

1. Data Structures

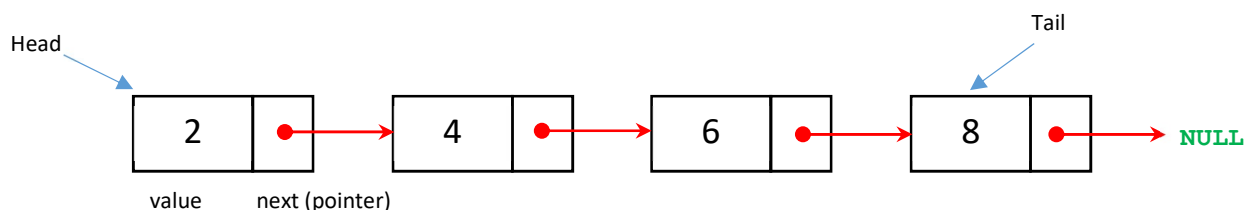
A data structure is a particular way of organizing data in a computer so that it can be used effectively. The idea is to reduce the space and time complexities of different tasks. Below is an overview of some popular linear data structures.

1.1 Linked List

Python's *list* class is highly optimized, and often a great choice for storage. With that said, there are some notable disadvantages:

- The length of a dynamic array might be longer than the actual number of elements that it stores.
- Amortized bounds for operations may be unacceptable in real-time systems.
- Insertions and deletions at interior positions of an array are expensive.

Linked lists provide an alternative to an array-based sequence (such as a Python list). Both array-based sequences and linked lists keep elements in a certain order, but using a very different style. An array provides the more centralized representation, with one large chunk of memory capable of accommodating references to many elements. A linked list, in contrast, relies on a more distributed representation in which a lightweight object, known as a node, is allocated for each element. Each node maintains a reference to its element and one or more references to neighboring nodes in order to collectively represent the linear order of the sequence. Unlike array-based sequences, linked list elements are not stored at contiguous location; the elements are linked using pointers.



Linked list is often compared to array-based sequences. Whereas an array is a fixed size of sequence, a linked list can have its elements to be dynamically allocated. Eventually, it comes with its pros and cons:

- PRO: A linked list saves memory since it only allocates the memory required for values to be stored. In arrays, you have to set an array size before filling it with values, which can potentially waste memory.
- PRO: Linked list nodes can live anywhere in the memory. Whereas an array requires a sequence of memory to be initiated, as long as the references are updated, each linked list node can be flexibly moved to a different address.

- **CON:** Linked lists have linear look up time. When looking for a value in a linked list, you have to start from the beginning of the chain, and check one element at a time for a value you are looking for. If the linked list is n elements long, this can take up to n time.

Below is the basic structure of a linked list in Python:

```
class Node:
    def __init__(self, value):
        self.value = value
        self.next = None
```



Initializing node object:

- Value assignments
- Initialize next as NULL

```
class LinkedList:
    def __init__(self):
        self.head = None
```



Initializing linked list object:

- Initialize head of list as NULL

Implementation of a linked list with 4 nodes:

```
class Node:
    def __init__(self, value):
        self.value = value
        self.next = None
```

```
class LinkedList:
    def __init__(self):
        self.head = None
```

```
    def printList(self):
        temp = self.head
        while (temp):
            print(temp.value, end=' ')
            temp = temp.next
```



printList() traverses the created list and prints the data of each node.

```
if __name__ == '__main__':
```

```
    linked_list = LinkedList()
```

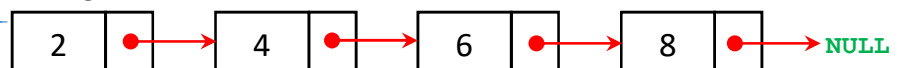
```
    linked_list.head = Node(2)
    node2 = Node(4)
    node3 = Node(6)
    node4 = Node(8)
```

Creating the 4 nodes:



```
    linked_list.head.next = node2
    node2.next = node3
    node3.next = node4
```

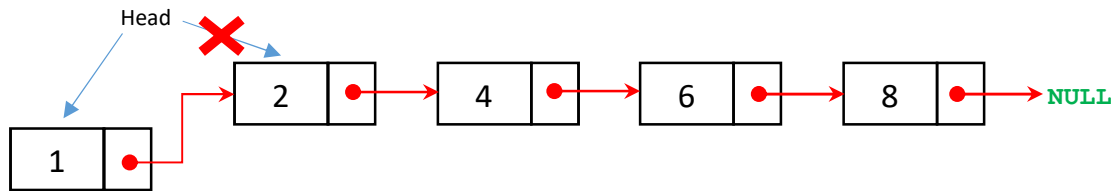
Linking nodes:



```
    linked_list.printList()
```

1.1.1 Inserting nodes

A node can be added to a linked list in three ways: (i) at the front of the linked list, (ii) after a given node or (iii) at the end of the linked list:

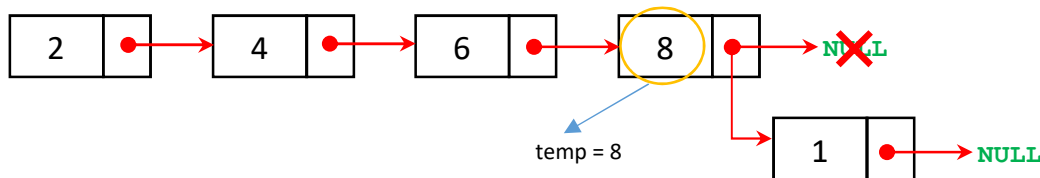


```
class LinkedList:
    def __init__(self):
        self.head = None

    def add_first(self, e):
        new_node = Node(e)
        new_node.next = self.head
        self.head = new_node
```

Annotations for `add_first`:

- `new_node = Node(e)`: create new node instance storing reference to element `e`
- `new_node.next = self.head`: set new node's next to reference the old head node
- `self.head = new_node`: set variable head to reference the new node



```
class LinkedList:
    def __init__(self):
        self.head = None

    def add_last(self, e):
        new_node = Node(e)
        new_node.next = None

        if self.head is None:
            self.head = new_node
            return

        last = self.head
        while (last.next):
            last = last.next

        last.next = new_node
```

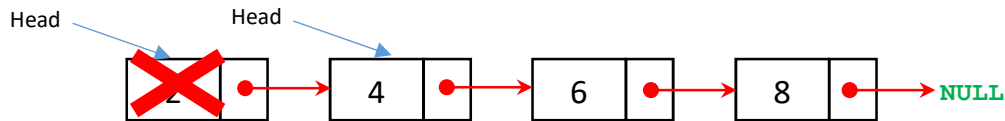
Annotations for `add_last`:

- `new_node = Node(e)`: create new node instance storing reference to element `e`
- `new_node.next = None`: set new node's next to reference the `None` object
- `if self.head is None: self.head = new_node; return`: if the linked list is empty, then make the `new_node` as head
- `last = self.head; while (last.next): last = last.next; last.next = new_node`: traverse list until the last node is reached, then change the next of last node to `new_node`

1.1.2 Deleting nodes

To delete a node from linked list, we need to do following steps:

- Find previous node of the node to be deleted
- Changed next of previous node
- Free memory for the node to be deleted (for C/C++)

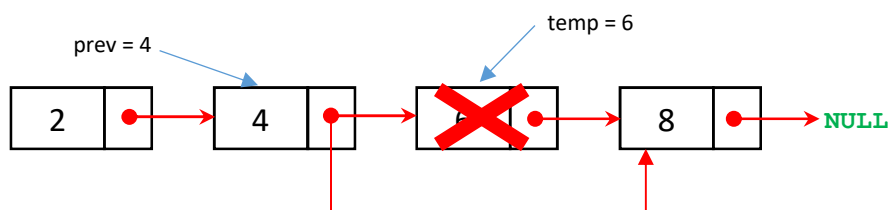


```
class LinkedList:
    def __init__(self):
        self.head = None

    def delete_first(self):
        if self.head is None:
            print("The list is empty")
            return
```

`self.head = self.head.next` → make head point to next node (or None)

Unfortunately, we cannot easily delete a given node or the last node of a singly linked list. Even if we maintain a tail reference directly to the last node of the list, we must be able to access the node before the last node in order to remove the last node. But we cannot reach the node before the tail by following next links from the tail. The only way to access this node is to start from the head of the list and search all the way through the list. But such a sequence of link-hopping operations could take a long time. If you want to support such an operation efficiently, you will need to make your list doubly linked.



```
class LinkedList:
    def __init__(self):
        self.head = None

    def delete_node(self, e):
        temp = self.head
```

```
        if (temp is not None):
            if (temp.value == e):
                self.head = temp.next
                temp = None
            return
```

if the head node itself holds the value to be deleted, make head point to next node and point temp to None

```

while(temp is not None):
    if temp.value == e:
        break
    prev = temp
    temp = temp.next

