

692. Top K Frequent Words

Medium Trie Hash Table String Bucket Sort Counting Sorting Heap (Priority Queue) Leetcode Link

Problem Description

The problem requires us to find the k most frequent strings in an array of strings `words`. This entails two primary tasks: calculating the frequency of each word and then finding the top k based on this frequency. The problem adds a layer of complexity by asking for these top k strings to be in a specific order. The words should first be sorted based on frequency, with the most frequent words appearing first. If there are multiple words with the same frequency, those words should be ordered lexicographically (like in a dictionary), which means alphabetically in ascending order.

The expected outcome is a list of strings representing the k most frequent words in the order described above.

Intuition

The intuition behind the solution is to employ data structures and sorting mechanisms that can efficiently handle the two-level sorting criteria described in the problem statement.

The solution uses the following steps to arrive at the desired result:

- Count the occurrences of each word using a data structure that maps words to their frequencies. Python's `Counter` class from the `collections` module is well-suited for this and helps us quickly determine the frequency of each word in `words`.
- Sort the words based on the two criteria using a custom sorting function:
 - First, by the frequency of each word in descending order so that the most frequent words come first.
 - Second, by the word itself in lexicographical order for words with the same frequency. The lexicographical order ensures that if two words have the same frequency, the word that comes first alphabetically will come first in the sorted list.
- To sort by multiple criteria, we can use a tuple where the first element is the negative frequency (to sort in descending order) and the second element is the word itself (to sort lexicographically as a tiebreaker).
- Finally, we slice the sorted list to only include the top k elements.

Using this approach, we can leverage the powerful built-in sorting available in Python, ensuring that both sorting criteria are met, and only k elements are returned as the final result.

Solution Approach

The implementation of the solution takes advantage of Python's libraries and language features to accomplish the task of finding the top k most frequent strings. Here's a step-by-step explanation of the methods used:

Step 1: Counting Frequencies

We first use the `Counter` class from Python's `collections` module, which is a specialized dictionary designed for counting hashable objects. It takes an iterable (in this case, the `words` list) and creates a `Counter` object that maps each unique word to its frequency count.

```
1 cnt = Counter(words)
```

Step 2: Custom Sorting

Once we have the frequency counts, the next step is to sort the items based on our two criteria. Python's `sorted` function is utilized here, which allows for custom sorting logic through its `key` parameter. The `key` parameter accepts a function that returns a tuple containing the sorting criteria for each element in the iterable.

```
1 sorted(cnt, key=lambda x: (-cnt[x], x))[:k]
```

In the lambda function:

- `x` is each word from the `Counter` keys.
- `-cnt[x]` is the frequency of the word, negated so that the list is sorted in descending order (since Python sorts in ascending order by default).
- `x` is the word itself, which ensures that words with identical frequencies are sorted lexicographically.

The tuple `(-cnt[x], x)` ensures that the sorting happens first by frequency and then, as a tiebreaker, by lexicographical order.

Step 3: Slicing

Lastly, after sorting, we slice the list to only select the first k elements, as we're only interested in the top k most frequent words.

```
1 [:k]
```

This represents the final step where the method only returns the first k sorted elements, effectively giving us the top k most frequent words in the required order.

Combining these steps, the solution applies well-known algorithms for counting (hash table via `Counter`) and sorting (with custom criteria via `sorted`) to efficiently resolve the problem at hand. The use of lambda expressions and list slicing demonstrates the expressive power of Python's syntax, allowing for concise and readable code.

Example Walkthrough

Let's consider the following small example array of strings `words` and a value for k to illustrate the solution approach:

```
words = ["apple", "banana", "cherry", "apple", "banana", "apple"] k = 2
```

We want to find the 2 most frequent strings in the `words` array.

