

295. Find Median from Data Stream

HardDesignTwo PointersDataStreamSortingHeap (Priority Queue)

Leetcode Link

Problem Description

The problem requires designing a class, `MedianFinder`, that can handle a stream of integers and provide a way to find the median at any given time. Median, by definition, is the middle value in an ordered list of numbers. If the list has an odd number of elements, the median is simply the middle element. If the list has an even number of elements, there is no single middle element, so the median is calculated by taking the average of the two middle numbers.

Intuition

To find the median efficiently, we need a data structure that allows quick access to the middle elements. Utilizing two heaps is an elegant solution: a max heap to store the smaller half of the numbers and a min heap to store the larger half. This way, the largest number in the smaller half or the smallest number in the larger half can easily give us the median.

In Python, the default `heapq` module provides a min heap implementation. To get a max heap behavior, we insert negatives of the numbers into a heap.

By balancing the heaps so that their size differs by at most one, we ensure that we either have a single middle element when the combined size is odd (this will be at the top of the larger heap), or two middle elements when the combined size is even (the top of each heap).

The `addNum` method works by adding a new number to the max heap (smaller half) first, by pushing its negative value. We then pop the top value from the max heap and push it to the min heap (larger half) to maintain the order and balance. If the larger half has more than one extra element compared to the smaller half, we move the top element from the larger half to the smaller half.

`findMedian` checks the current size of the heaps. If the heaps are of the same size, the median is the average of the top values of both heaps. If the sizes differ, the median is the top element of the larger heap.

Solution Approach

The solution is based on maintaining two heaps (`h1` and `h2`), which are used to simulate a max heap and a min heap respectively.

- `h1` (max heap) stores the smaller half of the numbers, and we simulate it by pushing negatives of the numbers into a min heap.
- `h2` (min heap) stores the larger half of the numbers as they are, since the `heapq` library already implements a min heap.

Here's how the methods of the `MedianFinder` class work:

- The `__init__` method simply initializes two empty lists, `h1` and `h2`, that we'll use as heaps.
- `addNum` adds a new number to the data structure:
 - We first add the number to `h1` (which is simulating a max heap) by pushing its negative.
 - Then, we balance the max heap by popping an element from `h1` and pushing it onto `h2` (min heap).
 - We check if `h2` (larger half) has more than one element than `h1` (smaller half). If it does, we balance the heaps by moving the top element from `h2` to `h1`.
- `findMedian` computes the median based on the current elements:
 - If `h2` has more elements than `h1`, the median is just the top element of `h2` (the smallest number in the larger half).
 - If `h1` and `h2` have the same number of elements, the median is the average of the top element of `h1` (the largest number in the smaller half, remember `h1` is storing negatives) and the top element of `h2` (the actual smallest number in the larger half).

The efficiency of adding numbers is $O(\log n)$ due to the heap operations, and finding the median is $O(1)$ since it involves only accessing the tops of the heaps. This solution is quite efficient and handles the streaming data well.

Example Walkthrough

Let's walk through a small example to illustrate the implementation of the `MedianFinder` class using the given solution approach:

- Initialize the `MedianFinder` class, so we start with two empty heaps, `h1` (max heap) and `h2` (min heap).
- We execute `addNum(3)`. Here's how it works:
 - We insert -3 into `h1` to keep track of the max heap.
 - We pop -3 from `h1` (which is 3 in its original form) and push it to `h2`.
 - Since `h1` is now empty and `h2` has only one element, no balancing is needed.
- Next, we execute `addNum(1)`:
 - We insert -1 into `h1`.
 - We pop -1 from `h1` (which is 1 in its original form) and push it to `h2`.
 - Now `h2` has two elements (1 and 3), we pop 1 from `h2` and push its negative into `h1` to balance the heaps.
- At this point, `h1` has the single element -1, and `h2` has the single element 3.
- We add another number by calling `addNum(5)`:
 - We insert -5 into `h1`.
 - We pop -5 from `h1` (which is 5 in its original form) and push it into `h2`.
 - `h2` now having 3 and 5, is larger than `h1` which only has -1. No additional balancing is needed.
- Now, when calling `findMedian()`:
 - Since we have a total of 3 elements and `h1` has 1 element and `h2` has 2 elements, the median is the top of `h2` which is 3.
- Let's add another number by executing `addNum(2)`:
 - We insert -2 into `h1`.
 - We pop -2 from `h1` (which is 2 in its original form) and push it to `h2`.
 - Since `h2` (with elements 2, 3, and 5) now has more elements, we pop 2 from `h2` and push its negative into `h1` to balance the heaps. Now `h1` has -2 and -1, and `h2` has 3 and 5.
- Calling `findMedian()` now gives us:
 - Since `h1` and `h2` have the same number of elements (2 each), the median is the average of the top element of `h1` (-1, which is 1 when we negate it) and the top element of `h2` (which is 3), giving us a median of $(1 + 3) / 2 = 2$.

