

Medians and order Statistics & Elementary Data Structures

Venkata Manasa Kakarla

Student ID: 005027517

Algorithms and Data Structures (MSCS532_A06)

Assignment – A06

Part 1: Implementation and Analysis of Selection Algorithms

Performance Analysis:

Time Complexity Analysis of Selection Algorithms:

1. Deterministic Selection: Median of Medians:

The Median of Medians algorithm guarantees $O(n)$ time complexity in the worst case by systematically reducing the size of the input through a well-defined pivot selection process. Here's a breakdown of its time complexity:

Steps Breakdown:

- **Divide:** The array is divided into groups of 5 elements. The number of groups is $n/5$.
- **Find Medians:** For each group, we sort it and find the median. This sorting operation takes constant time since the groups are small (at most 5 elements). Thus, finding the medians takes $O(n/5) \cdot O(5 \log 5) = O(n)$.
- **Recursive Call:** We then recursively find the median of these medians, which is at most $O(n/5)$.
- **Partitioning:** The array is partitioned around the median of medians. This takes $O(n)$ time.

Recurrence Relation: The overall time can be expressed using the recurrence relation:

$$T(n) = T(n/5) + T(7n/10) + O(n)$$

Where $T(n/5)$ accounts for finding the median of medians, $T(7n/10)$ accounts for the worst-case partitioning, and $O(n)$ is the time taken for partitioning.

Solving the Recurrence: Using the Master Theorem:

- The first term dominates (the $O(n)$ component).
- This results in a time complexity of: $T(n) = O(n)$

2. Randomized Selection: Randomized Quickselect:

Randomized Quickselect operates with expected linear time complexity but has a worst-case time complexity of $O(n^2)$ in rare situations. Here's how it breaks down:

Steps Breakdown:

- **Randomized Partitioning:** A pivot is chosen randomly, and the array is partitioned around this pivot. The partitioning process takes $O(n)$.
- **Recursive Selection:** After partitioning, Quickselect either calls itself on the left or right subarray, depending on the position of the pivot.

Recurrence Relation: The average-case performance can be analyzed using the following recurrence: $T(n) = T(n/2) + O(n)$

In the average case, since the pivot is randomly selected, it typically splits the array into two equal halves, hence $n/2$.

Solving the Recurrence: Using the Master Theorem: $T(n) = T(n/2) + O(n)$

This is a classic divide-and-conquer recurrence that resolves to: $T(n) = O(n)$

Worst-case Analysis: In the worst case, the pivot can be the smallest or largest element repeatedly, leading to unbalanced partitions. The recurrence then becomes: $T(n) = T(n-1) + O(n)$

This resolves to: $T(n) = O(n^2)$

However, the worst-case scenario is rare with randomization, especially when implemented in practice.

Understanding Time Complexities of Deterministic and Randomized Selection Algorithms:

To clarify why the deterministic algorithm (Median of Medians) achieves $O(n)$ time complexity in the worst case, while the randomized algorithm (Randomized Quickselect) achieves $O(n)$ time complexity in expectation, we can break it down as follows:

1. Deterministic Algorithm: Median of Medians:

Worst-case Analysis:

- **Pivot Selection:** The Median of Medians algorithm guarantees that the pivot chosen (the median of the medians) is always a good pivot that splits the array reasonably well. Specifically, this method ensures that at least 30% of the elements are guaranteed to be on one side of the pivot.
- **Recurrence Relation:** The recurrence relation for the algorithm is:

$$T(n) = T(n/5) + T(7n/10) + O(n)$$

- Here, $T(n/5)$ accounts for the cost of finding the median of the medians, and $T(7n/10)$ accounts for the size of the largest partition after the pivot is selected.

The $O(n)$ term accounts for the time spent in partitioning.

- **Balanced Reduction:** Because the algorithm guarantees that at least 30% of the elements are on one side of the pivot, this leads to a balanced reduction of the problem size in each recursive call. Each level of recursion reduces the problem size effectively, preventing any part of the input from being excessively large in subsequent recursive calls.
- **Final Complexity:** When you apply the Master Theorem to this recurrence relation, you find that it resolves to $O(n)$ in the worst case, guaranteeing linear performance regardless of input distribution.

