**Assignment 6: Medians and Order Statistics & Elementary Data Structures**

Syed Noor UI Hassan

University of the Cumberlands

Algorithms and Data Structures (MSCS-532-B01)

November 24, 2024

**Part 1: Implementation and Analysis of Selection Algorithms**

**Implementation**

       Implementing selection algorithms for finding the kth smallest element in an array involves two main approaches: a deterministic algorithm (Median of Medians) and a randomized algorithm (Randomized Quickselect). The deterministic algorithm achieves worst-case linear time complexity by carefully selecting a pivot element through a recursive process of finding the median of medians. This method guarantees that the array is divided into well-balanced partitions, ensuring $O(n)$ time complexity even in the worst case. On the other hand, the randomized algorithm relies on randomly choosing a pivot, which leads to an expected linear time complexity but may degrade to quadratic time in rare worst-case scenarios. Both algorithms use a divide-and-conquer strategy, recursively partitioning the array and narrowing down the search space until the kth element is found.

**Performance Analysis**

       Performance analysis reveals that the deterministic algorithm's $O(n)$ worst-case time complexity stems from its clever pivot selection method, which ensures that at least 30% of the elements are eliminated in each recursive step. This leads to a recurrence relation that solves linear time. The randomized algorithm achieves $O(n)$ expected time complexity due to the average-case scenario where the randomly chosen pivot divides the array roughly in half, resulting in a balanced partitioning. However, its worst-case performance can be $O(n^2)$ if consistently poor pivots are chosen, though this is highly unlikely. Space complexity for both algorithms is generally $O(\log n)$ due

to the recursion stack, with the randomized algorithm potentially using O(n) space in

the worst case.

**Empirical Analysis**

```python
import time
import random

def select(arr, k):
    if len(arr) <= 5:
        return sorted(arr)[k]

    medians = [sorted(arr[i:i+5])[len(arr[i:i+5])//2] for i in range(0, len(arr), 5)]
    pivot = select(medians, len(medians)//2)

    left = [x for x in arr if x < pivot]
    equal = [x for x in arr if x == pivot]
    right = [x for x in arr if x > pivot]

    if k < len(left):
        return select(left, k)
    elif k < len(left) + len(equal):
        return pivot
    else:
        return select(right, k - len(left) - len(equal))

def randomized_select(arr, k):
    if len(arr) == 1:
        return arr[0]

    pivot = random.choice(arr)
    left = [x for x in arr if x < pivot]
    equal = [x for x in arr if x == pivot]
    right = [x for x in arr if x > pivot]

    if k < len(left):
        return randomized_select(left, k)
    elif k < len(left) + len(equal):
        return pivot
    else:
        return randomized_select(right, k - len(left) - len(equal))

def measure_time(func, arr, k):
    start = time.time()
    result = func(arr.copy(), k)
    end = time.time()
    return end - start

sizes = [1000, 10000, 100000]
distributions = ['random', 'sorted', 'reverse_sorted']

with open('output.txt', 'w') as f:
    for size in sizes:
        for dist in distributions:
            if dist == 'random':
                arr = [random.randint(1, 1000000) for _ in range(size)]
            elif dist == 'sorted':
                arr = list(range(1, size+1))
            else:  # reverse_sorted
                arr = list(range(size, 0, -1))

            k = size // 2  # Find median

            det_time = measure_time(select, arr, k)
            rand_time = measure_time(randomized_select, arr, k)

            output = f"Size: {size}, Distribution: {dist}\n"
            output += f"Deterministic: {det_time:.6f} seconds\n"
            output += f"Randomized: {rand_time:.6f} seconds\n\n"

            f.write(output)
            print(output)

print("Results have been exported to output.txt")
```

```
Size: 1000, Distribution: random
Deterministic: 0.000479 seconds
Randomized: 0.000262 seconds

Size: 1000, Distribution: sorted
Deterministic: 0.000372 seconds
Randomized: 0.000242 seconds

Size: 1000, Distribution: reverse_sorted
Deterministic: 0.000383 seconds
Randomized: 0.000105 seconds

Size: 10000, Distribution: random
Deterministic: 0.004529 seconds
Randomized: 0.001524 seconds

Size: 10000, Distribution: sorted
Deterministic: 0.003164 seconds
Randomized: 0.001315 seconds

Size: 10000, Distribution: reverse_sorted
Deterministic: 0.003261 seconds
Randomized: 0.001210 seconds

Size: 100000, Distribution: random
Deterministic: 0.034789 seconds
Randomized: 0.010520 seconds

Size: 100000, Distribution: sorted
Deterministic: 0.026110 seconds
Randomized: 0.011433 seconds

Size: 100000, Distribution: reverse_sorted
Deterministic: 0.025953 seconds
Randomized: 0.007317 seconds
```