Step 1: Counting Frequencies

Using Python's `Counter` from the `collections` module, calculate the frequency of each word.

```
1 from collections import Counter
2 cnt = Counter(words)
3 # cnt = {'apple': 3, 'banana': 2, 'cherry': 1}
```

The `Counter` object now holds the frequency of each word: "apple" appears 3 times, "banana" 2 times, and "cherry" 1 time.

Step 2: Custom Sorting

Now we need to sort the words based on two criteria: the frequency (higher to lower) and lexicographic order. We use a custom sorting function with Python's `sorted`:

```
1 sorted_words = sorted(cnt, key=lambda x: (-cnt[x], x))
2 # sorted_words = ['apple', 'banana', 'cherry']
```

Here, `sorted_words` is sorted by frequency in descending order. Since "apple" and "banana" have different frequencies, no need for lexicographical sorting between them. "Cherry" has a lower frequency so comes last.

Step 3: Slicing

Finally, we only want the top k elements, so we slice the sorted list:

```
1 result = sorted_words[:k]
2 # result = ['apple', 'banana']
```

The result list `['apple', 'banana']` represents the 2 most frequent words in the `words` array, adhering to the frequency then lexicographical ordering rules.

Combining Step 1 through Step 3, the complete code to accomplish finding the k most frequent strings would look like this:

```
1 from collections import Counter
2
3 def k_most_frequent(words, k):
4     cnt = Counter(words)
5     sorted_words = sorted(cnt, key=lambda x: (-cnt[x], x))
6     return sorted_words[:k]
7
8 # Usage
9 words = ["apple", "banana", "cherry", "apple", "banana", "apple"]
10 k = 2
11 print(k_most_frequent(words, k)) # Output: ['apple', 'banana']
```

This example walkthrough demonstrates how the combined steps of counting, sorting, and slicing give us a clean solution using Python's powerful built-in features.

Python Solution

```
1 from collections import Counter
2 from typing import List
3
4 class Solution:
5     def topKFrequent(self, words: List[str], k: int) -> List[str]:
6         # Create a frequency counter for the words
7         word_count = Counter(words)
8
9         # Sort the unique words first by frequency in descending order, then alphabetically
10        sorted_words = sorted(word_count, key=lambda word: (-word_count[word], word))
11
12        # Return the first k elements of the sorted list which represent the top-k frequent words
13        return sorted_words[:k]
```

Java Solution

```
1 import java.util.*;
2
3 class Solution {
4     public List<String> topKFrequent(String[] words, int k) {
5         // Map to store the frequency count of each word
6         Map<String, Integer> wordCount = new HashMap<>();
7         // Calculate the frequency of each word
8         for (String word : words) {
9             wordCount.put(word, wordCount.getOrDefault(word, 0) + 1);
10        }
11
12        // Priority Queue to store words based on their frequency and lexicographical order
13        PriorityQueue<String> heap = new PriorityQueue<>((word1, word2) -> {
14            int frequencyDifference = wordCount.get(word1) - wordCount.get(word2);
15            // If frequencies are the same, compare words in reverse lexicographical order
16            if (frequencyDifference == 0) {
17                return word2.compareTo(word1);
18            }
19            // Otherwise, order by frequency
20            return frequencyDifference;
21        });
22
23        // Iterate over the distinct words and add them to the Priority Queue
24        // The heap will maintain the top k frequent elements on top
25        for (String word : wordCount.keySet()) {
26            heap.offer(word);
27            // If heap size exceeds k, remove the least frequent/current smallest element
28            if (heap.size() > k) {
29                heap.poll();
30            }
31        }
32
33        // LinkedList to store the result in correct order
34        LinkedList<String> topKWords = new LinkedList<>();
35        // Populating the result list in reverse order since we want the highest frequency on top
36        while (!heap.isEmpty()) {
37            topKWords.addFirst(heap.poll());
38        }
39
40        // Return the list of top k frequent words
41        return topKWords;
42    }
43 }
44
```