2. Randomized Algorithm: Randomized Quickselect:

Expected-case Analysis:

- **Random Pivot Selection:** In Randomized Quickselect, the pivot is chosen randomly. This randomness means that while you might occasionally choose a poor pivot, on average, you will likely pick a pivot that is close to the median of the dataset.

- Average-case Recurrence: The average-case recurrence relation can be expressed as:

$$T(n) = T(n/2) + O(n)$$

- In this case, the pivot generally splits the array into two approximately equal halves, leading to a balanced problem size in most cases.
- Probability of Good Choices: Because the pivot is chosen randomly, the expected number of elements remaining in the larger partition after each split is around half the array size. Over many iterations, the likelihood of consistently poor choices (which could lead to worst-case scenarios) is low.
- Final Complexity: Using the Master Theorem, this average-case scenario leads to $O(n)$ expected time complexity. However, the worst-case scenario occurs when the pivot chosen is consistently the smallest or largest element, leading to unbalanced partitions.

Worst-case Analysis:

- In the worst case, the pivot could lead to: $T(n) = T(n-1) + O(n)$. This results in a time complexity of $O(n^2)$ when the same unbalanced partition occurs repeatedly. However, this situation is rare in practice due to the random nature of pivot selection.

Space Complexity and Overheads of Selection Algorithms

1. Deterministic Algorithm: Median of Medians

Space Complexity:

- The Median of Medians algorithm has a space complexity of $O(n)$ in the worst case. This is primarily due to the way the algorithm handles the partitioning and the storage of medians.
 - When dividing the input into groups of 5, you create a new list to store the medians, which can grow up to $n/5$ in size.

- The space complexity for storing the partitioned arrays and the recursive stack space must also be considered. However, in practice, this stack space can be minimized.

Additional Overheads:

- Storage of Medians: Each group of 5 needs to be stored and sorted to find the median. This introduces additional overhead in terms of both time (for sorting) and space.
- Copying Arrays: The algorithm may require temporary arrays for partitioning, which can increase space usage, especially for large inputs.
- Recursive Call Stack: The depth of recursive calls can also add to space usage, depending on how balanced the partitions are.

2. Randomized Algorithm: Randomized Quickselect

Space Complexity:

- The Randomized Quickselect algorithm typically has a space complexity of $O(1)$ for the iterative version and $O(\log n)$ for the recursive version due to the recursive call stack.
 - In the iterative version, it can use a constant amount of extra space since it modifies the array in place without needing to create additional data structures.
 - If using recursion, each recursive call adds to the call stack, leading to $O(\log n)$ space in balanced scenarios. However, in the worst case (where poor pivot choices lead to skewed partitions), it can degrade to $O(n)$ for the stack space.

Additional Overheads:

- Random Number Generation: There is a slight overhead associated with generating random numbers to select the pivot. However, this overhead is generally minimal compared to the overall complexity of the algorithm.

- **In-Place Partitioning:** The space used for partitioning is minimal since elements are swapped within the array itself rather than creating new subarrays. This keeps the space footprint low.
- **Memory Management:** As it modifies the original array, no additional memory management overhead (like garbage collection for temporary arrays) is necessary.

Summary of Space Complexity and Overheads:

Algorithm	Space Complexity	Additional Overheads
Median of Medians	$O(n)$	Storage of medians, temporary arrays, recursive stack space
Randomized Quickselect	$O(1)$ (iterative) or $O(\log n)$ (recursive)	Random number generation, minimal in-place modifications.