```

traverse the list looking for the value to be deleted, keeping track of the previous node since we need to change the node.next value

```

if(temp == None):
    return

```

key was not present in linked list

```

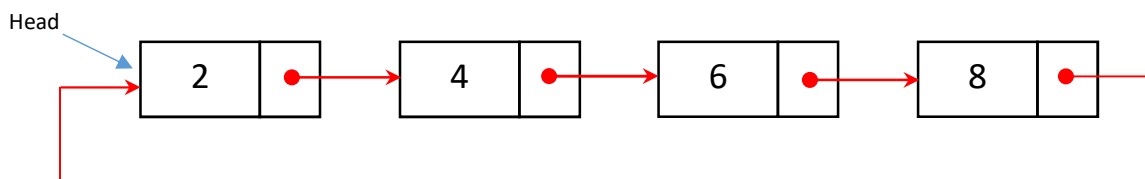
prev.next = temp.next
temp = None

```

unlinking the node from linked list

1.1.3 Circular Linked List

A circular linked list is a linked list where all nodes are connected to form a circle. There is no NULL at the end. A circular linked list can be a singly circular linked list or doubly circular linked list.



Some of the advantages of circular linked lists are:

- Any node can be a starting point. You can traverse the whole list by starting from any point. You just need to stop when the first visited node is visited again.
- Useful for implementation of queues. You do not need to maintain two pointers for front and rear if we use circular linked list. You can maintain a pointer to the last inserted node and front can always be obtained as next of last.
- Circular lists are useful in applications to repeatedly go around the list.

In a conventional linked list, you traverse the list from the head node and stop the traversal when you reach NULL. In a circular linked list, you stop traversal when we reach the first node again.

```

class Node:
    def __init__(self, value):
        self.value = value
        self.next = None

class CircularLinkedList:
    def __init__(self):
        self.head = None

    def add_first(self, e):
        new_node = Node(e)
        temp = self.head

        new_node.next = self.head

```



inserting a node at the beginning of the list:

- create a new_node
- make new_node.next = last_node.next
- make last_next = new_node
-

```

    if self.head is not None:
        while(temp.next != self.head):
            temp = temp.next
            temp.next = new_node
        else:
            new_node.next = new_node  → Only for the first node

    self.head = new_node

def printList(self):
    temp = self.head
    if self.head is not None:
        while (temp):
            print(temp.value, end=' ')
            temp = temp.next
            if (temp == self.head):
                break

if __name__ == '__main__':

    circular_list = CircularLinkedList()

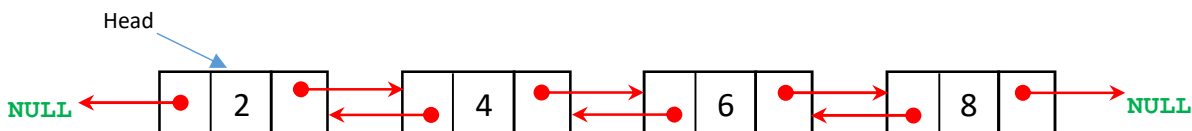
    circular_list.add_first(8)
    circular_list.add_first(6)
    circular_list.add_first(4)
    circular_list.add_first(2)

    circular_list.printList()    # 2 4 6 8

```

1.1.4 Doubly Linked List

A doubly linked list (DLL) contains an extra pointer, typically called *previous pointer*, together with next pointer and data which are there in singly linked list.



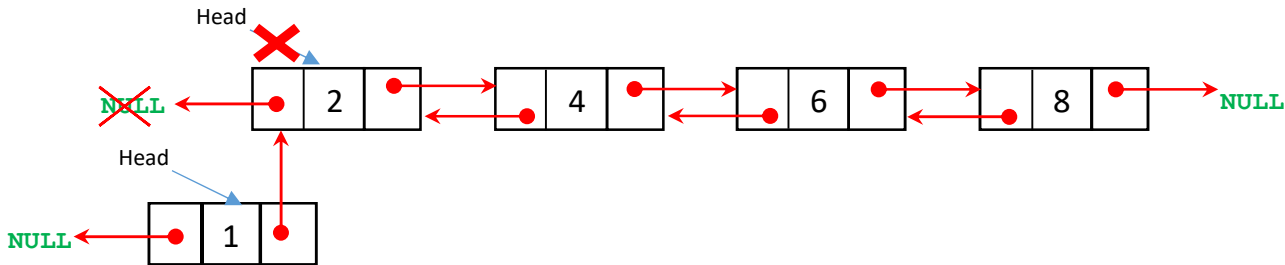
Some of the advantages of doubly linked lists over singly linked lists are:

- A doubly linked list can be traversed in both forward and backward direction.
- The delete operation in doubly linked list is more efficient if pointer to the node to be deleted is given. Remember that in singly linked list, to delete a node, pointer to the previous node is needed. To get this previous node, sometimes the list is traversed. In doubly linked list, we can get the previous node using previous pointer.

Some of the disadvantages of doubly linked lists over singly linked lists are:

- Every node of doubly linked list requires extra space for a previous pointer.
- All operations require an extra pointer previous to be maintained. For example, in insertion, we need to modify previous pointers together with next pointers.

A node can be added to a doubly linked list in four ways: (i) at the front of the list, (ii) after a given node, (iii) at the end of the list, (iv) before a given node.



```
class Node:
    def __init__(self, value):
        self.value = value
        self.next = None
        self.prev = None

class DoublyLinkedList:
    def __init__(self):
        self.head = None

    def add_first(self, e):
        new_node = Node(e)
        new_node.next = self.head  # Make next of new_node the head and previous as None (already None in the Node class)

        if self.head is not None:
            self.head.prev = new_node  # change previous of head node to the new_node

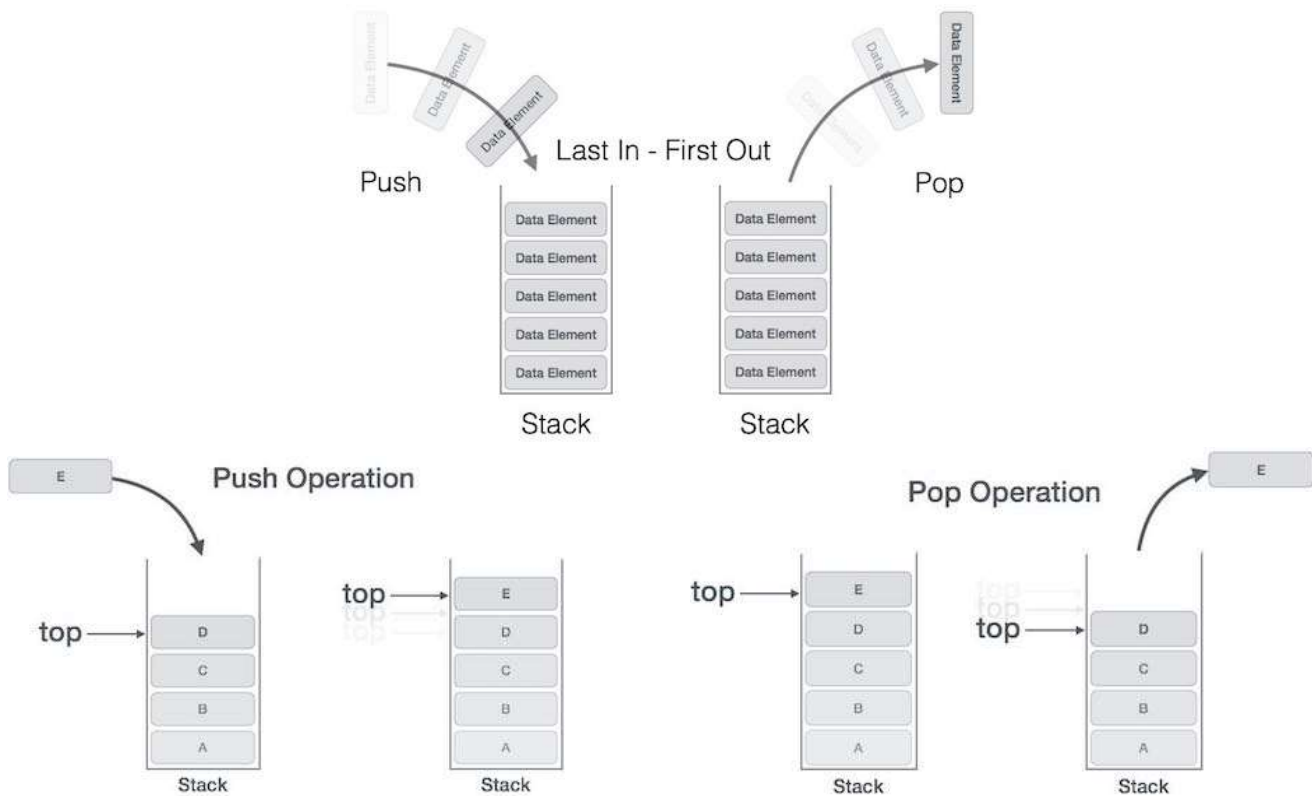
        self.head = new_node  # head points to the new_node
```

1.2 Stack

A stack is a linear data structure with three basic operations:

- **Push:** Adds an element on top of the stack. If the stack is full, then it is said to be an Overflow condition.
- **Pop:** Removes an element from the stack. The items are popped in the reversed order in which they are pushed. If the stack is empty, then it is said to be an Underflow condition.
- **isEmpty:** Returns true if stack is empty, else false.

