of the element.

One can call SWIM-UP instead of SWIM-DOWN to repair the heap. However, the former takes  $O(n \log n)$  time and the latter takes  $O(n)$  time.

1.  $i = 5$ :  
 SWAP( $A, 5, 0$ ): [ 1, 6, 3, 5, 4, 7]  
 SWIM-DOWN( $5 - 1$ ) (rebuild heap from index 0 to 4): [ 6, 5, 3, 1, 4 7]
2.  $i = 4$ :  
 SWAP( $A, 4, 0$ ): [ 4, 5, 3, 1, 6, 7]  
 SWIM-DOWN( $4 - 1$ ): [ 5, 4, 3, 1, 6, 7]
3.  $i = 3$ :  
 SWAP( $A, 3, 0$ ): [ 1, 4, 3, 5, 6, 7]  
 SWIM-DOWN( $3 - 1$ ): [ 4, 1, 3, 5, 6, 7]
4.  $i = 2$ :  
 SWAP( $A, 2, 0$ ): [ 3, 1, 4, 5, 6, 7]  
 SWIM-DOWN( $2 - 1$ ): [ 3, 1, 4, 5, 6, 7]
5.  $i = 1$ :  
 SWAP( $a, 1, 0$ ): [ 1, 3, 4, 5, 6, 7]  
 SWIM-DOWN( $1 - 1$ ): [ 1, 3, 4, 5, 6, 7]

### 3.4 Merge Sort

MERGE-SORT is a *divide-and-conquer* sorting algorithm that divides the array down to multiple lists with one element and merge them together.

```
def merge(arr, l, m, r):
    n1 = m - l + 1
    n2 = r - m
    L = [None] * (n1 + 1)
    R = [None] * (n2 + 1)

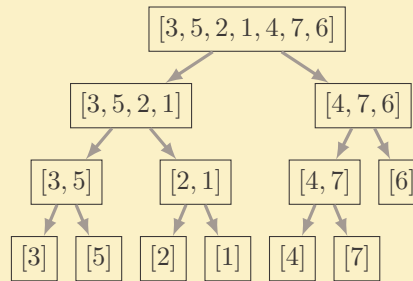
    # arrsign elements to each array
    for i in range(n1):
        L[i] = arr[l + i]
    for j in range(n2):
        R[j] = arr[m + j + 1]

    L[n1], R[n2] = float("inf"), float("inf")
    i, j = 0, 0

    for k in range(l, r + 1):
        if L[i] <= R[j]:
            arr[k] = L[i]
            i += 1
        else:
            arr[k] = R[j]
            j += 1

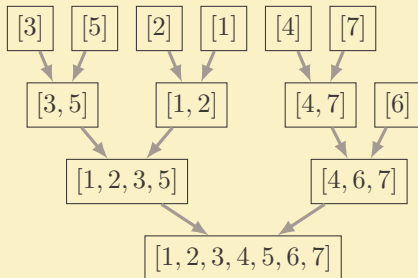
    return arr
```

The algorithm recursively calls itself with half of the current array, and once all the sub-arrays are of size 1, it begins merging them together. The diagram below shows how MERGE-SORT breaks  $A = [3, 5, 2, 1, 4, 7, 6]$  down.



```
def merge_sort(arr, l, r):
    if l < r:
        m = (l + r) // 2
        merge_sort(arr, l, m)
        merge_sort(arr, m + 1, r)
        merge(arr, l, m, r)
    return arr
```

MERGE algorithm merges two sorted arrays by “climbing the ladder”. The diagram on the right shows how MERGE algorithm merges the array we split earlier in the recursive step of the MERGE-SORT, and the diagram on the left shows the visualization of how MERGE algorithm merges two sorted array by “climbing the ladder.”



1	2	3	4	5	6	7
<u>1</u> : 4	<u>2</u> : 4	<u>3</u> : 4	5 : <u>4</u>	<u>5</u> : 6	$\infty$ : <u>6</u>	$\infty$ : <u>7</u>

1. Sub-arrays  $L$  and  $R$  are sorted, and their last element are set to  $\infty$ . In this example,  $L = [1, 2, 3, 5, \infty]$ ,  $R = [4, 6, 7, \infty]$
2.  $k = 0, i = 0, j = 0$ : Since  $L[0] = 1 < 4 = R[0]$ , place  $L[0]$  in  $A[0]$  and increment  $i$
3.  $k = 1, i = 1, j = 0$ : Since  $L[1] = 2 < 4 = R[0]$ , place  $L[1]$  in  $A[1]$  and increment  $i$
4.  $k = 2, i = 2, j = 0$ : Since  $L[2] = 3 < 4 = R[0]$ , place  $L[2]$  in  $A[2]$  and increment  $i$
5.  $k = 3, i = 3, j = 0$ : Since  $L[3] = 5 > 4 = R[0]$ , place  $R[0]$  in  $A[3]$  and increment  $j$
6.  $k = 4, i = 3, j = 1$ : Since  $L[3] = 5 < 6 = R[1]$ , place  $L[3]$  in  $A[4]$  and increment  $i$
7.  $k = 5, i = 4, j = 1$ : Since  $L[4] = \infty > 6 = R[1]$ , place  $R[1]$  in  $A[5]$  and increment  $j$
8. Because the sub-array  $L$  reached its  $\infty$  element, I will skip the rest of the steps and place rest of the elements in  $R$  to  $A$

## 3.5 Quick Sort

QUICK-SORT is another divide-and-conquer sorting algorithm.

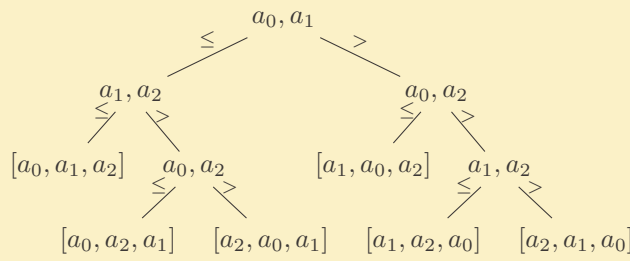
```
def quick_sort(arr, l, r):
    if l < r:
        m = partition(arr, l, r)
        quick_sort(arr, l, m - 1)
        quick_sort(arr, m + 1, r)
    return arr
```

### 3.5.1 Partition and Pivot

```
def partition(arr, l, r):
    p = arr[r]  # rightmost element as the pivot
    i = l - 1
    for j in range(l, r + 1):
        if arr[j] < p:
            i += 1
            arr[i], arr[j] = arr[j], arr[i]
    i += 1
    arr[i], arr[r] = arr[r], arr[i]
    return i
```

## 3.6 Decision Tree and the Lower Bound for Comparison Sorting Algorithm

So far, we have been discussing *comparison sorting algorithms* (sorting algorithms that only reads array elements through  $>$ ,  $=$ ,  $<$  comparison). We can draw a decision tree with all the permutations of how a comparison sorting algorithms would compare and sort the array  $[a_0, a_1, a_2]$ .



The decision tree for array size of 3 has  $3! = 6$  leaves, and it is trivial that a decision tree for an **array with  $n$  elements will have  $n!$  leaves**. The height of the decision tree represents the worst-case number of comparisons the algorithm has to make in order to sort the array.