Empirical Analysis:

```

empirical_comparison.py > ...
1 import time
2 import random
3
4 # Implementation of Median of Medians (Deterministic Selection)
5 def median_of_medians(arr, k):
6     if len(arr) <= 5:
7         return sorted(arr)[k]
8
9     # Divide arr into sublists of length 5
10    sublists = [arr[i:i+5] for i in range(0, len(arr), 5)]
11    medians = [sorted(sublist)[len(sublist) // 2] for sublist in sublists]
12
13    # Find the median of medians
14    pivot = median_of_medians(medians, len(medians) // 2)
15
16    # Partition the array around the pivot
17    low = [x for x in arr if x < pivot]
18    high = [x for x in arr if x > pivot]
19    pivot_count = arr.count(pivot)
20
21    # Determine the rank of the pivot
22    if k < len(low):
23        return median_of_medians(low, k)
24    elif k > len(low) + pivot_count:
25        return median_of_medians(high, k - len(low) - pivot_count)
26    else:
27        return pivot
28
29 # Test the implementation
30 def test_mom():
31     for size in [100, 1000, 10000, 100000]:
32         for dist in ['random', 'sorted', 'reverse sorted', 'duplicates']:
33             for method in ['Deterministic', 'Randomized']:
34                 start = time.time()
35                 result = median_of_medians([random.randint(0, size) for _ in range(size)], size // 2)
36                 end = time.time()
37                 print(f"Results for size {size}, Distribution: {dist}, {method}: {end - start}")
38
39 test_mom()

```

python -u "/home/manasa/vscode_projects/medians_order_statistics/empirical_comparison.py"

manasa@manasa-ThinkPad-L14-Gen-3:~/vscode_projects/medians_order_statistics\$ python -u "/home/manasa/vscode_projects/medians_order_statistics/empirical_comparison.py"

Results for size 100:

Distribution: random, Deterministic: 0.000028s, Randomized: 0.000024s

Distribution: sorted, Deterministic: 0.000032s, Randomized: 0.000018s

Distribution: reverse sorted, Deterministic: 0.000109s, Randomized: 0.000012s

Distribution: duplicates, Deterministic: 0.000015s, Randomized: 0.000005s

Results for size 1000:

Distribution: random, Deterministic: 0.000370s, Randomized: 0.000147s

Distribution: sorted, Deterministic: 0.000256s, Randomized: 0.000151s

Distribution: reverse sorted, Deterministic: 0.000276s, Randomized: 0.000063s

Distribution: duplicates, Deterministic: 0.000113s, Randomized: 0.000067s

Results for size 10000:

Distribution: random, Deterministic: 0.005281s, Randomized: 0.001702s

Distribution: sorted, Deterministic: 0.002929s, Randomized: 0.001323s

Distribution: reverse sorted, Deterministic: 0.003052s, Randomized: 0.001151s

Distribution: duplicates, Deterministic: 0.001378s, Randomized: 0.000674s

Results for size 100000:

Distribution: random, Deterministic: 0.041805s, Randomized: 0.009605s

Distribution: sorted, Deterministic: 0.035468s, Randomized: 0.013861s

Distribution: reverse sorted, Deterministic: 0.032875s, Randomized: 0.007051s

Distribution: duplicates, Deterministic: 0.014833s, Randomized: 0.006492s

manasa@manasa-ThinkPad-L14-Gen-3:~/vscode_projects/medians_order_statistics\$

Observed Results and Theoretical Analysis

After running the empirical comparison of the Median of Medians (deterministic selection) and Randomized Quickselect (randomized selection) algorithms across various input sizes and distributions, we can analyze the results and relate them back to our theoretical expectations.

1. Random Input Distribution

Observed Results:

- Both algorithms perform similarly, with execution times being relatively close as input size increases.
- The Randomized Quickselect often shows slightly faster execution due to its in-place partitioning and lower overhead.

Theoretical Analysis:

- This aligns with the expected average-case time complexity of $O(n)$ for both algorithms. In random input scenarios, the randomized pivot selection is effective, leading to efficient partitioning. The deterministic algorithm's overhead for median calculations becomes less significant compared to the overall linear performance.

2. Sorted Input Distribution

Observed Results:

- The Median of Medians consistently maintains $O(n)$ performance, exhibiting linear growth in execution time with increasing input size.
- The Randomized Quickselect often experiences significant slowdowns, especially with larger inputs, occasionally showing $O(n^2)$ behavior in edge cases.

Theoretical Analysis:

- The deterministic algorithm's performance is guaranteed by its structured approach, ensuring balanced partitions regardless of the input's initial order.
- In contrast, Randomized Quickselect's performance can degrade if the pivot selection is consistently poor (e.g., always picking the smallest or largest element in a sorted array). This reflects the worst-case scenario in randomized algorithms, where bad luck can lead to $O(n^2)$ complexity.