One extra method available to a stack is the ability to “peek”. This method returns the top element of the stack without removing it. In a stack, elements are inserted and removed according to the LIFO (Last In, First Out) or FILO (First In, Last Out) principle. A user may insert objects into a stack at any time, but may only access or remove the most recently inserted object that remains (at the so-called “top” of the stack).



Stacks are useful to trace back previous elements/operations. For example, undo operations in editors are like popping a code change that was just pushed in the stack of edit history, or back operations in browsers are like popping a site visit that was just pushed in the stack of browser history. Stacks also come in handy when matching recursive elements/operation. For example, balancing of symbols, tree traversals, or the recursive algorithm to implement the Towers of Hanoi. Stacks can be implemented using array-based sequences or linked lists. The code below shows a linked list implementation of a stack:

```
class StackNode:
    def __init__(self, value):
        self.value = value
        self.next = None

class Stack:
    def __init__(self):
        self.stack = None

    def is_empty(self):
        return True if self.stack is None else False

    def push(self, e):
        newNode = StackNode(e)
        newNode.next = self.stack
        self.stack = newNode
        print("%d pushed in stack" %e)
```

Stack is empty when the root of the list is empty

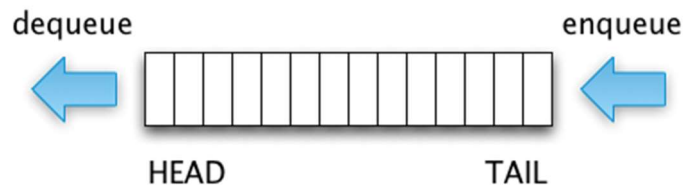
Although this implementation can grow and shrink according to the runtime needs, it requires extra memory due to the use of pointers.

1.3 Queue

Another fundamental data structure is the queue. Like stacks, queues are a linear structure which follows a particular order in which the operations are performed. In a queue, its elements are inserted and removed according to the FIFO (First In, First Out) or LILO (Last In, Last Out) principle. A queue has two basic operations:

- Enqueue: append an element to the tail of the queue
- Dequeue: remove an element from the head of the queue

One extra method available to a queue is to now who is at the head of the queue. This method returns a reference to the element at the front of queue without removing it.



Queues are used whenever you want to process things one at a time as they come in. Some examples are, uploading bunch of images, printing multiple documents, and processing thousands of requests to a web server. Queues can also be implemented using array-based sequences or linked lists. The code below shows an array-based implementation of a queue:

```
class Queue:
    def __init__(self):
        self.queue = []

    def enqueue(self, val):
        self.queue.append(val)

    def dequeue(self):
        if self.is_empty():
            return None
        else:
            return self.queue.pop(0)
```

1.3.1 Priority Queues

Priority Queue is an extension of queue with following properties:

- Every item has a priority associated with it.
- An element with high priority is dequeued before an element with low priority.
- If two elements have the same priority, they are served according to their order in the queue.

A typical priority queue supports following operations:

- insert(item, priority): Inserts an item with given priority.
- getHighestPriority(): Returns the highest priority item.
- deleteHighestPriority(): Removes the highest priority item.

A priority queue can implemented using array-based sequences or linked lists.

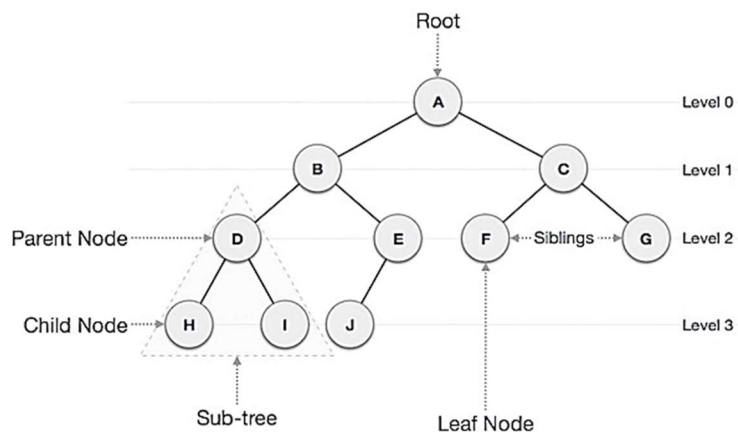
1.4 Tree

Unlike arrays, linked lists, stack and queues, which are linear data structures, trees are hierarchical data structures that are widely used, with a root value and subtrees of children with a parent node, represented as a set of linked nodes. Unlike trees in nature, the tree data structure is upside down: the root of the tree is on top. A tree consists of nodes and its connections are called edges. The bottom nodes are also named leaf nodes or external nodes. A tree does not have a cycle. Additionally:

- Nodes with the same parent are called siblings
- The depth of a node is the number of edges from the root to the node.
- The height of a node is the number of edges from the node to the deepest leaf.
- The height of a tree is a height of the root.

1.4.1 Binary trees

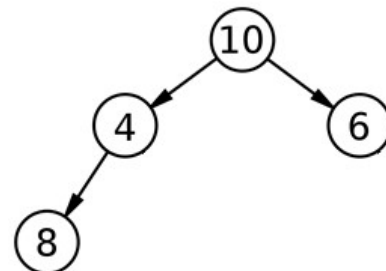
A binary tree is a data structure where every node has at most two children. The root of a tree is on top. Since each element in a binary tree can have only 2 children, we typically name them the left and right child. Since Python does not have built-in support for trees, you need to define a class tree which has a left and right attribute.



```
class Node:
    def __init__(self, key):
        self.left = None
        self.right = None
        self.value = key
```

```
root = Node(10)
root.left = Node(4)
root.right = Node(6)
root.left.left = Node(8)
```

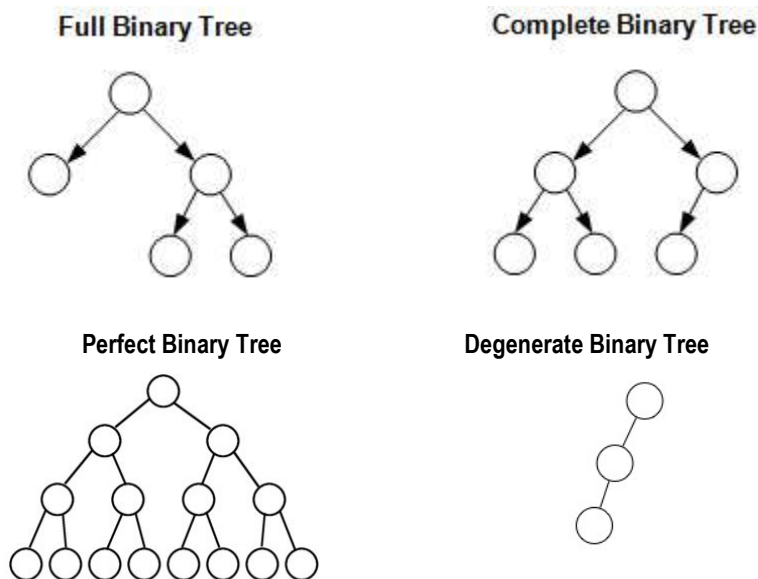
```
print(root.value)
print(root.left.value)
```



1.4.1.1 Types of binary trees

The most common types of binary trees are:

- Full binary tree: a binary tree in which each node is either a leaf or possesses exactly two child nodes.
- Complete binary tree: a binary tree which is completely filled, with the possible exception of the bottom level, which is filled from left to right.
- Perfect binary tree: a binary tree in which all internal nodes have two children and all leaves are at same level.
- Degenerate binary tree: a binary tree where every internal node has one child. Such trees are performance-wise same as a linked list.