A tree with the height  $h$  has at most  $2^h$  leaves. Using this property, we can have the lower Big-omega bound for height of the decision tree for comparison sorting algorithms.

$$\begin{aligned}
 2^h &\geq n! \therefore h \geq \log_2(n!) && (h \text{ is the height of the decision tree}) \\
 h &\geq \log_2(n!) = \log_2(1 \cdot 2 \cdots (n-1) \cdot n) \\
 &= \log_2(1) + \log_2(2) + \cdots + \log_2(n-1) + \log_2(n) \\
 &= \sum_{i=1}^n \log_2(i) = \sum_{i=1}^{\frac{n}{2}-1} \log_2(i) + \sum_{i=\frac{n}{2}}^n \log_2(i) \\
 &\geq 0 + \sum_{i=\frac{n}{2}}^n \log_2(i) \geq 0 + \sum_{i=\frac{n}{2}}^n \log_2\left(\frac{n}{2}\right) && (i \text{ in the former expression is } \geq \frac{n}{2}) \\
 &= \frac{n}{2} \log_2\left(\frac{n}{2}\right) && (i = \frac{n}{2} \text{ to } n \text{ is exactly } \frac{n}{2} \text{ iterations}) \\
 &\in \Theta(n \log_2(n))
 \end{aligned}$$

Because  $h \geq \log_2(n!) \in \Theta(n \log_2(n))$ ,  $h \in \Omega(n \log_2(n))$ . Therefore, we can conclude that any **comparison sorting algorithms cannot run faster than  $O(n \log_2(n))$  time** in the worst-case scenario.



### 3.7 Bucket Sort

BUCKET-SORT is a sorting algorithm where elements are divided into each "bucket" and a different sorting algorithm is called for each bucket. While BUCKET-SORT is a comparison sorting algorithm, it's an attempt to reduce the runtime by decreasing the number of elements that the sorting algorithm has to sort.

### 3.8 Counting Sort

COUNTING-SORT is a *non-comparison* sorting algorithm. It uses the extra array *count*, where its index initially represents the value of each element in *A* (e.g., if there are three 5's in *A*,  $count[5] = 3$  before the "accumulation" step to determine the final index), to sort the array.

```
def countingSort(arr):
    n = len(arr)

    # Step 1  $O(n)$ : Find the maximum element and initialize count array
    k = max(arr)
    cnt = [0] * (k + 1)

    # Step 2  $O(n)$ : Find the number of occurrences in each element in the array
    for i in range(n):
        cnt[arr[i]] += 1

    # Step 3  $O(k)$ : Convert count into a prefix sum array of itself
    for i in range(1, k + 1):
        cnt[i] += cnt[i - 1]

    # Step 4  $O(n)$ : Use the count array to find the index of each element
    out = [None] * n
    for i in range(n - 1, -1, -1):
        out[cnt[arr[i]] - 1] = arr[i]
        cnt[arr[i]] -= 1

    return out
```

1. Suppose we have an array  $A = [2, 5, 3, 0, 2, 3, 0, 3]$ .  $k = \text{MAX}(A) = 5$ .
2. Initialize *count*, the array size  $5 + 1$ , with 0's.  $count = [0, 0, 0, 0, 0, 0]$ .
3. Count number of occurrence.  $count = [2, 0, 2, 3, 0, 1]$  (e.g., 2 occurred 2 times)
4. Accumulate values of *count* from left to right.  $count = [2, 2, 4, 7, 7, 8]$  (e.g.,  $count[1] = 2 + 0$ ,  $count[2] = 2 + 0 + 2$ , ...)
5. Initialize *out*, the array size  $n = 8$ , with nil's.  $out = [nil, nil, nil, nil, nil, nil, nil, nil]$ .
6. Place each element to the *out* array using *count* array
  - (a) When  $i = n - 1 = 7$ :  $A[7] = 3$  and  $count[3] = 7 \Rightarrow out[7 - 1] := A[7] = 3$  and  $count[3] := 7 - 1$   
 $out = [nil, nil, nil, nil, nil, nil, 3, nil]$   
 $count = [2, 2, 4, 6, 7, 8]$
  - (b) When  $i = n - 2 = 6$ :  $A[6] = 0$  and  $count[0] = 2 \Rightarrow out[2 - 1] := A[6] = 0$  and  $count[0] := 2 - 1$   
 $out = [nil, 0, nil, nil, nil, nil, 3, nil]$   
 $count = [1, 2, 4, 6, 7, 8]$
  - (c) When  $i = n - 3 = 5$ :  $A[5] = 3$  and  $count[3] = 6 \Rightarrow out[6 - 1] := A[5] = 3$  and  $count[3] := 6 - 1$

```
out = [nil, 0, nil, nil, nil, 3, 3, nil]
count = [1, 2, 4, 5, 7, 8]
```

(d) ...

```
out = [0, 0, 2, 2, 3, 3, 3, 5]
count = [0, 2, 2, 4, 7, 7]
```

In-place?	Stable?
False	True

Because of its use for RADIX-SORT, COUNTING-SORT must be stable, and it indeed is. If there are items with the same value, it will be moved to the *out* array in order in the last (third) for loop.

Runtime	Space Usage
$O(n + k)$	$O(n + k)$

As the algorithm iterates both the size of the array  $n$  and the maximum element in the array  $k$ , **the algorithm runs in  $O(n + k)$  time and uses  $O(n + k)$  space.**

## 3.9 Radix Sort

RADIX-SORT is a non-comparative sorting algorithm for elements with more than one significant digits. It utilizes a stable sorting algorithm such as COUNTING-SORT to sort elements lexicographically.

---

```

1: function RADIX-SORT( $A, k$ )                                ▷  $A$  is an array where the maximum dimension of an element is  $d$ 
2:   for  $i = d$  down to 1 do
3:     | Call a stable sorting algorithm at dimension  $i$ 
4:   return  $A$ 

```

---

### 3.9.1 Lexicographic Order

$$(x_1, x_2, \dots, x_d) < (y_1, y_2, \dots, y_d) \Leftrightarrow (x_i < y_i) \vee (x_1 = y_1 \wedge (x_2, \dots, x_d) < (y_2, \dots, y_d))$$

## **Chapter 4**

# **Hash Tables**

### **4.1 Division Method**

### **4.2 Multiplication Method**

### **4.3 Collision**

#### **4.3.1 Chaining**

#### **4.3.2 Open Addressing**

## **Chapter 5**

# **Search Tree**

### **5.1 Binary Search Tree and Its Limit**

### **5.2 2-3 Tree**

### **5.3 Red-Black Tree**

### **5.4 Left-Leaning Red-Black Tree**

#### **5.4.1 Deletion in LLRBT**

## Chapter 6

# Undirected Graph

*Graph* is a set of vertices  $V$  and a collection of edges  $E$  that connect a pair of vertices.

*Undirected graph* is a graph where edges do not have direction. *Degree of a vertex* representing how many edges is this vertex connected to.

Undirected graph has a few properties:

- **Handshaking Theorem:**  $\sum_{v \in V} \deg(v) = 2 \cdot |E|$  (Sum of degrees of all vertices equal to the number of edges)
- Maximum degree of a vertex:  $\deg(v) \leq |V| - 1$  (connected to every vertices except for itself)
- Maximum edge count:  $|E| \leq \frac{|V|(|V|-1)}{2}$
- **Complete graph:** A graph is said to be complete when each vertex pair is connected by a unique edge. Id est, a complete graph has the maximum number of edges ( $|E| = \frac{|V|(|V|-1)}{2}$ ) (implies that  $\forall v \in V (\deg(v) = |V| - 1)$ )

Path and Cycle:

- Path: Sequence of vertices connected by edges
  - Euler Path: A path that visits every edge exactly once
- Cycle: A path that starts and ends on the same vertex
  - An Euler path that starts and ends on the same vertex

## 6.1 Adjacency Matrix and List

```
class Graph:
    # e.g.,: Graph(n = 3, edges = [[0,1],[1,2],[0,2]])
    # creates a graph with 3 nodes indexed 0, 1, 2,
    # with each connected to each other
    def __init__(self, n: int, edges: list[list[int]]):
        self.n = n

        self.adjList = [[] for _ in range(n)]
        for u, v in edges:
            self.adjList[u].append(v)
            self.adjList[v].append(u)

        self.adjMatrix = [[float("inf") for _ in range(n)] for _ in range(n)]
        for i in range(n):
            self.adjMatrix[i][i] = 0
        for u, v in edges:
```

```

self.adjMatrix[u][v] = 1
self.adjMatrix[v][u] = 1

```

-	Matrix	List
Space	$O( V )$	$O( V  +  E )$
Adding an edge	$O(1)$	$O(1)$
Checking adjacent vertices (to vertex $v$ )	$O(1)$	$O(\deg(v))$
Iteration	$O( V )$	$O(\deg(v))$

## 6.2 BFS

For finding a shortest path in a undirected graph.