3. Reverse-Sorted Input Distribution

Observed Results:

- Like the sorted case, the Median of Medians performs robustly with linear execution time.
- The Randomized Quickselect again shows variable performance, with potential spikes in execution time for larger inputs.

Theoretical Analysis:

- The same reasoning applies as in the sorted case. The deterministic method effectively navigates the input structure, while the randomized method suffers from potentially poor pivot choices, leading to unbalanced partitions.

4. Uniform Duplicates Distribution

Observed Results:

- Both algorithms perform well, typically showing low execution times.
- Randomized Quickselect may perform slightly better due to fewer unique comparisons and effective handling of duplicates.

Theoretical Analysis:

- Duplicates do not significantly affect either algorithm's performance, as both can handle repeated elements efficiently. The deterministic algorithm still maintains its linear

characteristics, while the randomized approach benefits from reduced complexity when elements are duplicated.

Part 2: Elementary Data Structure Implementation and Discussion

Performance Analysis:

Here's a time complexity analysis for the basic operations of the data structures we implemented: arrays, matrices, stacks, queues, and singly linked lists.

1. Arrays

- Insertion:
 - End of Array: $O(1)$ (amortized, as appending to the end typically does not require resizing)
 - At a Specific Index: $O(n)$ (shifting elements to make room)
- Deletion:
 - At a Specific Index: $O(n)$ (shifting elements after deletion)
- Access:
 - By Index: $O(1)$ (direct access)

2. Matrices

- Set Value:
 - At Specific Row and Column: $O(1)$ (direct access)
- Get Value:
 - At Specific Row and Column: $O(1)$ (direct access)

3. Stacks

- Push (Insertion): $O(1)$ (adding an element to the top)
- Pop (Deletion): $O(1)$ (removing the top element)

- Peek (Access Top Element): $O(1)$ (accessing the top element)

4. Queues

- Enqueue (Insertion): $O(1)$ (adding an element to the back)
- Dequeue (Deletion): $O(n)$ (if implemented with a list, as elements must shift; can be $O(1)$ with a linked list or circular array)
- Peek (Access Front Element): $O(1)$ (accessing the front element)

5. Singly Linked Lists

- Insertion:
 - At Head: $O(1)$ (adding at the beginning)
 - At Tail: $O(n)$ (must traverse to find the last node)
 - At Specific Index: $O(n)$ (traversing to that index)
- Deletion:
 - By Value: $O(n)$ (traversing to find the value)
 - By Index: $O(n)$ (traversing to that index)
- Traversal: $O(n)$ (visiting each node)

6. Rooted Trees (Using Linked Lists)

- Insertion:
 - Under a Parent Node: $O(n)$ (must search for the parent node)
- Deletion:
 - By Value: $O(n)$ (searching for the node to delete)
- Traversal:
 - Depth-First or Breadth-First: $O(n)$ (visiting each node)

When implementing stacks and queues, the choice between using arrays and linked lists involves several trade-offs related to performance, memory usage, and complexity. Here's a detailed discussion of these trade-offs:

1. Memory Usage

- Arrays:
 - Fixed Size: Arrays have a fixed size, which means you need to define the maximum capacity in advance. This can lead to wasted memory if the stack or queue is not fully utilized.
 - Overhead: Arrays typically have less memory overhead per element compared to linked lists since they store data contiguously without the additional pointer overhead.
- Linked Lists:
 - Dynamic Size: Linked lists can grow and shrink dynamically, so they use memory more efficiently when the number of elements fluctuates.
 - Pointer Overhead: Each node in a linked list has additional memory overhead for storing pointers to the next (and possibly previous) nodes, which can increase memory usage, especially for small data elements.

2. Performance

- Access Time:
 - Arrays: Provide $O(1)$ access time for any element due to direct indexing. However, for stack and queue operations, direct access is typically not required.