1.4.1.2 Properties of binary trees

Binary trees have the following properties:

- A binary tree of n nodes has $n-1$ edges
- For every $k \geq 0$, there are no more than 2^k nodes in level k
- A binary tree with k levels has at most 2^{k-1} leaves
- The number of nodes on the last level is at most the sum of the number of nodes on all other levels plus 1
- A binary tree with k levels has no more than $2^k - 1$ nodes



Important terms

Below are important terms with respect to a tree:

- path: refers to the sequence of nodes along the edges of a tree
- root: the node at the top of the tree. There is only one root per tree and one path from the root node to any node
- parent: any node except the root node has one edge upward to a node called parent
- child: the node below a given node connected by its edge downward
- leaf: the node which does not have any child node
- subtree: represents the descendants of a node
- visiting: refers to checking the value of a node when control is on the node
- traversing: passing through nodes in a specific order
- levels: level of a node represents the generation of a node
- key: represents a value of a node based on which a search operation is to be carried out for a node

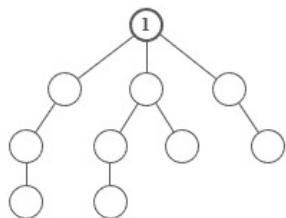
1.4.2 Traversal of binary trees

A traversal is a process that visits all the nodes in the tree. Since a tree is a nonlinear data structure, there is no unique traversal. A tree is typically traversed in two ways:

- Breadth-First Traversal (Or Level Order Traversal)
- Depth-First Traversals:
 - ✓ Inorder Traversal (Left- Root -Right)
 - ✓ Preorder Traversal (Root-Left-Right)
 - ✓ Postorder Traversal (Left-Right- Root)

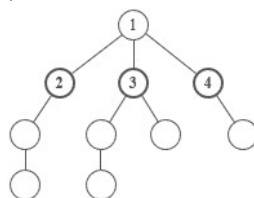
1.4.2.1 Breadth-First Search

Breadth-first search (BFS) is a method for exploring a tree or graph. In a BFS, you first explore all the nodes one step away, then all the nodes two steps away, etc. Here's how a BFS would traverse this tree, starting with the root:

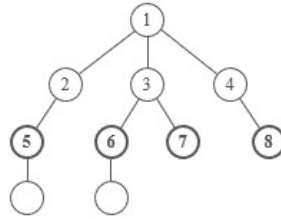


start at the root node (top of the tree)

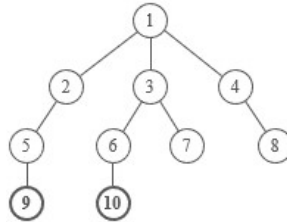
We would visit all the immediate children (all the nodes that are one step away from our starting node):



Then we would move on to all those nodes' children (all the nodes that are two steps away from our starting node):



And so on until we reach the end:



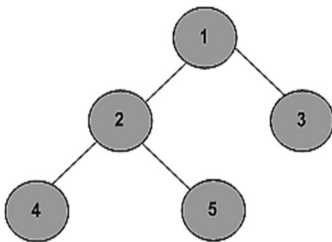
Advantages:

- BFS will find the shortest path between the starting point and any other reachable node.

Disadvantages:

- BFS on a binary tree generally requires more memory than a DFS.

There is only one kind of breadth-first traversal known as the level order traversal. This traversal visits nodes by levels from top to bottom and from left to right:

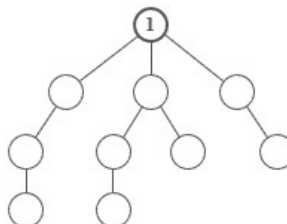


Level order traversal is:

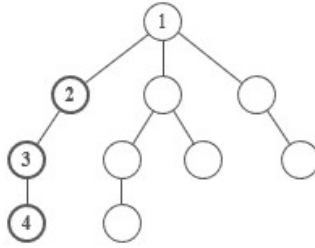
1 2 3 4 5

1.4.2.2 Depth-First Search

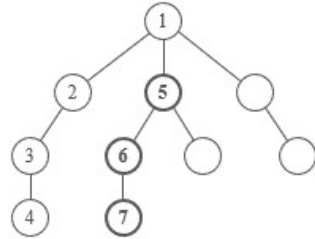
Depth-first search (DFS) is another method for exploring a tree or graph. In a DFS, you go as deep as possible down one path before backing up and trying a different one. Depth-first search is like walking through a corn maze. You explore one path, hit a dead end, and go back and try a different one. Here's a how a DFS would traverse this tree, starting with the root:



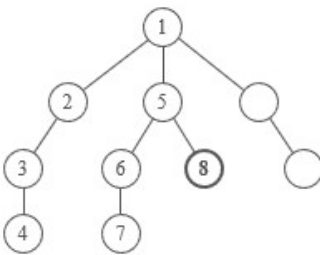
We would go down the first path we find until we hit a dead end:



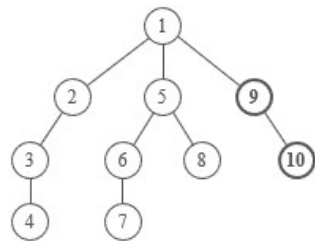
Then we would do the same thing again—go down a path until we hit a dead end:



And again:



And again:



Until we reach the end.

Advantages:

- DFS on a binary tree generally requires less memory than breadth-first.
- DFS can be easily implemented with recursion.

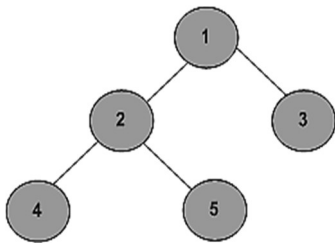
Disadvantages:

- DFS does not necessarily find the shortest path to a node

There are three commonly used depth-first traversals to visit all the nodes in a tree. The difference between these patterns is the order in which each node is visited. These three traversals are:

- Inorder Traversal (Left- Root -Right): we recursively do a Left- Root -Right traversal on the left subtree, visit the root node, and finally do a recursive Left-Root-Right traversal of the right subtree.

- Preorder Traversal (Root-Left-Right): we visit the root node first, then recursively do a DFS traversal of the left subtree, followed by a recursive preorder traversal of the right subtree.
- Postorder Traversal (Left-Right- Root): we recursively do a Left-Right-Root traversal of the left subtree and the right subtree followed by a visit to the root node.



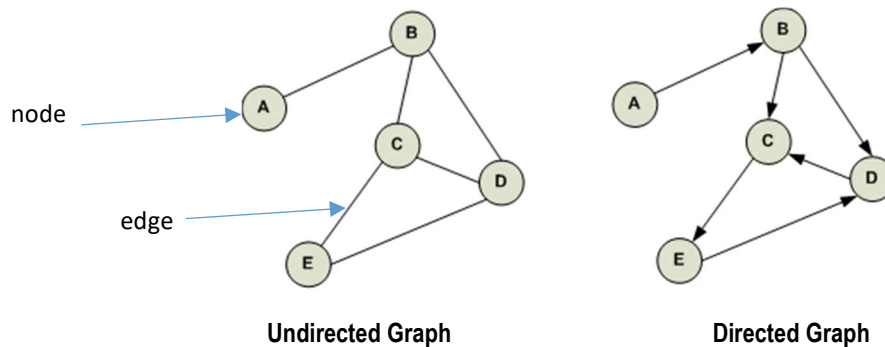
Inorder traversal: 4 2 5 1 3

Preorder traversal: 1 2 4 5 3

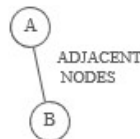
Postorder traversal: 4 5 2 3 1

1.5 Graph

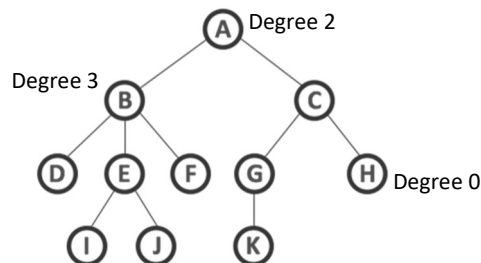
A graph is another abstract data structure with nodes (or vertices) that are connected by edges. Graphs can be directed or undirected. In directed graphs, edges point from the node at one end to the node at the other end. In undirected graphs, the edges simply connect the nodes at each end. The edges could represent distance or weight.