Runtime:  $O(|V| + |E|)$

```

def bfs(self, v: int): # v is the starting node
    q = deque()
    dist = [-1] * self.n
    path = [[] for _ in range(self.n)]
    visited = [False] * self.n

    q.append(v)
    dist[v] = 0
    path[v] = [0]
    visited[v] = True

    while q:
        u = q.pop()

        for w in self.adjList[u]:
            if not visited[w]:
                q.append(w)
                dist[w] = dist[u] + 1
                path[w] = path[u] + [w]
                visited[w] = True

    return dist, path, visited

```

## 6.3 DFS

For connectivity.

Iterative implementation using a stack:

```

def dfs_iterative(self, v: int): # v is the starting node
    st = []
    dist = [-1] * self.n
    path = [[] for _ in range(self.n)]
    visited = [False] * self.n

    st.append(v)
    dist[v] = 0
    path[v] = [0]
    visited[v] = True

    while st:

```

```

    u = st.pop()

    for w in self.adjList[u]:
        if not visited[w]:
            st.append(w)
            dist[w] = dist[u] + 1
            path[w] = path[u] + [w]
            visited[w] = True

    return dist, path, visited

```

Using the recursive implementation to count the number of connected components. You can, in fact, change the DFS\_RECURSIVE\_DRIVER to BFS and the algorithm still works as expected.

```

# helper function for finding connected components using recursive DFS
def dfs_recursive_driver(self, v: int,
                        visited,
                        connectedComponents,
                        currConnectedCompIdx):
    visited[v] = True
    connectedComponents[currConnectedCompIdx].append(v)
    for w in self.adjList[v]:
        if not visited[w]:
            self.dfs_recursive_driver(
                w, visited, connectedComponents, currConnectedCompIdx)

def find_connected_comp(self):
    visited = [False] * self.n

    connectedComponents = []
    currConnectedCompIdx = -1
    for v in range(self.n):
        if not visited[v]:
            currConnectedCompIdx += 1
            connectedComponents.append([])
            self.dfs_recursive_driver(
                v, visited, connectedComponents, currConnectedCompIdx)

    return connectedComponents

```

## Chapter 7

# Directed Graphs

The overall constructor is identical to that of undirected graph except that for an edge directing  $u \rightarrow v$ , only  $v$  is added to the adjacency list/matrix of  $u$  and not the other way around.

I also added a recursive implementation of functional DFS method that modifies external VISITED and PATH arrays.

```
class Graph:
    # e.g.,: Graph(n = 3, edges = [[0,1],[1,2],[0,2]])
    # creates a graph with 3 nodes indexed 0, 1, 2,
    # with 0 -> 1, 0 -> 2, 1 -> 2
    def __init__(self, n: int, edges: list[list[int]]):
        self.n = n
        self.edges = edges

        self.adjList = [[] for _ in range(n)]
        for u, v in edges:
            self.adjList[u].append(v)

        self.adjMatrix = [[float("inf") for _ in range(n)] for _ in range(n)]
        for i in range(n):
            self.adjMatrix[i][i] = 0
        for u, v in edges:
            self.adjMatrix[u][v] = 1

    def dfs(self, v, visited, path):
        visited[v] = True
        path.append(v)

        for w in self.adjList[v]:
            if not visited[w]:
                self.dfs(w, visited, path)
```

## 7.1 Strong Connectivity

A directed graph is strongly connected if every vertex is reachable from any other vertex. To test the strong connectivity, use the following algorithm, which takes  $O(|V| + |E|)$ .

- 
- 1: **function** IS-STRONGLY-CONNECTED( $G$ )
  - 2:     Pick a vertex  $v$  in  $G$
  - 3:     Perform DFS( $G, v$ )
  - 4:     If there is an unvisited vertex, return false



---

```

5: | Compute  $G^T$  (transpose of  $G$ , flip all the edges)
6: | Perform DFS( $G^T, v$ )
7: | IF there is an unvisited vertex, return false
8: | return true

```

---

### 7.1.1 Strongly Connected Components and Kosaraju's Algorithm

$O(|V| + |E|)$

1. Perform DFS with a global stack. For each run, append the them in the reversal order of how they would be appended to the path in the regular DFS (this is essentially picking a vertex with the lowest in-degrees)
2. Iterating through the stack, run the DFS on a transposed graph
3. The set of visited vertex from this run will form a strongly connected component

```

def kosaraju(self):
    # Run DFS to make a stack of vertices based on exit time
    visited = [False] * self.n
    indegreePriorityStack = []

    def visit(v):
        visited[v] = True

        for w in self.adjList[v]:
            if not visited[w]:
                visit(w)
        indegreePriorityStack.append(v)

    for v in range(self.n):
        if not visited[v]:
            visit(v)

    # DFS on the transposed graph, in the order of exit time
    transposedG = Graph(self.n, [[v, u] for u, v in self.edges])
    visitedTranspose = [False] * self.n
    scc = []
    while indegreePriorityStack:
        v = indegreePriorityStack.pop()
        if not visitedTranspose[v]:
            path = []
            transposedG.dfs(v, visitedTranspose, path)
            scc.append(sorted(path))

    return scc

```

## 7.2 Articulation Points

An articulation point is a vertex such that removing it from the graph increases the number of connected components. This is applicable for undirected graph as well.

---

```

1: function AP( $G, s$ , discovery)
2: | Mark  $s$  as visited
3: | discovery[ $s$ ]  $\leftarrow d$ 

```

$\triangleright s$  is the starting vertex

---

---

```

4:   low[s] ← d
5:   for w adj to s do
6:       if v is not visited then AP(g, s, d)
7:       low[s] ← MIN(low[s], low[w])
8:       if low[s] ≥ discovery[s] then
9:           if s is not root or s has more than 1 child in the DFS tree then
10:  PRINT(s is the AP)

```

---

## 7.3 Directed Acyclic Graphs

DAG is a directed graph without a cycle.