Through these steps, we can see how the `MedianFinder` class handles the addition of numbers and maintains a structure that allows us to easily find the median at any point. The use of two heaps is instrumental in efficiently managing the median calculation for a stream of data.

Python Solution

```
1 from heapq import heappush, heappop
2
3 class MedianFinder:
4     def __init__(self):
5         """
6         Initialize the MedianFinder data structure.
7         self.small - a max heap storing the smaller half of the numbers
8         self.large - a min heap storing the larger half of the numbers
9         """
10        self.small = [] # Max heap (simulated with negative values)
11        self.large = [] # Min heap
12
13    def addNum(self, num: int) -> None:
14        """
15        Add a number into the data structure.
16        :type num: int
17        :rtype: None
18        """
19        # Add to max heap
20        heappush(self.small, -num)
21        # Ensure the smallest element in max heap is not greater than
22        # the largest element in min heap
23        heappush(self.large, -heappop(self.small))
24
25        # Balance the heaps so that the min heap is not larger
26        # than the max heap by more than one element
27        if len(self.large) > len(self.small):
28            heappush(self.small, -heappop(self.large))
29
30    def findMedian(self) -> float:
31        """
32        Find and return the median of all elements so far.
33        :rtype: float
34        """
35        # If the heaps are of different sizes, the larger heap has the median
36        if len(self.large) > len(self.small):
37            return self.large[0]
38        # If the heaps are the same size, the median is the average of the
39        # two middle values
40        return (self.large[0] - self.small[0]) / 2.0
41
42
43 # How to use the MedianFinder class:
44 # medianFinder = MedianFinder()
45 # medianFinder.addNum(1)
46 # medianFinder.addNum(2)
47 # median = medianFinder.findMedian() -> returns 1.5
48
```

Java Solution

```
1 import java.util.PriorityQueue;
2 import java.util.Collections;
3
4 class MedianFinder {
5     private PriorityQueue<Integer> minHeap = new PriorityQueue<>(); // Min-heap to store the larger half of the numbers
6     private PriorityQueue<Integer> maxHeap = new PriorityQueue<>(Collections.reverseOrder()); // Max-heap to store the smaller half c
7
8     /** Initialize MedianFinder. */
9     public MedianFinder() {
10        // The constructor is kept empty as there's nothing to initialize outside the declarations.
11    }
12
13    /** Adds a number into the data structure. */
14    public void addNum(int num) {
15        minHeap.offer(num); // Add the number to the min-heap
16        maxHeap.offer(minHeap.poll()); // Balance the heaps by moving the smallest number of min-heap to max-heap
17
18        // Ensure max-heap has equal or one more element than the min-heap
19        if (maxHeap.size() > minHeap.size() + 1) {
20            minHeap.offer(maxHeap.poll()); // Move the maximum number of max-heap to min-heap
21        }
22    }
23
24    /** Returns the median of current data stream. */
25    public double findMedian() {
26        if (maxHeap.size() > minHeap.size()) {
27            // If max-heap has more elements, the median is the top of the max-heap
28            return maxHeap.peek();
29        }
30        // Otherwise, the median is the average of the tops of both heaps
31        return (minHeap.peek() + maxHeap.peek()) / 2.0;
32    }
33 }
34
35 /**
36  * Example of usage:
37  * MedianFinder medianFinder = new MedianFinder();
38  * medianFinder.addNum(num); // Add a number
39  * double median = medianFinder.findMedian(); // Get the current median
40  */
41
```