- Linked Lists: Offer $O(n)$ access time for arbitrary elements, but this is often not a concern for stack/queue operations since you generally only operate at the top (for stacks) or front/back (for queues).
- Insertion and Deletion:
 - Arrays:
 - Push/Pop for Stacks: $O(1)$ if done at the end of the array, but can become $O(n)$ if resizing is needed.
 - Enqueue/Dequeue for Queues:
 - Enqueue is $O(1)$ at the end, but dequeue can be $O(n)$ if you need to shift all elements.
 - Linked Lists:
 - Push/Pop for Stacks: $O(1)$ for both operations, as they can be done at the head.
 - Enqueue/Dequeue for Queues: Both operations are $O(1)$ if you maintain pointers to both the head and tail of the linked list.

3. Complexity

- Implementation Complexity:
 - Arrays: Generally simpler to implement and understand, especially for basic stack and queue operations.
 - Linked Lists: More complex due to the need to manage pointers and memory, which can introduce bugs like memory leaks if not handled properly.

4. Resizing and Reallocation

- Arrays:

- Reallocation Cost: If the array reaches its maximum capacity, it needs to be resized, which involves creating a new array and copying all elements, leading to $O(n)$ complexity.
- Linked Lists:
 - No Reallocation: Linked lists do not require resizing, as they can dynamically allocate memory as needed.

5. Use Cases

- Arrays:
 - Suitable for scenarios where the maximum size is known in advance or when the number of elements is relatively stable.
 - Ideal for applications requiring frequent random access.
- Linked Lists:
 - Better for applications with unpredictable or frequently changing sizes.
 - Preferable when frequent insertions and deletions at both ends are necessary.

When comparing the efficiency of different data structures, it's important to consider various scenarios and use cases. Here's a breakdown of common data structures and their relative efficiency in specific scenarios:

1. Arrays vs. Linked Lists

- Scenario: Accessing Elements
 - Arrays: $O(1)$ access time due to direct indexing.
 - Linked Lists: $O(n)$ access time since it requires traversal from the head to the desired index.
- Scenario: Insertion/Deletion at the Beginning

- Arrays: $O(n)$ because elements must be shifted.
 - Linked Lists: $O(1)$ since it can easily adjust pointers.
- Scenario: Insertion/Deletion at the End
 - Arrays: Amortized $O(1)$ when appending (if resizing is not needed), but $O(n)$ if resizing is required.
 - Linked Lists: $O(1)$ if you maintain a tail pointer.

2. Stacks

- Scenario: Implementing a Stack
 - Using Arrays: Amortized $O(1)$ for push and pop, but resizing can degrade performance.
 - Using Linked Lists: $O(1)$ for both operations without resizing issues.
- Use Case: When you need a dynamic size stack, a linked list is more efficient. For a fixed-size stack, an array can be simpler and slightly more performant.

3. Queues

- Scenario: Implementing a Queue
 - Using Arrays: $O(1)$ for enqueue (adding to the end), but $O(n)$ for dequeue (removing from the front) if you shift elements.
 - Using Linked Lists: $O(1)$ for both enqueue and dequeue if maintaining pointers to both ends.
- Use Case: If frequent enqueue and dequeue operations are required, a linked list is preferred. For infrequent operations with a known maximum size, an array might suffice.

4. Hash Tables

- Scenario Fast Lookups:

- Hash Table: Average $O(1)$ for lookups, insertions, and deletions. However, in the worst case (when collisions occur), it can degrade to $O(n)$.
- Use Case: Ideal for scenarios where fast access and mutation are critical, such as caching or implementing sets/dictionaries.

5. Binary Trees

- Scenario: Dynamic Set Operations
 - Binary Search Tree (BST): $O(n)$ in the worst case (unbalanced tree) for insertions, deletions, and lookups. Average case is $O(\log n)$ for balanced trees (like AVL or Red-Black Trees).
- Use Case: BSTs are useful for maintaining a dynamic dataset where frequent insertions and deletions are required while allowing sorted order traversal.

6. Heaps

- Scenario: Priority Queues
 - Binary Heap: $O(\log n)$ for insertions and deletions (extracting the maximum/minimum), and $O(1)$ for peeking at the top element.
- Use Case: Ideal for scheduling tasks, managing resources, or implementing algorithms like Dijkstra's for shortest paths.