### 7.3.1 Topological Sort

---

```

1: function KAHN-TOPOLOGICAL-SORT(G) ▷ G is the graph object
2:   H ← a specialized graph able to keep track of in-degrees of vertices, initialized with G
3:   Ordering ← array of size |V|
4:   n ← 0
5:   while H is not empty do
6:       v ← a vertex with in-degree 0 from H
7:       Ordering[n] = v
8:       P.remove(v and all connect edges)
9:       n ← n + 1

```

---

If *G* is not a DAG (i.e., contains cycle(s)), *H* will still have edge(s) left in it because vertices in a cycle will never reach in-degree 0. This means that you can check if a graph is DAG by running the topological sort and asserting *Ordering*.length = |*V*|

Here is an actual implementation of the algorithm using queue.

```

def kahnTopologicalSort(self):
    indeg = [0] * self.n

    for u, v in self.edges:
        indeg[v] += 1

    q = deque()
    for v in range(self.n):
        if indeg[v] == 0:
            q.append(v)

    ordering = []
    while q:
        u = q.pop()
        ordering.append(u)

        for w in self.adjList[u]:
            indeg[w] -= 1 # "removing" the edge
            if indeg[w] == 0:
                q.append(w)

```

```
# if we do not have all the nodes in the ordering list,  
# it means that the graph is not DAG.  
# It is because if the graph is DAG, all nodes will eventually have  
# zero indegree after edge removal  
return [] if len(ordering) != self.n else ordering
```

## Chapter 8

# Weighted Graphs

The constructor for weighted undirected graph takes an edge list with an extra information, namely the weight of the edge. When building an adjacency list/matrix, a tuple containing the destination vertex and the weight is pushed.

```
class Graph:
    # e.g.,: Graph(n = 3, edges = [[0,1,3],[1,2,1],[0,2,6]])
    # creates a graph with 3 nodes indexed 0, 1, 2,
    # with 0 - 1 (weight 3), 0 - 2 (weight 6), 1 - 2 (weight 1)
    def __init__(self, n: int, edges: list[list[int]]):
        self.n = n
        self.edges = edges

        self.adjList = [[] for _ in range(n)]
        for u, v, weight in edges:
            self.adjList[u].append((v, weight))
            self.adjList[v].append((u, weight))

        self.adjMatrix = [[float("inf") for _ in range(n)] for _ in range(n)]
        for i in range(n):
            self.adjMatrix[i][i] = 0
        for u, v, w in edges:
            self.adjMatrix[u][v] = w
            self.adjMatrix[v][u] = w
```

### 8.1 Shortest Path

#### 8.1.1 Initialization and Edge Relaxation for Single-source Shortest-path Algorithms

---

<pre>1: <b>function</b> INITIALIZE-SINGLE-SOURCE(<math>G, s</math>) 2:   <math>dist \leftarrow</math> array size <math> V </math> 3:   <math>prev \leftarrow</math> array size <math> V </math> 4:   <b>for</b> <math>v \in V</math> <b>do</b> 5:     <b>if</b> <math>v = s</math> <b>then</b> <math>dist[v] \leftarrow 0</math> 6:     <b>if</b> <math>v \neq s</math> <b>then</b> <math>dist[v] \leftarrow \infty</math> 7:   <math>prev[u] \leftarrow nil</math>    <b>return</b> <math>dist, prev</math></pre>	<p><math>\triangleright G</math> is the graph, <math>s</math> is the starting vertex</p> <p><math>\triangleright</math> or <math>-1</math></p>
--	--

---

```

1: function RELAX( $u, v$ )
2:    $d \leftarrow \text{dist}[u] + \text{weight}(u, v)$ 
3:   if  $\text{dist}[v] > d$  then
4:      $\text{dist}[v] = d$ 
5:      $\text{prev}[v] = u$ 

```

Simply put, relaxation is a process of greedily updating the `DIST` value of a vertex  $v$  to the shorter of the current `DIST[v]` or newly discovered edge. These two are easier to understand in action.

### 8.1.2 Bellman-Ford Algorithm

Bellman-Ford algorithm computes shortest-path by performing relaxation for every edge  $|V| - 1$  times (maximum number of edges in a simple path). Then it attempts to relax each edge once more, and any additional relaxation would indicate the existence of a negative weight cycle (cycle involving an edge with a negative weight; one can travel this over and over to reduce the cost of path, rendering shortest path meaningless).

```

1: function BELLMAN-FORD-SHORTEST-PATH( $G, V$ )                                ▷  $G$  is the graph,  $V$  is the vertex list
2:   ▷ Initialize-Single-Source:  $O(|V|)$                                        ◁
3:    $\text{dist} \leftarrow \text{array size } |V|$ 
4:    $\text{prev} \leftarrow \text{array size } |V|$ 
5:   for  $v \in V$  do
6:     if  $v = s$  then  $\text{dist}[v] \leftarrow 0$ 
7:     if  $v \neq s$  then  $\text{dist}[v] \leftarrow \infty$ 
8:      $\text{prev}[u] \leftarrow -1$ 
9:   ▷ Visiting each edges and relaxing:  $O(|V||E|)$                              ◁
10:  for  $i = 1$  to  $|V| - 1$  do
11:    for  $e \in E$  do                                                         ▷ Edge  $e$  connects vertex  $u$  and  $v$ 
12:      if  $\text{weight}[e] + \text{dist}[u] < \text{dist}[v]$  then
13:         $\text{dist}[v] = \text{dist}[u] + \text{weight}[e]$ 
14:         $\text{prev}[v] = u$ 
15:      if no relaxation happened in the inner loop then
16:        ▷ This means that every edges are relaxed to its shortest path      ◁
17:        return true,  $\text{dist}$ ,  $\text{prev}$ 
18:  ▷ Visiting each edge to check for negative weight cycle:  $O(|E|)$            ◁
19:  for  $e \in E$  do                                                         ▷ Edge  $e$  connects vertex  $u$  and  $v$ 
20:    if  $\text{weight}[e] + \text{dist}[u] < \text{dist}[v]$  then
21:      output Negative weight edge cycle detected
22:    return false
23:  return true,  $\text{dist}$ ,  $\text{prev}$ 

```

### 8.1.3 Dijkstra's Algorithm

Dijkstra is much like BFS but instead of regular queue, it uses min-heap priority queue to visit a vertex with minimum cost. The following pseudocode is a textbook implementation of Dijkstra's algorithm, using a min-heap priority queue that supports modifying priorities of an element after they have been enqueued.

```

1: function DIJKSTRA( $G, s$ )                                                  ▷  $s$  is the starting vertex
2:   ▷ Initialize-Single-Source:  $O(|V|)$                                        ◁
3:    $\text{dist} \leftarrow \text{array size } |V|$ 

```

```

4:  prev ← array size |V|
5:  for  $v \in V$  do
6:      if  $v = s$  then  $\text{dist}[v] \leftarrow 0$ 
7:      if  $v \neq s$  then  $\text{dist}[v] \leftarrow \infty$ 
8:       $\text{prev}[u] \leftarrow -1$ 

9:   $Q \leftarrow$  a min-heap priority queue
10: for all  $v \in G.V$  do
11:     ▷ This inserts  $(0, s)$  for the starting vertex  $s$  and  $(\infty, v)$  for every other vertices
12:      $Q.\text{ENQUEUE}(\text{dist}[v], v)$ 

13: while  $Q \neq \emptyset$  do
14:      $(\_, u) \leftarrow Q.\text{DEQUEUE}$ 
15:     for all  $v \in G.\text{ADJACENCY LIST OF } u$  do
16:         ▷ Edge relaxation
17:          $d \leftarrow \text{dist}[u] + G.\text{WEIGHT}((u, v))$ 
18:         if  $d < \text{dist}[v]$  then
19:              $\text{dist}[v] = d$ 
20:              $\text{prev}[v] = u$ 
21:              $Q.\text{SET-PRIORITY}((d, v))$ 

```

In real life, you seldom find an implementation of priority queue with SET-PRIORITY function. Instead, you can simply insert the new tuple into the queue when relaxing an edge.

```

def dijkstra(self, v):
    dist = [float("inf")] * self.n
    prev = [-1] * self.n

    dist[v] = 0
    q = [(0, v)]

    while q:
        priority, u = heappop(q)

        for w, weight in self.adjList[u]:
            candidate = dist[u] + weight
            if candidate < dist[w]:
                dist[w] = candidate
                prev[w] = u
                heappush(q, (candidate, w))

    return dist, prev

```

This version might cause duplicate vertices in the queue. For example, if a vertex  $v$  gets relaxed to the distance value of 9 by a vertex  $u$ , then it later gets relaxed to 6 by a vertex  $w$ , then there will be both  $(6, v)$  and  $(9, v)$  present in the queue. It will not cause any problem since by the time  $(9, v)$  is dequeued, vertices adjacent to  $v$  would have been relaxed by  $(6, v)$ .

