

Documentation

The hIndex problem is a classic algorithmic challenge that involves calculating a researcher's h-index from a sorted list of citation counts. In this context, the h-index is defined as the maximum number (h) such that the researcher has published (h) papers with each having at least (h) citations. This metric helps quantify the impact of a researcher's work by balancing the quantity of publications with their quality in terms of citation count. The unique requirement in this problem is to achieve a solution with logarithmic time complexity, leveraging the fact that the citation list is sorted in ascending order. This sorting allows us to optimize the search process significantly using binary search techniques, making it feasible to handle large lists efficiently.

To approach this problem, it's essential to understand the relationship between the index of a paper in the sorted array and the number of papers cited at least as many times as the citation count at that index. Given the sorted nature of the list, we can interpret any citation count at index i in terms of how many papers have equal to or greater than this citation count. Specifically, if there are (n) papers in total, then $n - i$ gives the count of papers with citations equal to or higher than $\text{citations}[i]$. This value, $n - i$, represents a potential h-index candidate if the citation count at that position meets or exceeds $n - i$, thus establishing a valid balance between publication quantity and quality.

The optimal solution utilizes binary search, a powerful tool for logarithmic-time searches in sorted lists. By initializing two pointers, left and right, at the start and end of the array, respectively, we can repeatedly divide the array to zero in on the correct h-index value. For each midpoint index, we calculate the potential h-index $n - \text{mid}$, where mid is the midpoint of the current search range. We then check if $\text{citations}[\text{mid}]$ (the number of citations for the paper at mid) is sufficient to make $n - \text{mid}$ a valid h-index. If it is, this indicates that we might be able to increase the h-index, so we adjust our search to the left half by setting right to $\text{mid} - 1$. Otherwise, we move our search to the right half by setting left to $\text{mid} + 1$, thereby ignoring citation counts that are insufficient for a valid h-index.

The primary insight of the binary search approach is that it identifies the first position in the array where the citation count meets the h-index criteria without needing to sequentially check every paper. By the end of the search, the left will point to the minimum index satisfying $\text{citations}[\text{left}] \geq n - \text{left}$, providing the largest possible h-index. The h-index can then be derived as $n - \text{left}$, since left will have converged to the smallest index meeting the h-index condition. This approach effectively ensures the solution's logarithmic time complexity, crucial for handling large datasets up to the problem's upper limit.

The solution's time complexity is $(O(\log n))$, owing to the binary search method, which divides the search space in half at each step. This efficiency is especially important given that the number of citations can reach 100,000, making any linear-time solution too slow. Additionally, the space complexity is $(O(1))$, as the algorithm operates directly on the input list without needing additional data structures. This design ensures the algorithm is both time-efficient and space-efficient, suitable for applications with large volumes of citation data.

In summary, the binary search approach to finding the h-index from a sorted list of citations is an elegant solution that leverages the sorted nature of the data to achieve optimal performance. By focusing on the relationship between the citation counts and the number of papers with equal or higher citations, the algorithm efficiently narrows down the potential h-index, arriving at the correct answer in logarithmic time. This method highlights the power of binary search for optimizing problems that might initially seem to require linear checks, and it provides a practical solution for determining h-indices in bibliometric analyses.