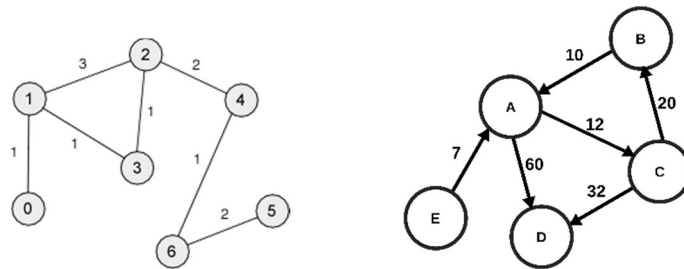
Graph are useful in cases where you have things that connect to other things. Nodes and edges could, for example, respectively represent cities and highways, routers and Ethernet cables, or Facebook users and their friendships. Two nodes connected by an edge are *adjacent* or neighbors.



The *degree* of a node is the number of edges connected to the node.



If a graph is weighted, each edge has a weight. The weight could, for example, represent the distance between two locations, or the cost or time it takes to travel between the locations.



A graph is cyclic if it has at least one cycle (an unbroken series of nodes with no repeating nodes or edges that connects back to itself). Graphs without cycles are acyclic.

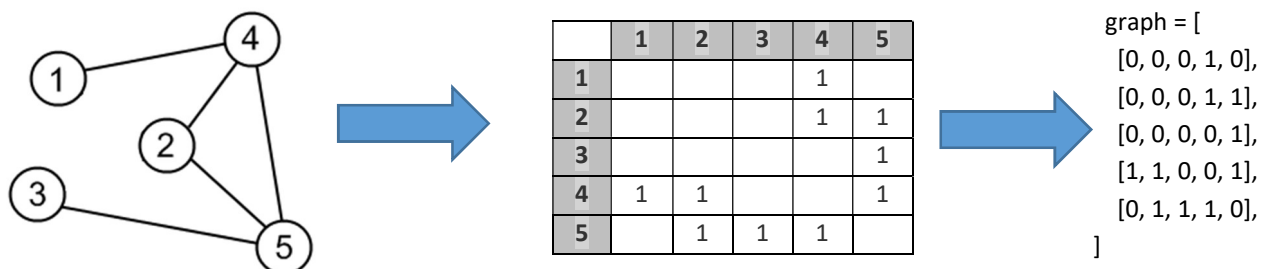


1.5.1 Graph Representation

There are two well-known implementations of a graph, the adjacency matrix and the adjacency list.

1.5.1.1 Adjacency Matrix

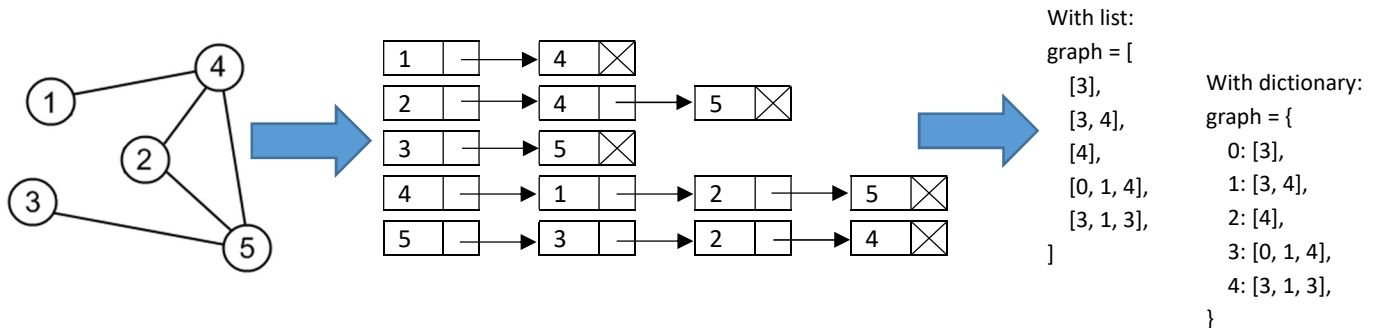
One of the easiest ways to implement a graph is to use a two-dimensional matrix of size $V \times V$ where V is the number of vertices in a graph. In this matrix implementation, each of the rows and columns represent a vertex in the graph. The value that is stored in the cell at the intersection of row v and column w indicates if there is an edge from vertex v to vertex w . A value in a cell represents the weight of the edge from vertex v to vertex w (for unweighted graphs we use 1).



The advantage of the adjacency matrix is that it is simple, and for small graphs it is easy to see which nodes are connected to other nodes. However, notice that most of the cells in the matrix are empty. Because most of the cells are empty we say that this matrix is “sparse”. A matrix is not a very efficient way to store sparse data. The adjacency matrix is a good implementation for a graph when the number of edges is large (in the order of $|V|^2$). A matrix is full when every vertex is connected to every other vertex. There are few real problems that approach this sort of connectivity.

1.5.1.2 Adjacency List

A more space-efficient way to implement a sparsely connected graph is to use an adjacency list. In an adjacency list implementation we keep a master list of all the vertices in the Graph object and then each vertex object in the graph maintains a list of the other vertices that it is connected to.



The advantage of the adjacency list implementation is that it allows us to compactly represent a sparse graph. The adjacency list also allows us to easily find all the links that are directly connected to a particular vertex.

Python does not have a graph data type. To use graphs you can either use the *networkx* module (pip install networkx) or implement it yourself

```
import networkx as nx
```

```
G=nx.Graph()
G.add_node("A")
G.add_node("B")
G.add_node("C")
G.add_edge("A","B")
G.add_edge("B","C")
G.add_edge("C","A")
```

```
print("Nodes: " + str(G.nodes()))
print("Edges: " + str(G.edges()))
```

Nodes: ['A', 'C', 'B']

Edges: [('A', 'C'), ('A', 'B'), ('C', 'B')]

For more details on networkx, read https://networkx.github.io/documentation/stable/downloads/networkx_reference.pdf

```
class Graph(object):
    def __init__(self, graph_dict=None):
        if graph_dict == None:
            graph_dict = {}
        self.__graph_dict = graph_dict

    def add_vertex(self, vertex):
        if vertex not in self.__graph_dict:
            self.__graph_dict[vertex] = []

    def add_edge(self, edge):
        edge = set(edge)
        (vertex1, vertex2) = tuple(edge)
        if vertex1 in self.__graph_dict:
            self.__graph_dict[vertex1].append(vertex2)
```

initializes the graph object on adjacency list mode using a dictionary. If no dictionary or None is given, an empty dictionary will be used

edges can be provided using sets, lists or tuples

```

else:
    self.__graph_dict[vertex1] = [vertex2]

def vertices(self):
    return list(self.__graph_dict.keys())

def edges(self):
    return self.__generate_edges()

```

1.5.2 BFS and DFS for a graph

Given a graph G and a starting vertex s , a breadth first search proceeds by exploring edges in the graph to find all the vertices in G for which there is a path from s . The remarkable thing about a breadth first search is that it finds *all* the vertices that are a distance k from s before it finds *any* vertices that are a distance $k+1$.

BFS for a graph is similar to Breadth First Traversal of a tree. The only catch here is, unlike trees, graphs may contain cycles, so we may come to the same node again. To avoid processing a node more than once, we use a boolean visited array.

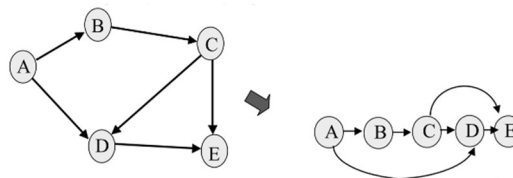
DFS explores the graph depth wise, that is we move down the graph from one node to another till the time we are out of all unexplored edges at a node, once we reach that condition, we backtrack and move to node one before and so on. In this traversal, we start with a node s , and mark it as visited. Now s be our current node u :

- Move to v where there is an edge (u,v) .
- If v is already visited, we move back to u .
- If v is not visited, mark v as visited and make v as current node and repeated above steps.
- At some point all the edges at u will lead to already visited node, then we drop u and move to node which was visited before u .

This backtracking will make us reach the start node s again, and when there is no edge left to be explored at s , the graph traversal is done. DFS for a graph is similar to preorder traversal of a tree. The only catch here is, unlike trees, graphs may contain cycles, so we may come to the same node again. To avoid processing a node more than once, we also use a boolean visited array.

1.5.3 Topological Sorting

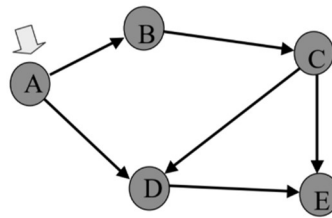
A topological sort takes a directed acyclic graph (DAG) and produces a linear ordering of all its vertices such that if the graph G contains an edge (v,w) then the vertex v comes before the vertex w in the ordering. DAGs are used in many applications to indicate the precedence of events. For example, software project schedules, precedence charts for optimizing database queries, and multiplying matrices. Any linear ordering in which all the arrows go to the right is a valid solution. Topological sorting for a graph is not possible if the graph is not a DAG.