To optimize this, we can check if the priority we dequeued matches the value in the `dist` array.

```

def dijkstra(self, v):
    dist = [float("inf")] * self.n
    prev = [-1] * self.n

    dist[v] = 0
    q = [(0, v)]

```

```

while q:
    priority, u = heappop(q)

    if priority == dist[u]: # or p <= dist[u]
        for w, weight in self.adjList[u]:
            candidate = dist[u] + weight
            if candidate < dist[w]:
                dist[w] = candidate
                prev[w] = u
                heappush(q, (candidate, w))

return dist, prev

```

Because we only enqueue when a vertex is relaxed,  $p < \text{dist}[u]$  would never happen unless there is a negative weight cycle.

### 8.1.4 Floyd-Warshall All Pairs Shortest Path Algorithm

Dijkstra and Bellman-Ford are examples of "single source" shortest path algorithm, computing shortest path from a given vertex to other nodes. Floyd-Warshall algorithm is an "all pairs" shortest path algorithm, computing the shortest path between all pairs of vertices. It uses an adjacency matrix to compute shortest paths in  $O(|V|^3)$  time.

```

def floydWarshall(self):
    # Making a copy of adjacency matrix
    dist = [row[:] for row in self.adjMatrix]

    for k in range(self.n):
        for i in range(self.n):
            for j in range(self.n):
                dist[i][j] = min(dist[i][j], dist[i][k] + dist[k][j])

    return dist

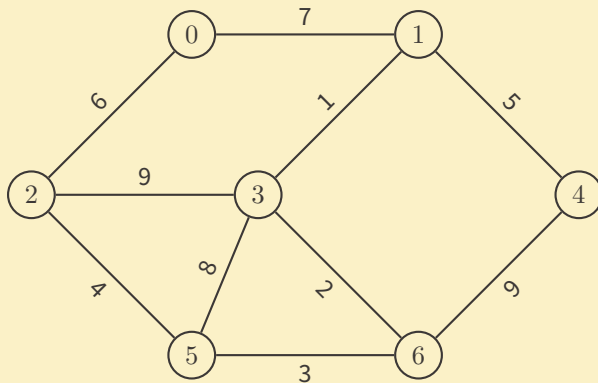
```

## 8.2 Minimum Spanning Tree

- Minimum:  $\sum \text{weight}$  is minimum
- Spanning: All vertices in the graph are connected
- Tree: No cycle

The subset of edges that connects all vertices in the minimum cost.

### 8.2.1 Cycle and Cut Properties



There are two fundamental properties of MST:

1. *Cycle Property*: For any cycle  $C$  in the graph, if the weight of an edge  $e \in C$  is higher than any of individual weights of all other edges in  $C$ , then its edge cannot belong in the MST.  
In the example above,  $(6, 4)$  in the cycle  $2 - 0 - 1 - 4 - 6 - 5$  and  $(2, 5)$  in the cycle  $2 - 3 - 5$  cannot be in the MST.
2. *Cut Property*: For any cut (subdivision of graph with disjoint)  $C$  in the graph, if the weight of an edge  $e$  in the cut-set of  $C$  is strictly smaller than the weights of all other edges of the cut-set of  $C$ , then this edge belongs to all MST of the graph.  
If we remove  $(0, 1)$ ,  $(2, 3)$ ,  $(2, 5)$ , then graph separates into  $(0, 1)$  and rest of the vertices. Thus, the lowest cost edge  $(2, 5)$  must be in the MST.

### 8.2.2 Prim's Algorithm

Prim's algorithm is essentially Dijkstra but for finding MST.

```
def primMST(self, s, visited):
    dist = [float("inf")] * self.n
    prev = [-1] * self.n

    dist[s] = 0
    q = [(0, s)]

    cost = 0
    mstEdges = []

    while q:
        priority, u = heappop(q)

        # Ensure that the dequeued vertex has not been relaxed yet.
        # For the algorithm with implementation without set_priority method
        # See Dijkstra section for more information.
        if priority == dist[u]:
            # update the MST information
            visited[u] = True
            cost += priority
            if prev[u] != -1:
                mstEdges.append((prev[u], u, priority))

        for v, weight in self.adjList[u]:
```



```

        if not visited[v]:
            if weight < dist[v]:
                prev[v] = u
                dist[v] = weight
                heappush(q, (weight, v))

    return cost, mstEdges

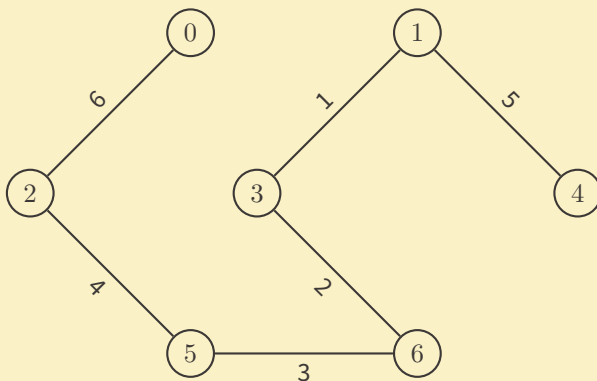
```

Running the algorithm on the graph in the section 8.2.1, we get the following MST.

```

g = Graph(7, [
    [0, 1, 7],
    [0, 2, 6],
    [1, 4, 5],
    [1, 3, 1],
    [2, 3, 9],
    [2, 5, 4],
    [3, 5, 8],
    [3, 6, 2],
    [4, 6, 9],
    [5, 6, 3]
])
visited = [False] * g.n
print(g.primMST(0, visited))
...
(21, [(0, 2, 6), (2, 5, 4), (5, 6, 3), (6, 3, 2), (3, 1, 1), (1, 4, 5)])

```



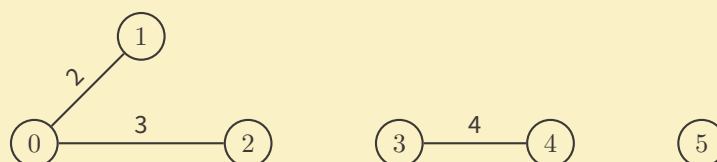
The Prim's MST algorithm runs in  $O(|E| \log(|V|))$ .

Assuming the graph is fully connected, the starting node  $s$  can be any node. Since the algorithm traverses on connected edges, running the algorithm on a graph with disconnected components will result in an incomplete MST.

```

g = Graph(6, [
    [0, 1, 2],
    [0, 2, 3],
    [3, 4, 4]
])
# Node 5 is completely disconnected

```



```

visited = [False] * g.n
print(g.primMST(0, visited))
...
(5, [(0, 1, 2), (0, 2, 3)])

```



You can introduce a loop to iterate all vertices to ensure that the correct MST is formed, or you can use Kruskal MST algorithm.

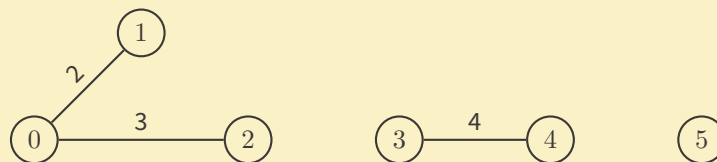
```

def primAllVertices(self):
    cost = 0
    mstEdges = []
    visited = [False] * self.n
    for v in range(self.n):
        if not visited[v]:
            localCost, localMstEdges = self.primMST(v, visited)
            cost += localCost
            mstEdges.extend(localMstEdges)

    return cost, mstEdges

print(g.primAllVertices)
...
(9, [(0, 1, 2), (0, 2, 3), (3, 4, 4)])

```



### 8.2.3 Union-Find

Union-Find, or disjoint-set data structure, store sets in terms of their relations to common elements.

```

class UnionFind:
    def __init__(self, n):
        # Initially, each node is parent to itself
        self.parent = list(range(n))

    def findBasic(self, p):
        if self.parent[p] == p:
            return p
        return self.findBasic(self.parent[p])

    def unionBasic(self, p, q):
        rootP = self.find(p)
        rootQ = self.find(q)
        self.parent[rootP] = rootQ

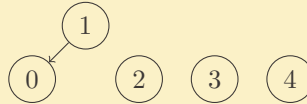
```

The data structure assigns the root/representative of each disjoint set to distinguish each set.