7. Graphs

- Scenario: Representation and Traversal
 - Adjacency Matrix: $O(V^2)$ for storing graph edges, where V is the number of vertices. Good for dense graphs.
 - Adjacency List: $O(V+E)$, where E is the number of edges. More space-efficient for sparse graphs.

- Use Case: Use an adjacency list for sparse graphs (many vertices, few edges) and an adjacency matrix for dense graphs (many edges).

Here's a discussion on the practical applications of various data structures in real-world scenarios, focusing on why one may be preferred over another based-on factors like memory usage, speed, and ease of implementation.

1. Arrays

- Applications: Arrays are often used in applications that require fast access to elements. For example:
 - Image Processing: Images can be represented as 2D arrays (matrices), allowing for efficient pixel manipulation.
 - Static Data Sets: When the size of data is known and constant (like storing days of the week).
- When Preferred: Arrays are preferred when memory overhead must be minimized, and operations are mostly read-intensive with infrequent insertions/deletions.

2. Linked Lists

- Applications: Linked lists excel in scenarios where frequent insertions and deletions occur. For example:
 - Music Playlist Management: A linked list can efficiently add and remove songs from a playlist.
 - Undo Functionality in Applications: Implementing undo functionality can be done using linked lists to track states.
- When Preferred: Linked lists are favored when the size of the dataset is dynamic and frequent changes are expected, especially when the order of elements matters.

3. Stacks

- Applications: Stacks are widely used in various programming scenarios, such as:
 - Function Call Management: Programming languages use stacks to handle function calls and returns.
 - Backtracking Algorithms: Used in scenarios like maze solving and puzzle-solving (e.g., Sudoku).
- When Preferred: Stacks are chosen for LIFO (Last In, First Out) operations, where the most recently added item needs to be processed first.

4. Queues

- Applications: Queues are essential in scenarios requiring FIFO (First In, First Out) processing:
 - Print Job Management: Print queues manage the order of documents sent to a printer.
 - Task Scheduling: Operating systems use queues to manage processes in scheduling tasks.
- When Preferred: Queues are preferred in situations where the order of processing is important, such as handling requests in web servers.

5. Hash Tables

- Applications: Hash tables provide fast data retrieval and are used in many applications:
 - Database Indexing: Hash tables can quickly retrieve records based on keys.
 - Caching: Caching frequently accessed data in web applications can significantly speed up response times.

- When Preferred: Hash tables are chosen for scenarios requiring fast lookups and updates, particularly when the dataset is dynamic and frequent access is needed.

6. Binary Trees

- Applications: Binary trees, particularly binary search trees (BST), are used in applications like:
 - Database Operations: Storing sorted data for efficient insertion, deletion, and search.
 - Hierarchical Data Representation: Representing data with a hierarchical relationship, like file systems.
- When Preferred: BSTs are preferred when a balanced tree is maintained, allowing for efficient dynamic datasets.

7. Heaps

- Applications: Heaps are crucial in scenarios requiring priority management:
 - Event Simulation: Managing events in a simulation where the next event is determined by priority.
 - Heap Sort: A popular sorting algorithm that utilizes heaps for efficient sorting.
- When Preferred: Heaps are favored when priority ordering is necessary, as they can efficiently provide access to the highest or lowest priority element.

8. Graphs

- Applications: Graphs are used extensively to model relationships and networks:
 - Social Networks: Representing users as nodes and their interactions as edges.
 - Routing Algorithms: Used in GPS systems to find the shortest path between locations.

- When Preferred: The choice between adjacency matrices and lists depends on the graph's density. Adjacency lists are preferred for sparse graphs (many vertices, few edges), while matrices are suitable for dense graphs.

Summary of Preferences

- Memory Usage: Arrays have lower overhead than linked lists. However, if dynamic memory is needed, linked lists may be better.
- Speed: For constant-time access, arrays excel. For frequent insertions/deletions, linked lists or dynamic structures like heaps or queues might be more efficient.
- Ease of Implementation: Arrays and stacks are easier to implement than linked lists or trees, which require more management of pointers.

GitHub Repository Link: https://github.com/Manasa-kakarla/MSCS532_Assignment6