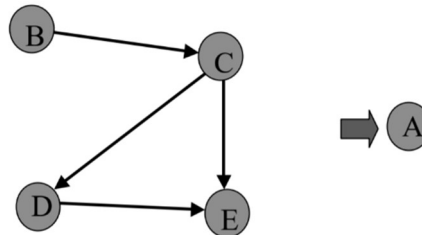
Topological sorting can be accomplished in two different ways:

- Identifying incoming edges:

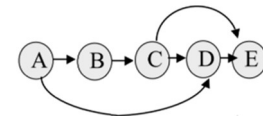
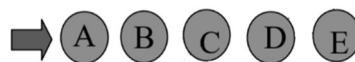
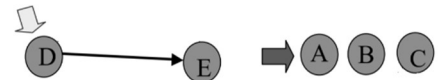
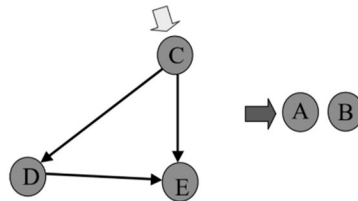
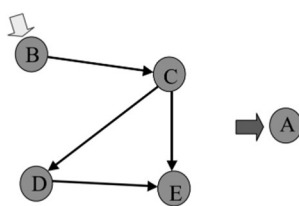
Step 1- Identify vertices that have no incoming edge (in-degree is zero). If no such edges exist, graph is not DAG. If there are several vertices of in-degree zero, select one.



Step 2- Delete this vertex of in-degree zero and all its outgoing edges from the graph. Place it in the output.



Step 3- Repeat Steps 1 and Step 2 until graph is empty

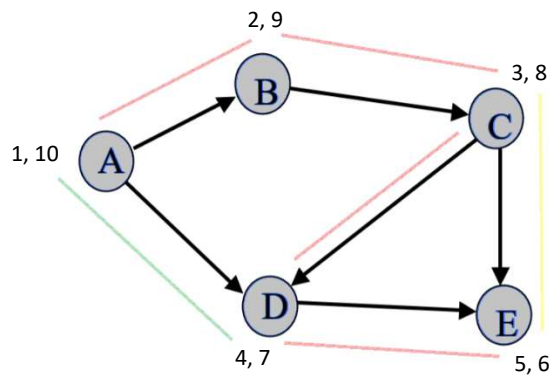


- Using DFS:

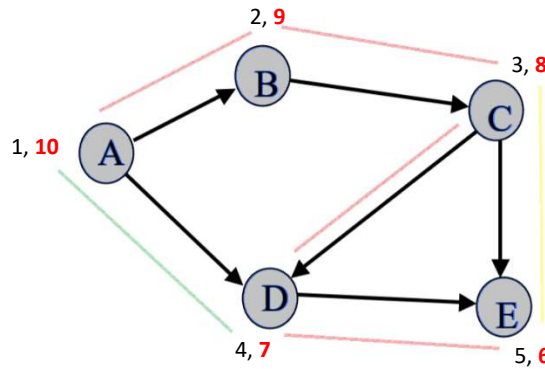
The topological sort is a simple but useful adaptation of a depth first search:

Step 1- Perform DFS for some graph G . (we call DFS to compute the finish times for each of the vertices)

DFS starting from node A



Step 2- Store the vertices in a list in decreasing order of finish time



Step 3- Return the ordered list as the result of the topological sort



1.6 Hash Tables

A hash table is a collection of items which are stored in such a way as to make it easy to find them later. Each position of the hash table, often called a *slot*, can hold an item and is named by an integer value starting at 0. For example, we will have a slot named 0, a slot named 1, a slot named 2, and so on. Initially, the hash table contains no items so every slot is empty. We can implement a hash table by using a list with each element initialized to the special Python value `None`.

hash table of
size $m=11$

0	1	2	3	4	5	6	7	8	9	10
None	None	None	None	None	None	None	None	None	None	None

The mapping between an item and the slot where that item belongs in the hash table is called the hash function. In simple terms, a hash function maps a big number or string to a small integer that can be used as index in hash table. A good hash function must be efficiently computable and should uniformly distribute the keys (each table position equally likely for each key). Assume that we have the set of integer items 32, 53, 79, 15, 93, and 67. The modular method is the most common hash function, it simply takes an item and divides it by the table size, returning the remainder as its hash value ($\text{hash}(\text{item}) = \text{item} \% 11$).

Item	Hash Value ($\text{item} \% 11$)
32	10
53	9
79	2
15	4
93	5
67	1

Once the hash values have been computed, we can insert each item into the hash table at the designated position. Note that 6 of the 11 slots are now occupied. This is referred to as the *load factor*, and is commonly denoted by $\lambda = \frac{\text{\# of items}}{\text{table size}}$ (For this example, $\lambda = \frac{6}{11}$).

0	1	2	3	4	5	6	7	8	9	10
None	67	79	None	15	93	None	None	None	53	32

Now when we want to search for an item, we simply use the hash function to compute the slot name for the item and then check the hash table to see if it is present. This searching operation takes constant time, since a constant amount of time is required to compute the hash value and then index the hash table at that location. However, this technique is going to work only if each item maps to a unique location in the hash table. For example, if we want to add the item 42 to our collection, it would have a hash value of 9 ($42\%11=9$). Since 53 also had a hash value of 9, we would have a problem. According to the hash function, two or more items would need to be in the same slot. This is referred to as a collision, and it creates problems for the hashing technique.

1.6.1 Hash Functions

Given a collection of items, a hash function that maps each item into a unique slot is referred to as a *perfect hash function*. If we know the items and the collection will never change, then it is possible to construct a perfect hash function. Unfortunately, given an arbitrary collection of items, there is no systematic way to construct a perfect hash function. Thus, the goal is to create a hash function that minimizes the number of collisions, is easy to compute, and evenly distributes the items in the hash table. There are several ways to extend the simple modular method.

1.6.1.1 Folding method

This method for constructing hash functions begins by dividing the item into equal-size pieces (the last piece may not be of equal size). These pieces are then added together to give the resulting hash item. For example, if our item was the phone number 436-555-4601, we would take the digits and divide them into groups of 2 (43,65,55,46,01). After the addition, $43+65+55+46+01$, we get 210. If we assume our hash table has 11 slots, then we need to perform the extra step of dividing by 11 and keeping the remainder. In this case $210\%11$ is 1, so the phone number 436-555-4601 hashes to slot 1. Some folding methods go one step further and reverse every other piece before the addition. For the above example, we get $43+56+55+64+01=219$ which gives $219\%11=10$.

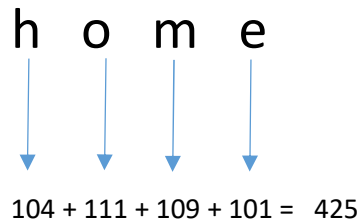
1.6.1.2 Mid-square method

This method for constructing hash functions begins by squaring the item, and then extract some portion of the resulting digits. For example, for the item 44, we would first compute $44^2=1,936$. By extracting the middle two digits, 93, and performing the remainder step, we get $93\%11=5$.

Item	Hash Value (item%11)	Mid-square
32	10	2
53	9	3
79	2	2
15	4	3
93	5	9
67	1	4

We can also create hash functions for character-based items such as strings. The word “home” can be thought of as a sequence of ordinal values. We can then take these four ordinal values, add them up, and use the modular method to get a hash value:

```
>>> ord('h')
104
>>> ord('o')
111
>>> ord('m')
109
>>> ord('e')
101
```



$$425 \% 11 = 7$$

`ord(c)`: Given a string representing one Unicode character, return an integer representing the Unicode code point of that character

You may be able to think of a number of additional ways to compute hash values for items in a collection. The important thing to remember is that the hash function has to be efficient so that it does not become the dominant part of the storage and search process.

1.6.2 Collision Resolution

When two items hash to the same slot, we must have a systematic method for placing the second item in the hash table. This process is called collision resolution. As we stated earlier, if the hash function is perfect, collisions will never occur. However, since this is often not possible, collision resolution becomes a very important part of hashing.

1.6.2.1 Open Addressing

One method for resolving collisions looks into the hash table and tries to find another open slot to hold the item that caused the collision. A simple way to do this is to start at the original hash value position and then move in a sequential manner through the slots until we encounter the first slot that is empty (we may need to go back to the first slot (circularly) to cover the entire hash table). This collision resolution process is referred to as *open addressing*, where we try to find the next open slot or address in the hash table. By systematically visiting each slot one at a time, we are performing an open addressing technique called *linear probing*.