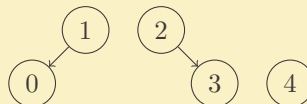
- In the beginning, each element is a separate set.



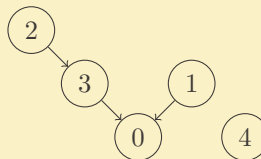
- UNION(1, 0): the root of each element is themselves, so 0 becomes the parent of 1



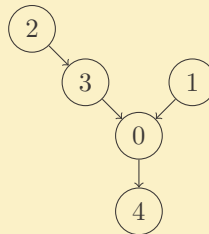
- UNION(2, 3)



- UNION(2, 1): the root of 2 is 3, the root of 1 is 0, so 0 becomes the parent of 3



- UNION(0, 4)



There are two ways to optimize Union-Find data structure.

1. Path Compression (improving find)
2. Union by rank or by size (improving UNION)

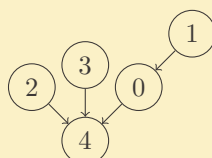
### Path Compression

```

def findPathCompression(self, p):
    if self.parent[p] != p:
        # path compression
        self.parent[p] = self.findPathCompression(self.parent[p])
    return self.parent[p]
  
```

Currently, when find is called on 2, it has to travel multiple elements to get to its root, 4. Path compression optimizes the process by "flattening" the tree during the find process.

After find(2), the tree becomes:



## Union by Rank

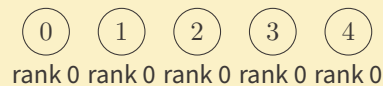
When union-ing two elements, union by rank (height of the tree if path compression was not used) chooses to append the root node with shorter height to the taller one to prevent total tree height from increasing.

```
class UnionFind:
    def __init__(self, n):
        # Initially, each node is parent to itself
        self.parent = list(range(n))
        self.rank = [0] * n

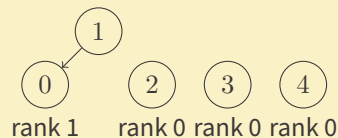
    def unionByRank(self, p, q):
        rootP = self.find(p)
        rootQ = self.find(q)
        if rootP == rootQ:
            return

        if self.rank[rootP] < self.rank[rootQ]:
            self.parent[rootP] = rootQ
        elif self.rank[rootP] > self.rank[rootQ]:
            self.parent[rootQ] = rootP
        else:
            # you can switch P and Q if you wish
            self.parent[rootP] = rootQ
            self.rank[rootQ] += 1
```

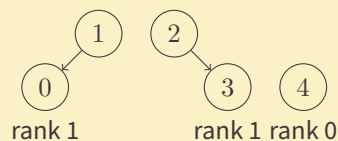
- In the beginning, each element is a separate set with rank 0



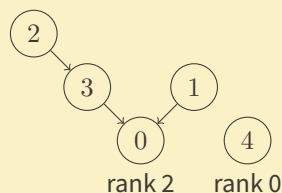
- UNION(1, 0): the root of each element is themselves with equal ranks, we choose 0 to be the root of 1 and increase the rank of 0



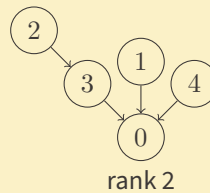
- UNION(2, 3)



- UNION(2, 1): find(2) = 3 and find(1) = 0. Since they are equal in rank, we arbitrarily choose 0 to be the parent of 3



- `UNION(0, 4)`: Since 4 has lower rank than 0 (the result of `FIND(0)`), 0 becomes the parent of 4.



### Union by Size

I am too lazy to cover it, but here is my Python implementation. It is essentially union by rank, but it compares the sizes of the trees instead of their heights.

```

class UnionFind:
    def __init__(self, n):
        # Initially, each node is parent to itself
        self.parent = list(range(n))
        self.size = [1] * n

    def unionBySize(self, p, q):
        rootP = self.find(p)
        rootQ = self.find(q)
        if rootP == rootQ:
            return
        if self.size[rootP] < self.size[rootQ]:
            self.parent[rootP] = rootQ
            self.size[rootQ] += self.size[rootP]
        else:
            self.parent[rootQ] = rootP
            self.size[rootP] += self.size[rootQ]

```

### 8.2.4 Kruskal MST Algorithm

Kruskal's algorithm uses Union-Find data structure to find MST.

```

from unionFind import UnionFind

class Graph:
    # ...

    def kruskalMST(self):
        # Sort edges based on their weights
        sortedEdges = sorted(self.edges, key=lambda edge: edge[2])

        uf = UnionFind(self.n)
        cost = 0
        mstEdges = []
        for u, v, weight in sortedEdges:
            if uf.findPathCompression(u) != uf.findPathCompression(v):
                mstEdges.append((u, v, weight))
                cost += weight
                uf.unionByRank(u, v)

        return cost, mstEdges

```

Kruskal MST algorithm tests whether adding the edge to the MST tree will create a cycle by checking the respective roots in the Union-Find data structure. If each end vertices of the edge has the same root in the Union-Find, they both belong to the MST tree already. Since we iterate edges in sorted order, it means that more cost-efficient edge has been added to the MST tree.

The algorithm runs in  $O(|E| \log |E|)$  time, and as it traverses the list of sorted edges, disconnected components will be considered when MST is built.

## Chapter 9

# Strings

For a string  $T$ , substring of  $T$  is denoted by  $T[i, j]$  where  $0 \leq i \leq j < M$ .

**Prefix** Substring from 0 to  $i$  where  $i < M$

**Suffix** Substring from  $i$  to  $M - 1$  where  $i \geq 0$

### 9.1 Brute-force String Pattern Matching

Suppose we want to check if a string  $P$  (stands for "pattern") is a substring of a longer string  $T$ . If you are lazy, you can use the following brute-force pattern-matching algorithm.

```
def bruteForceMatch(T, P): # P is the "pattern" string we wish to find in T
    n, m = len(T), len(P)
    for i in range(n - m):
        j = 0
        while j < m and T[i + j] == P[j]:
            j += 1
        if j == m:
            return i
    return -1
```

The algorithm goes through  $T$  and checks compares characters of  $P$  until we reach the end  $P$ , returning the earliest match. The algorithm is highly inefficient, taking  $O(nm)$  time in the worst case scenario.

T	h	e	o
e			
	e		
		e	o

### 9.2 KMP Algorithm

Knuth-Morris-Pratt algorithm is an efficient pattern matching algorithm with  $O(n + m)$  time complexity. It uses the pre-computed LPS (Longest Proper Prefix, where proper prefix is a prefix of  $T$  that does not include  $T$  itself) array. The name LPS is slightly misleading, as what FAILURE function is actually calculating is **the length of the longest matching prefix & suffix for each substring**.

```

def failure(pattern):
    m = len(pattern)
    lps = [0] * m
    j = 0
    for i in range(1, m):
        while j > 0 and pattern[i] != pattern[j]:
            j = lps[j - 1]
        if pattern[i] == pattern[j]:
            j += 1
        lps[i] = j

    return lps

```

For example, for the string "ABABCAB", the FAILURE function produces the following result.