Empirical analysis comparing these algorithms on various input sizes and distributions (random, sorted, and reverse-sorted) reveals interesting insights. For smaller input sizes, both algorithms perform comparably. As the input size increases, the randomized algorithm tends to outperform the deterministic one, especially on random distributions. However, the deterministic algorithm shows more consistent performance across different input distributions. The randomized algorithm's performance can vary, excelling on random data

but potentially slowing down on sorted or reverse-sorted arrays. These observations align with the theoretical analysis, highlighting the trade-offs between guaranteed worst-case performance and average-case efficiency. In practice, the choice between these algorithms depends on the specific requirements of the application, with the randomized algorithm often preferred for its generally faster average-case performance. In contrast, the deterministic algorithm is favored in scenarios demanding consistent, predictable performance regardless of input distribution.

## Part 2: Elementary Data Structures Implementation and Discussion

### Implementation

Implementing elementary data structures in Python provides a foundation for understanding fundamental concepts in computer science and software development. Arrays and matrices, represented using Python lists, offer efficient random access and are particularly useful for tasks requiring direct indexing. The Matrix class implementation demonstrates creating a 2D structure with methods for getting, setting, and manipulating rows. This structure is invaluable in various applications, from image processing to scientific computing, where organized data representation is crucial.

### Performance Analysis

While conceptually different, stacks and queues can be implemented using Python lists. The Stack class, with its push, pop, and peek operations, exemplifies the Last-In-FirstOut (LIFO) principle. This data structure finds extensive use in scenarios such as managing function calls in programming languages, implementing undo mechanisms in applications, and parsing expressions in compilers. On the other hand, the Queue

class, implementing enqueue and dequeue operations, represents the First-In-First-Out (FIFO) principle. Queues are essential in task scheduling, breadth-first search algorithms, and managing data buffers in various systems.

Linked lists, implemented through a chain of Node objects, offer a dynamic approach to data storage. The LinkedList class demonstrates operations like insertion at the beginning and end, deletion, and traversal. This structure excels in scenarios requiring frequent insertions and deletions, especially at the beginning of the list. Linked lists are fundamental in implementing more complex data structures like hash tables. They are often used in scenarios where the data size is unknown or frequently changing, such as inmemory allocation systems or certain types of caches.

Performance analysis of these data structures reveals their strengths and limitations. Arrays and matrices provide $O(1)$ access time but suffer from $O(n)$ complexity for insertions and deletions in arbitrary positions. Stacks implemented with arrays offer $O(1)$ complexity for all operations, making them highly efficient. When implemented with simple arrays, Queues face $O(n)$ complexity for dequeue operations, which can be optimized using circular arrays or linked lists. Linked lists provide $O(1)$ insertion at the beginning but require $O(n)$ time for most other operations due to the need for traversal.

The trade-offs between using arrays and linked lists for implementing stacks and queues are significant. Arrays benefit from better cache locality and constant-time access, making them generally preferred for stack implementations. However, they may face issues with fixed size in some languages and inefficient insertions or deletions in the middle. Linked lists, while requiring extra memory for pointers and

suffering from poor cache locality, offer dynamic sizing and efficient insertions and deletions at both ends, making them often more suitable for queue implementations.

**Discussion**

In practical applications, the choice of data structure depends heavily on the specific requirements of the problem at hand. Arrays and matrices are indispensable in fields like image processing, scientific computing, and game development, where efficient numerical computations and representation of multi-dimensional data are crucial. Stacks find their use in managing function calls, implementing undo mechanisms and parsing expressions. Queues are vital in operating systems for task scheduling, in graph algorithms for breadthfirst search, and in managing data streams. Linked lists, with their dynamic nature, are excellent for implementing hash tables, managing music playlists, and representing data with frequently changing size or structure.

In conclusion, understanding these elementary data structures and their implementations is crucial for any software developer or computer scientist. Each structure has its unique properties, making it more or less suitable for different scenarios. The ability to choose the right data structure for a given problem can significantly impact the efficiency and effectiveness of software solutions. As technology evolves, these fundamental concepts continue to form the backbone of more complex systems and algorithms, underlining their enduring importance in the field of computer science.