C++ Solution

```
1 #include <vector>
2 #include <string>
3 #include <unordered_map>
4 #include <algorithm> // Needed for sort function
5
6 using std::vector;
7 using std::string;
8 using std::unordered_map;
9
10 class Solution {
11 public:
12     // This function returns the top k frequent words from the vector of strings "words".
13     vector<string> topKFrequent(vector<string>& words, int k) {
14         // Hash map to store the frequency count for each word.
15         unordered_map<string, int> wordCount;
16         // Count frequency of each word.
17         for (const auto& word : words) {
18             ++wordCount[word];
19         }
20
21         // Vector to store the unique words for sorting.
22         vector<string> uniqueWords;
23         // Extract the unique words (keys) from the map.
24         for (const auto& pair : wordCount) {
25             uniqueWords.emplace_back(pair.first);
26         }
27
28         // Custom sort the uniqueWords vector based on the frequency and alphabetical order.
29         sort(uniqueWords.begin(), uniqueWords.end(), [&](const string& a, const string& b) {
30             // If frequencies are equal, sort by alphabetical order.
31             if (wordCount[a] == wordCount[b]) {
32                 return a < b;
33             }
34             // Otherwise, sort by frequency in descending order.
35             return wordCount[a] > wordCount[b];
36         });
37
38         // After sorting, keep only the top k elements in the uniqueWords vector.
39         uniqueWords.erase(uniqueWords.begin() + k, uniqueWords.end());
40
41         // Return the top k frequent words.
42         return uniqueWords;
43     }
44 };
45
```

Typescript Solution

```
1 import { sortBy } from 'lodash';
2
3 // Type alias for word and its frequency.
4 type WordFrequency = { [key: string]: number };
5
6 // This function counts the frequency of words in an array and returns an object where keys are words and values are frequencies.
7 function countWordFrequencies(words: string[]): WordFrequency {
8     const wordCount: WordFrequency = {};
9     words.forEach(word => {
10         wordCount[word] = (wordCount[word] || 0) + 1;
11     });
12     return wordCount;
13 }
14
15 // This function returns the top k frequent words from the array of strings "words".
16 function topKFrequent(words: string[], k: number): string[] {
17     // Use the countWordFrequencies function to get a frequency count for each word.
18     const wordCount: WordFrequency = countWordFrequencies(words);
19
20     // Extract the unique words (keys) from the wordCount object and sort them.
21     let uniqueWords: string[] = Object.keys(wordCount);
22     uniqueWords = sortBy(uniqueWords, [
23         // Sort by frequency in descending order.
24         word => -wordCount[word],
25         // If frequencies are equal, sort by alphabetical order.
26         word => word
27     ]);
28
29     // After sorting, keep only the top k elements in uniqueWords array.
30     return uniqueWords.slice(0, k);
31 }
32
33 // Insert the necessary imports or dependencies if required.
34 // Note: The 'sortBy' function is utilized from lodash for convenience and needs to be installed.
35
```

Time and Space Complexity

Time Complexity

The time complexity of the code can be analyzed in a few steps:

- `Counter(words)` has a time complexity of $O(n)$ where n is the number of words, as it requires one pass to count the frequency of each word.
- `sorted(cnt, key=lambda x: (-cnt[x], x))` has a time complexity of $O(n \log n)$ on average, as it sorts the list of words by their frequencies, and in case of ties, it sorts them lexicographically. Here, " n " refers to the unique words in the list. `3[:k]` slicing has a time complexity of $O(k)$ because it needs to copy the first k elements from the sorted array.

Overall, the time complexity of the solution is dominated by the sorting step, which, in the worst case, is $O(n \log n)$.

Space Complexity

For space complexity:

- The `Counter` uses $O(u)$ space where u is the number of unique words in the input list, to store the word counts.
- The sorted list will also use $O(u)$ space.
- The slice of the top k elements will require $O(k)$ space, but this does not add to the overall space complexity since k is no larger than u and so is subsumed by it.

Therefore, the overall space complexity of the function is $O(u)$.