<code>lps = failure("ababcab")</code>	
<code>lps[0] = 0</code>	( <code>"a"</code> )
<code>lps[1] = 0</code>	( <code>"ab"</code> )
<code>lps[2] = 1</code>	( <code>"ābā"</code> )
<code>lps[3] = 2</code>	( <code>"āb āb"</code> )
<code>lps[4] = 0</code>	( <code>"ababc"</code> )
<code>lps[5] = 1</code>	( <code>"ābabcā"</code> )
<code>lps[6] = 2</code>	( <code>"ābabcāb"</code> )

The KMP driver function uses the LPS array to find the earliest index where the pattern is found.

```

def kmp(T, pattern):
    n, m = len(T), len(pattern)
    lps = failure(T, pattern)
    i, j = 0, 0
    while i < n:
        if T[i] == pattern[j]:
            if j == m - 1:
                return i - j
            else:
                i += 1
                j += 1
        elif j > 0:
            j = lps[j - 1]
        else:
            i += 1
    return -1

```

a	b	a	b	c	a	a	b	a	c	c	a	b	a	b	c	a	b
a	b	a	b	c	a	b											

Since the length of longest prefix/suffix is 1 ("A"), start matching at index 1.

					a	b											
--	--	--	--	--	---	---	--	--	--	--	--	--	--	--	--	--	--

Since the length of longest prefix/suffix is 0, start matching at index 0.

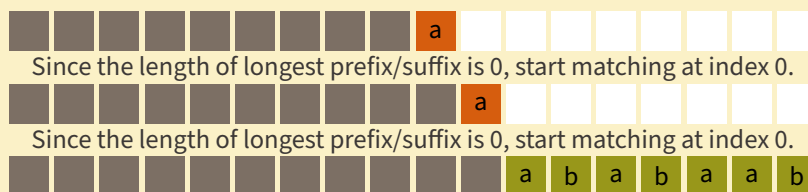
						a	b	a	b								
--	--	--	--	--	--	---	---	---	---	--	--	--	--	--	--	--	--

Since the length of longest prefix/suffix is 1 ("A"), start matching at index 1.

								a	b								
--	--	--	--	--	--	--	--	---	---	--	--	--	--	--	--	--	--

Since the length of longest prefix/suffix is 0, start matching at index 0.



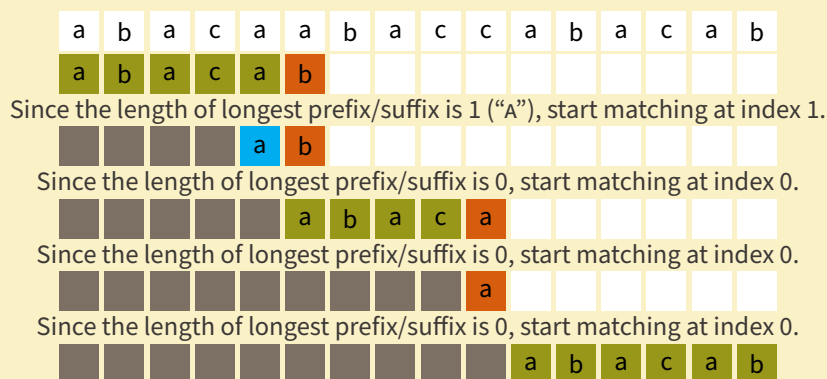


Okay, kind of a bad example, here is a simpler example with  $T = abacaabaccabacab$  and  $P = abacab$ .

```

lps = failure("abacab")
lps[0] = 0                                ("a")
lps[1] = 0                                ("ab")
lps[2] = 1                                ("ābā")
lps[3] = 0                                ("abac")
lps[4] = 1                                ("ābaca")
lps[5] = 2                                ("ābacāb")

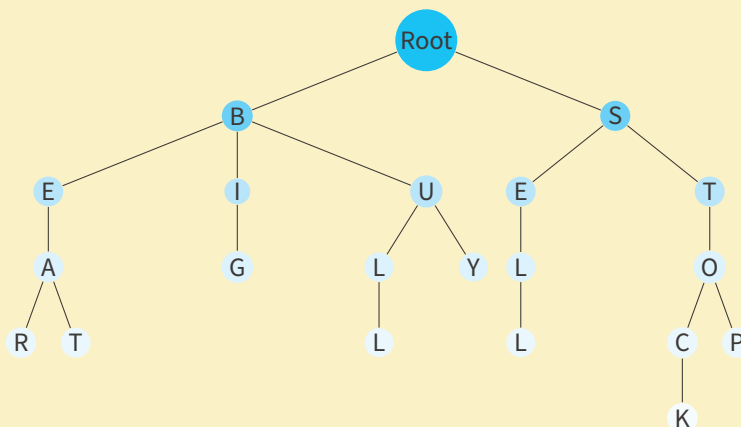
```



Okay, that was also not a good example, but hope you get the gist; when a failure in matching occurs, KMP algorithm skips searching the part where the prefix and suffix is known to match.

## 9.3 Trie

A trie, or a prefix tree is a specialized search tree for storing characters of strings.



Trie with the following words: BEAR, BEAT, BIG, BULL, BUY, SELL, STOCK, STOP

Following is the implementation of Trie in Python and the basic visualization method.

```
class TrieNode:
    def __init__(self):
        self.children = {}
        self.end = False
        self.freq = 0

class Trie:
    def __init__(self):
        self.root = TrieNode()

    def insert(self, word):
        """
        Inserts the given string in the Trie
        :param word: string to be inserted
        """
        node = self.root
        for c in word:
            if c not in node.children:
                node.children[c] = TrieNode()
            node = node.children[c]
            node.freq += 1
        node.end = True

    def visualize(self):
        def dfs(node, prefix="", indent=""):
            for char, child in sorted(node.children.items()):
                marker = "*" if child.end else ""
                print(f"{indent}-{char}-{marker}")
                dfs(child, prefix + char, indent + " ")

        print("ROOT")
        dfs(self.root)
```

The visualization, unfortunately, is not in the traditional top-down tree form, as the implementation of such visualization is vastly more complicated. Instead, the Trie above would be shown as the following:

```
trie = Trie()
trie.insert("BEAR")
trie.insert("BEAT")
trie.insert("BIG")
trie.insert("BULL")
trie.insert("BUY")
trie.insert("SELL")
trie.insert("STOCK")
trie.insert("STOP")
trie.visualize()

...
B
E
A
R*
T*
```

```

I
  G*
    U
      L
        L*
          Y*
S
  E
    L
      L*
T
  O
    C
      K*
        P*

```

With each indentation level indicating the level of the tree.

If you want to check if a word exists in the Trie, you can use the following method.

```

def searchAsWholeWord(self, word):
    """
    Checks if the given string exists in the Trie as a whole word
    :param word: string to be searched
    :return: Boolean indicating whether the word exists or not
    """
    node = self.root
    for c in word:
        if c not in node.children:
            return False
        node = node.children[c]
    return node.end

```

This, however, is not where Trie truly shines. Trie is very useful with anything involving prefixes. Consider the following two methods.

```

def searchFreqAsPrefix(self, pfx):
    """
    Checks how frequent the given pattern appears as prefixes
    in other words in Trie
    :param pfx: string to be searched
    :return: Integer denoting the frequency of the parameter
    """
    node = self.root
    for c in pfx:
        if c not in node.children:
            return 0
        node = node.children[c]
    return node.freq

def getLongestCommonPrefix(self):
    longestPfx, longestPfxLen = "", 0
    node = self.root
    # If node has more than one child, it means the trie diverges.
    # Therefore the common prefix ends there.
    while not node.end and len(node.children) <= 1:
        c = next(iter(node.children)) # node's only child

```

```

    longestPfxLen += 1
    longestPfx += c
    node = node.children[c]
    return longestPfx, longestPfxLen

```

SEARCHFREQASPREFIX checks how frequently the given string is used as prefixes of other words, and GETLONGESTCOMMONPREFIX returns the longest common prefix of the words currently in the Trie.