C++ Solution

```
1 #include <queue>
2 #include <vector>
3
4 class MedianFinder {
5 public:
6     MedianFinder() {
7         // Constructor doesn't need any code,
8         // as the priority queues are automatically initialized.
9     }
10
11    // Inserts a number into the data structure.
12    void addNum(int num) {
13        // Add the new number to the max heap.
14        maxHeap.push(num);
15
16        // Now balance the heaps by always having the top of the max heap
17        // ready to move to the min heap to maintain the ordering.
18        minHeap.push(maxHeap.top());
19        maxHeap.pop();
20
21        // If the min heap has more elements than the max heap,
22        // move the top element back to the max heap to maintain
23        // size property of heaps.
24        if (minHeap.size() > maxHeap.size()) {
25            maxHeap.push(minHeap.top());
26            minHeap.pop();
27        }
28    }
29
30    // Finds and returns the median of all numbers inserted.
31    double findMedian() {
32        // If the max heap is larger, the median is at the top of the max heap.
33        if (maxHeap.size() > minHeap.size()) {
34            return maxHeap.top();
35        }
36        // If both heaps have the same size, the median is the average of
37        // the tops of both heaps.
38        return (double) (minHeap.top() + maxHeap.top()) / 2;
39    }
40
41 private:
42    // Max heap for the lower half of the numbers.
43    std::priority_queue<int> maxHeap;
44    // Min heap for the upper half of the numbers
45    // (using greater<> to make it a min heap).
46    std::priority_queue<int, std::vector<int>, std::greater<int>> minHeap;
47 };
48
49 /**
50  * The MedianFinder object will be instantiated and called as such:
51  * MedianFinder* obj = new MedianFinder();
52  * obj->addNum(num); // Method to add a number into the MedianFinder.
53  * double median = obj->findMedian(); // Method to find and return the median.
54  */
55
```

Typescript Solution

```
1 // An array to store the numbers in sorted order.
2 let numbersList: number[] = [];
3
4 /**
5  * Insert a number into the array in a sorted position.
6  * @param {number} num The number to be added.
7  */
8 function addNumber(num: number): void {
9     let left = 0;
10    let right = numbersList.length;
11    // Binary search to find the correct insertion position.
12    while (left < right) {
13        const mid = (left + right) >> 1;
14        if (numbersList[mid] < num) {
15            left = mid + 1;
16        } else {
17            right = mid;
18        }
19    }
20    // Insert the number at the determined position.
21    numbersList.splice(left, 0, num);
22 }
23
24 /**
25  * Calculates and returns the median of the list of numbers.
26  * @return {number} The median value.
27  */
28 function findMedian(): number {
29    const count = numbersList.length;
30    // Check if the count of numbers is odd.
31    if ((count & 1) === 1) {
32        // If odd, the middle is the middle element.
33        return numbersList[count >> 1];
34    }
35    // If even, the median is the average of the two middle elements.
36    return (numbersList[count >> 1] + numbersList[(count >> 1) - 1]) / 2;
37 }
38
39 // Usage example, calling the methods directly since they're global:
40 addNumber(5);
41 addNumber(1);
42 console.log(findMedian()); // Should output the current median of the list
43
```

Time and Space Complexity

The provided class `MedianFinder` has two methods: `addNum` and `findMedian`, used for adding numbers and finding the median, respectively. The class maintains two heaps: `h1` is the min-heap that stores the larger half of the numbers, and `h2` is the max-heap that stores the smaller half of the numbers (as negatives).

addNum Method

Time Complexity:

- Inserting an element into a heap (`heappush`) has a time complexity of $O(\log n)$, where `n` is the number of elements in the heap.
- Removing the smallest element from a min heap (`heappop` on `h1`) or the largest element from a max heap (`heappop` on `h2`) has a time complexity of $O(\log n)$.

- The `addNum` method does at most two push and pop operations each time it is called:

- `heappush(self.h1, num)`: Insert `num` into the min-heap `h1`.
- `heappush(self.h2, -heappop(self.h1))`: Pop the minimum value from `h1`, negate it (to simulate a max-heap), and push it into `h2`.

- `heappush(self.h1, -heappop(self.h2))`: This happens only if the size of `h2` exceeds the size of `h1` by more than one, ensuring the difference in the number of elements in the two heaps never exceeds one.

Considering the above steps, the time complexity for `addNum` is $O(\log n)$ because of the heap operations.

Space Complexity:

- The space complexity of the `MedianFinder` class is $O(n)$, where `n` is the number of elements inserted into the `MedianFinder`. This is due to the storage of all elements across two heaps.

findMedian Method

Time Complexity:

- Accessing the top element of a heap (minimum value in `h1` and maximum value in `h2`) has a time complexity of $O(1)$.
- `findMedian` performs at most one arithmetic operation, which is done in constant time.

Therefore, the time complexity for `findMedian` is $O(1)$.

Space Complexity:

- The `findMedian` method does not use any additional space that depends on the input size, so its space complexity is $O(1)$.