When we attempt to place 60 into slot 5, a collision occurs. Under linear probing, we look sequentially, slot by slot, until we find an open position. In this case, we find slot 6. Again, 71 should go in slot 5 but must be placed in slot 7 since it is the next open position. The final value of 20 hashes to slot 9. Since slot 9 is full, we begin to do linear probing. We visit slot 10 and finally find an empty slot at position 0.


0	1	2	3	4	5	6	7	8	9	10
20	67	79	None	15	93	60	71	None	53	32

Once we have built a hash table using open addressing and linear probing, it is essential that we utilize the same methods to search for items. Assume we want to look up the item 93. When we compute the hash value, we get 5. Looking in slot 5 reveals 93, and we can return True. However, when we look for item 20,

the hash value is 9, and slot 9 is currently holding 53. We cannot simply return False since we know that there could have been collisions. We are now forced to do a sequential search, starting at position 10, looking until either we find the item 20 or we find an empty slot.

A disadvantage to linear probing is the tendency for *clustering*; items become clustered in the table. This means that if many collisions occur at the same hash value, a number of surrounding slots will be filled by the linear probing resolution. This will have an impact on other items that are being inserted, for example, when we try to add the item 26, a cluster of values hashing to 5 had to be skipped to finally find an open position.

0	1	2	3	4	5	6	7	8	9	10
20	67	79	None	15	93	60	71	None	53	32



cluster

One way to deal with clustering is to extend the linear probing technique so that instead of looking sequentially for the next open slot, we skip slots, thereby more evenly distributing the items that have caused collisions. This will potentially reduce the clustering that occurs. Below there is a hash table where the collision resolution is done with a +3 probe. This means that once a collision occurs, we will look at every third slot until we find one that is empty.

0	1	2	3	4	5	6	7	8	9	10
71	67	79	None	15	93	None	20	60	53	32

The general name for this process of looking for another slot after a collision is *rehashing*. In general, $hash_index = (hash_index + skip) \% hash_table_size$. It is important to note that the size of the “skip” must be such that all the slots in the table will eventually be visited. Otherwise, part of the table will be unused. To ensure this, it is often suggested that the table size be a prime number.

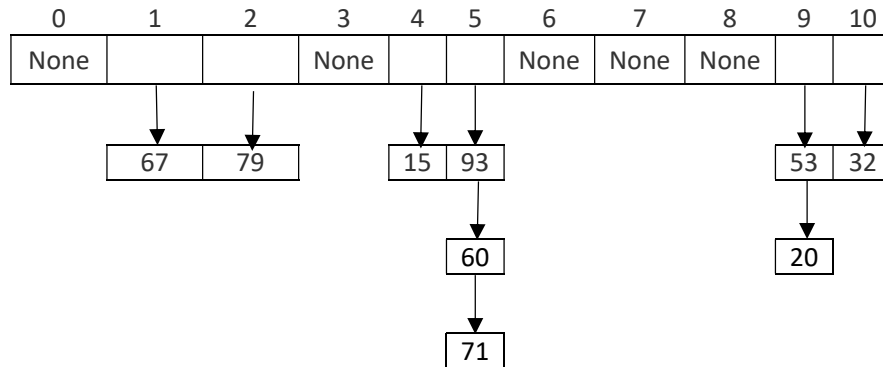
A variation of the *linear probing* idea is called *quadratic probing*. This method is similar to linear probing and the only difference is that if you were to try to insert into a space that is filled you would first check $1^2 = 1$ element away then $2^2 = 4$ elements away, then $3^2 = 9$ elements away then $4^2 = 16$ elements away and so on.

0	1	2	3	4	5	6	7	8	9	10
None	67	79	71	15	93	60	20	None	53	32

With linear probing we know that we will always find an open spot if one exists (it might be a long search but we will find it). However, this is not the case with quadratic probing unless you take care in the choosing of the table size. In order to guarantee that the quadratic probes will hit every single available spot eventually, the hash table size must be a prime number and never be more than half full (even by one element).

1.6.2.2 Chaining

An alternative method for handling the collision problem is to allow each slot to hold a reference to a collection (or chain) of items. Chaining allows many items to exist at the same location in the hash table. When collisions happen, the item is still placed in the proper slot of the hash table. As more and more items hash to the same location, the difficulty of searching for the item in the collection increases.



When we want to search for an item, we use the hash function to generate the slot where it should reside. Since each slot holds a collection, we use a searching technique to decide whether the item is present. The advantage is that on the average there are likely to be many fewer items in each slot, so the search is perhaps more efficient.

It is important to mention that if the load factor (λ) is small, then there is a lower chance of collisions, meaning that items are more likely to be in the slots where they belong. If λ is large, meaning that the table is filling up, then there are more and more collisions. This means that collision resolution is more difficult, requiring more comparisons to find an empty slot.

1.6.3 Hashing Implementation

Python has useful data types that could help us to develop our own hashing algorithm. For example, a dictionary that can store key-data pairs could be used to look up a data value based on a given key. However, we are living in a digitally disrupted world, where sensitive information flows among users and organizations. As a result, information security has become very important in most organizations. The main reason for this is that access to information and the associated resources has become easier because of the developments in distributed processing, for example the Internet and electronic commerce. The result is that organizations need to ensure that their information is properly protected and that they maintain a high level of information security.

Hash functions are used inside some cryptographic algorithms, in digital signatures, message authentication codes, manipulation detection, fingerprints, checksums (message integrity check), hash tables, password storage and much more. As a Python programmer you may need these functions to check for duplicate data or files, to check data integrity when you transmit information over a network, to securely store passwords in databases, or maybe some work related to cryptography. The Federal Information Processing Standards, as well as the Internet Engineering Task Force have defined and developed different hash algorithms that are widely use all over the Internet. Python's hashlib module implements a common

interface to many different secure hash and message digest algorithms. Included are the FIPS secure hash algorithms SHA1, SHA224, SHA256, SHA384, and SHA512 (defined in FIPS 180-2) as well as RSA's MD5 algorithm (defined in Internet RFC 1321).

```
class HashTable:
    def __init__(self):
        self.size = 11
        self.slots = [None] * self.size
        self.data = [None] * self.size

    def hashfunction(self, key, size):
        return key % size

    def rehash(self, oldhash, size):
        return (oldhash + 1) % size
```



Class structure to develop
our own hashing
techniques

```
list(map(hash, [0, 1, 2, 3]))
# Output: [0, 1, 2, 3]
```



hashing using Python's built-in hash function. Python is
using different hash() function depending on the type
of data

```
list(map(hash, ['0', '1', '2', '3']))
# Output: [3512700405625046622, -2154559771013955726, 4558503884423229281, -5090714507869946363]
```

```
hash('0')
# Output: 3512700405625046622
```

```
import hashlib
```



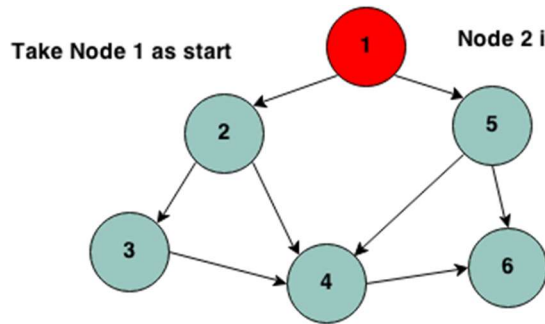
hashing using Python's hashlib module

```
print(hashlib.algorithms_available)
# Output: {'sha', 'sha224', 'sha3_256', 'DSA', 'shake_128', 'sha512', 'SHA224', 'RIPEMD160',
'dsaEncryption', 'ripemd160', 'sha3_512', 'MD5', 'SHA384', 'SHA256', 'md5', 'blake2b',
'sha3_384', 'sha384', 'sha256', 'sha3_224', 'ecdsa-with-SHA1', 'whirlpool', 'SHA', 'shake_256',
'sha1', 'md4', 'SHA1', 'DSA-SHA', 'SHA512', 'blake2s', 'MD4', 'dsaWithSHA'}
```

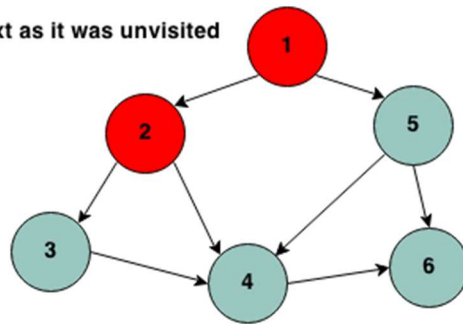
```
hash_object = hashlib.md5(b'Hello World')
print(hash_object.hexdigest())
# Output: b10a8db164e0754105b7a99be72e3fe5
```

```
value = input('Enter string to hash: ')
hash_object = hashlib.sha256(value.encode())
print(hash_object.hexdigest())
# Output:
Enter String to hash: Hello World
a591a6d40bf420404a011733cfb7b190d62c65bf0bcda32b57b277d9ad9f146e
```