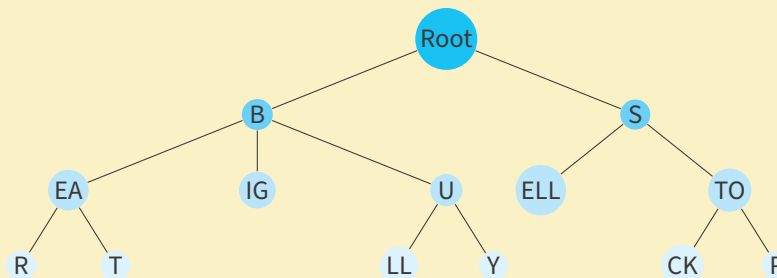


The longest common prefix in this Trie is "BEA".

Trie is used in operations like auto-completion and spell check. Instead of comparing the misspelled word to every words in the spell check dictionary, Trie only requires a few traversals of the nodes. In many spell checking algorithms, after morphological processing to handle grammatical variations, hash table is first used to check if a word exists in the dictionary. Then, the word is searched in the Trie, and possible candidates are sorted based on the frequency in the language and Levenshtein distance.

## 9.4 PATRICIA

PATRICIA, or radix tree is a version of Trie where node can store words and not just characters.



PATRICIA with the following words: BEAR, BEAT, BIG, BULL, BUY, SELL, STOCK, STOP

## 9.5 Huffman Coding

Huffman Coding is a technique for lossless data compression. The algorithm uses a specialized binary tree called Huffman tree to assign shorter "prefix code" more frequently used characters.

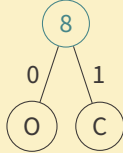
Supposed you have a data with the following frequency:

Char	E	D	C	O	A	G
Frequency	17	10	5	3	15	6

First, pick two nodes with the lowest frequencies and merge them into a binary tree, with the root node indicating the sum of the frequency. When forming a subtree, assign 0 to the edge to left node and 1 to the edge to the right node. Push the root back into the priority queue, pick two nodes or subtrees with the lowest frequencies, and merge them into a tree. Repeat this process.

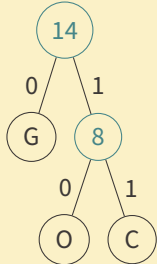
1. Initially,  $Q = \{(O, 3), (C, 5), (G, 6), (D, 10), (A, 15), (E, 17)\}$ .

Merging O and C:



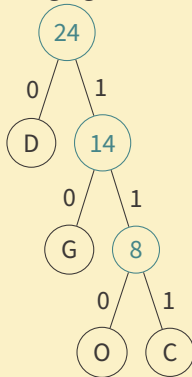
2.  $Q = \{(G, 6), (\text{Subtree with } (O, C), 8), (D, 10), (A, 15), (E, 17)\}$ .

Merging G and the subtree with O, C:



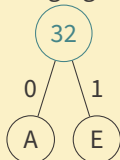
3.  $Q = \{(D, 10), (\text{Subtree with } (O, C, G), 14), (A, 15), (E, 17)\}$ .

Merging D and the subtree with O, C, G:



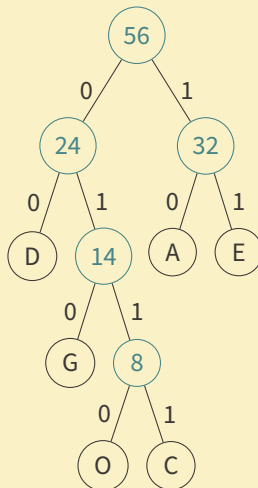
4.  $Q = \{(A, 15), (E, 17), (\text{Subtree with } (O, C, G, D), 24)\}$ .

Merging A and E:



5.  $Q = \{(\text{Subtree with } (O, C, G, D), 24), (\text{Subtree with } (A, E), 32)\}$ .

Merging two subtrees:



To determine the prefix code for each character, follow the 0's and 1's assigned to the edge.

Char	E	D	C	O	A	G
Prefix Code	11	00	0111	0110	10	010

To encode a string, say “GOOD” simply replace the characters with prefix code: “010 0110 0110 00” (spaces are added for the readability).

Here is the implementation in Python:

```

import heapq
from collections import Counter

class HuffNode:
    def __init__(self, character=None, frequency=0):
        self.character = character
        self.frequency = frequency
        self.left = None
        self.right = None

    def __gt__(self, other):
        return self.frequency > other.frequency

    def __str__(self):
        return f"{self.character}: {self.frequency}"

def buildHuffmanTree(freq: dict):
    """
    Given the dictionary of frequency of each character, build a Huffman Tree
    :return: dict of character to respective code generated by Huffman coding
    """
    min_heap = [HuffNode(char, f) for char, f in freq.items()]
    heapq.heapify(min_heap)

    while len(min_heap) > 1:
        left = heapq.heappop(min_heap)
        right = heapq.heappop(min_heap)

```

```

        mergedSubtree = HuffNode(frequency=left.frequency + right.frequency)
        mergedSubtree.left = left
        mergedSubtree.right = right

        heapq.heappush(min_heap, mergedSubtree)

    return min_heap[0]  # root

def getCodeMap(root: HuffNode):
    """
    Given the root of the Huffman Tree, create a dictionary mapping
    character to the respective code generated by the tree.
    Left nodes get 0, right nodes get 1.
    :param root: root of the Huffman Tree, generated by buildHuffmanTree()
    :return: dict of character to respective code generated by Huffman coding
    """
    codeMap = {}

    def preorder(node, code=""):
        if node is None:
            return

        if node.character is not None:
            codeMap[node.character] = code
            preorder(node.left, code=code + '0')
            preorder(node.right, code=code + '1')

    preorder(root)
    return codeMap

def encode(pattern, codeMap):
    """
    Given a string and a Huffman dictionary, encode the string
    :param pattern: String to code
    :param codeMap: dict mapping character to code, generated by getCodeMap()
    :return: coded string
    """
    return "".join([codeMap[c] for c in pattern])

strr = "EEEEEEEEEEEEEEEEEDDDDDDDDDCCCCOOOAAAAAAAAAAAAAGGGGGG"
root = buildHuffmanTree(Counter(strr))
print(encode("GOOD", getCodeMap(root)))

```

[Bonus] To decode a Huffman string, we can traverse the Huffman tree until we encounter a leaf node.

```

def decode(pattern, root):
    """
    Given Huffman string and a Huffman tree, decode the string
    :param pattern: Huffman string to decode
    :param root: Root of the respective Huffman tree
    :return: decoded string
    """

```

```

n = len(pattern)

curr = root
res = ""
for i in range(n):
    if pattern[i] == '0':
        curr = curr.left
    else:
        curr = curr.right

    # Reached a leaf node
    # 1. if curr.character is not None
    # 2. if curr.frequency is None
    # 3. if curr.left is None and curr.right is None
    # above three are all logically same in this context,
    # as nodes in Huffman tree contains either frequency or character
    # and leaf nodes are always the ones that contain the character
    if curr.character is not None:
        res += curr.character
        curr = root # reset to the root

return res

root = buildHuffmanTree({"E": 17, "D": 10, "C": 5, "O": 3, "A": 15, "G": 6})
print(decode("000110010", root))
# ...
"DOG" # !!

```