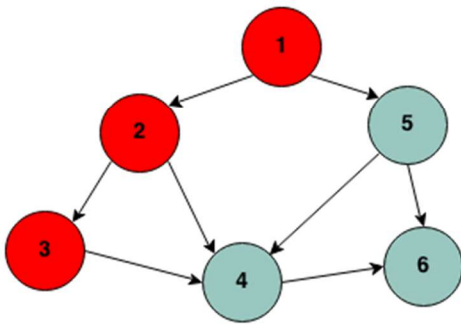
Appendix 1. Depth First Search Example



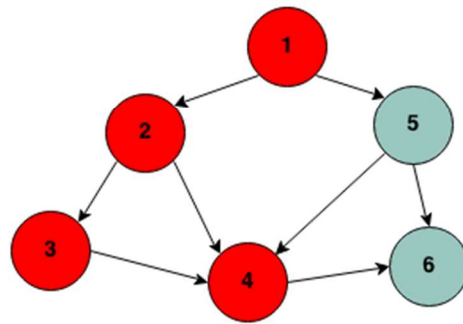
Node 2 is visited next as it was unvisited



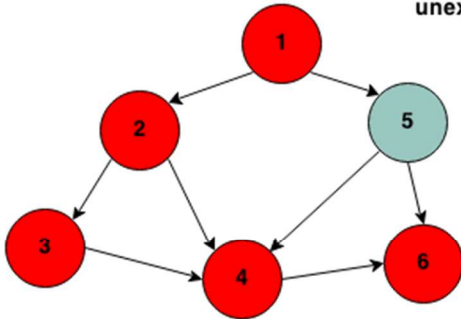
Node 3 is visited next as it was unvisited



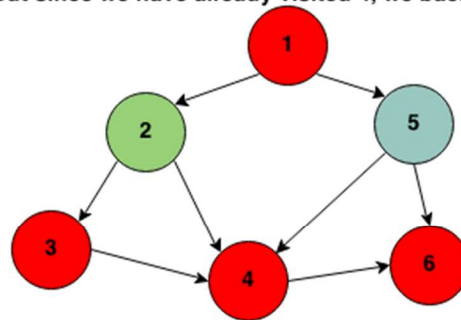
Node 4 is visited next as it was unvisited



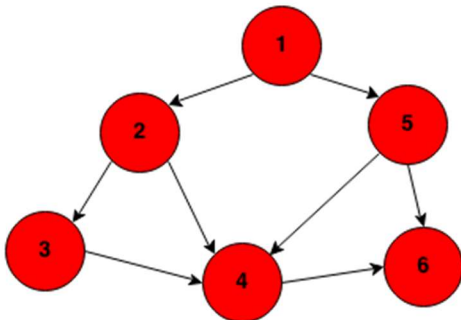
Node 6 is visited next as it was unvisited



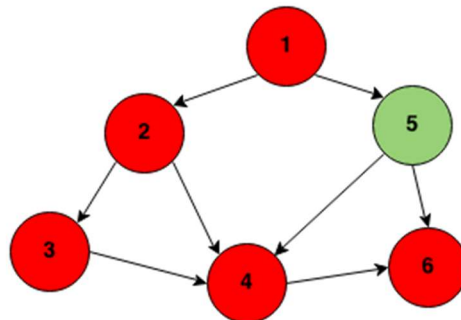
As there is no more nodes which can be traversed from 6, we back track, first node where there is unexplored edge is 2. But since we have already visited 4, we backtrack further



Node 5 from node 1 is visited next as it was unvisited

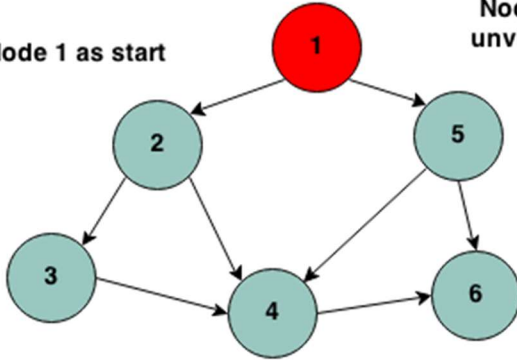


Again all neighbor nodes from 5 are already visited, hence we back track to node 1

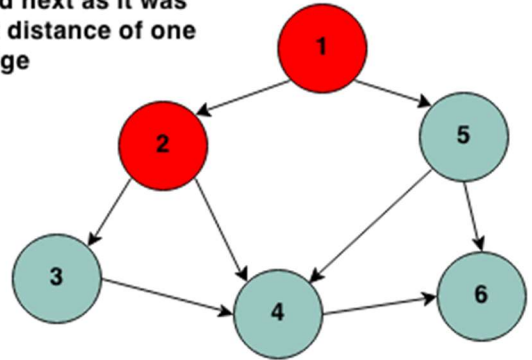


Appendix 2. Breadth First Search Example

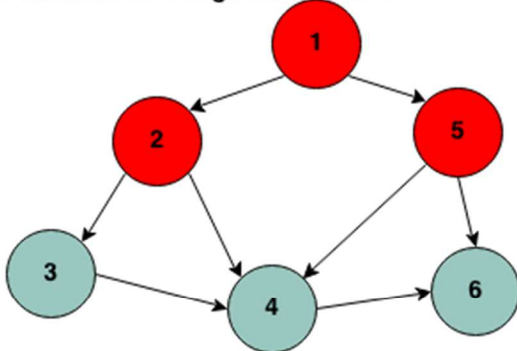
Take Node 1 as start



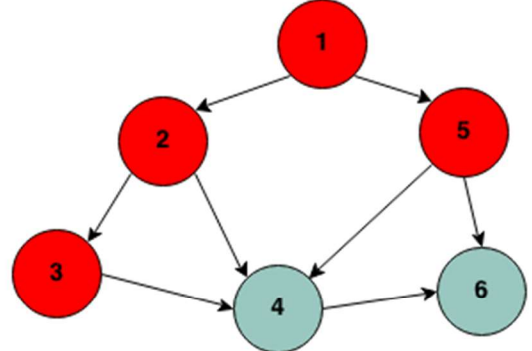
Node 2 is visited next as it was unvisited and at distance of one edge



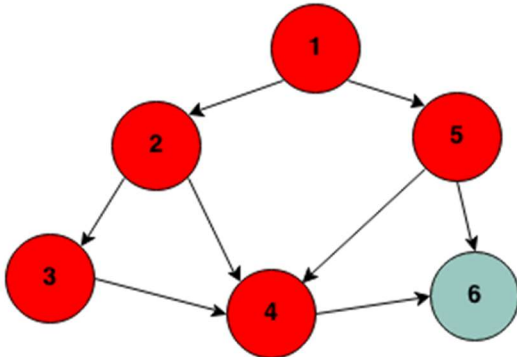
Node 5 is visited next as it was unvisited and at distance of 1 edge from node 1



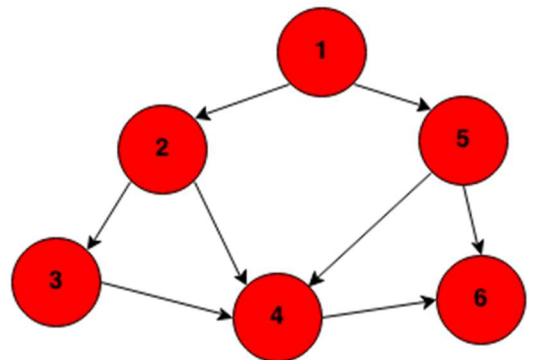
Node 3 is visited next as it was unvisited and at distance of one edge from node 2



Node 4 is visited next as it was unvisited and at distance of one edge from node 2



Node 6 is visited next as it was unvisited and at distance of one edge from node 5



References

Data Structures and Algorithms in Python. Michael H. Goldwasser, Michael T. Goodrich, and Roberto Tamassia. Chapters 6, 7, 8, 10, 14

https://www.tutorialspoint.com/data_structures_algorithms/stack_algorithm.htm

http://en.wikipedia.org/wiki/Priority_queue

https://www.tutorialspoint.com/data_structures_algorithms/tree_data_structure.htm

<https://docs.python.org/3/library/hashlib.html>