

DAY-11 JAVA ASSIGNMENT

Task 3: Implementing the KMP Algorithm

Code the Knuth-Morris-Pratt (KMP) algorithm in JAVA for pattern searching which pre-processes the pattern to reduce the number of comparisons. Explain how this pre-processing improves the search time compared to the naive approach.

KMP Algorithm Explanation

The KMP algorithm improves the search time for pattern matching by avoiding unnecessary comparisons. It does this through a preprocessing step where it builds a partial match (also known as "prefix") table. This table is used to skip characters in the text string that we know will match the pattern string, thus reducing the number of comparisons needed.

Key Concepts

1. **Partial Match Table (LPS Array):** The LPS (Longest Prefix which is also Suffix) array is built for the pattern. For each position in the pattern, the LPS array contains the length of the longest proper prefix which is also a suffix.
2. **Search Process:** While searching through the text, if a mismatch occurs after a certain number of matches, the LPS array is used to determine the next positions to check in the pattern without rechecking characters that we know will match.

Preprocessing to Build LPS Array

The preprocessing step constructs the LPS array which helps to determine the next positions to check in case of a mismatch.

How Preprocessing Improves Search Time

In the naive approach, every mismatch after a match requires the pattern to be shifted one position to the right and start again from the beginning of the pattern. This results in a worst-case time complexity of $O(M * N)$, where M is the length of the pattern and N is the length of the text.

The KMP algorithm, however, utilizes the LPS array to skip unnecessary comparisons. When a mismatch occurs after j matches, the algorithm uses the LPS array to avoid rechecking characters that are known to match. This optimization results in an improved worst-case time complexity of $O(M + N)$.

- **Time Complexity:**
 - Preprocessing: $O(M)$ to build the LPS array.
 - Searching: $O(N)$ to search the pattern in the text using the LPS array.

Thus, the overall time complexity of the KMP algorithm is $O(M + N)$, making it significantly more efficient than the naive approach, especially for larger texts and patterns.

```
KMPAlgorithm.java x Console x
1 package com.wipro.patterns;
2
3 public class KMPAlgorithm {
4
5     public static void main(String[] args) {
6         String text = "ABABDABACDABABCABAB";
7         String pattern = "ABABCABAB";
8         search(text, pattern);
9     }
10
11     public static void search(String text, String pattern) {
12         int[] lps = computeLPSArray(pattern);
13         int i = 0; // index for text
14         int j = 0; // index for pattern
15         int M = pattern.length();
16         int N = text.length();
17
18         while (i < N) {
19             if (pattern.charAt(j) == text.charAt(i)) {
20                 i++;
21                 j++;
22             }
23
24             if (j == M) {
25                 System.out.println("Pattern found at index " + (i - j));
26                 j = lps[j - 1];
27             } else if (i < N && pattern.charAt(j) != text.charAt(i)) {
28                 if (j != 0) {
29                     j = lps[j - 1];
30                 } else {
31                     i++;
32                 }
33             }
34         }
35     }
36
37     private static int[] computeLPSArray(String pattern) {
38         int M = pattern.length();
39         int[] lps = new int[M];
40         int length = 0; // length of the previous longest prefix suffix
41         int i = 1;
42         lps[0] = 0; // lps[0] is always 0
43
44         while (i < M) {
45             if (pattern.charAt(i) == pattern.charAt(length)) {
46                 length++;
47                 lps[i] = length;
48                 i++;
49             } else {
50                 if (length != 0) {
51                     length = lps[length - 1];
52                 } else {
53                     lps[i] = 0;
54                     i++;
55                 }
56             }
57         }
58         return lps;
59     }
60 }
```

<terminated> KMPAlgorithm [Java A
Pattern found at index 10

Task 4: Rabin-Karp Substring Search

Implement the Rabin-Karp algorithm for substring search using a rolling hash. Discuss the impact of hash collisions on the algorithm's performance and how to handle them.

Rabin-Karp Algorithm Explanation

The Rabin-Karp algorithm uses hashing to find a pattern in a text. It computes a hash value for the pattern and a hash value for each substring of the text of the same length as the pattern. If the hash values match, it then checks the actual substring to confirm a match. This approach allows for efficient pattern searching.

Key Concepts

1. **Hash Function:** A good hash function distributes values uniformly and minimizes collisions.
2. **Rolling Hash:** To efficiently compute hash values for the sliding window of substrings in the text, the rolling hash technique is used. It allows the hash value of the next substring to be computed in constant time.

Handling Hash Collisions

- **Hash Collisions:** Two different strings may have the same hash value (a collision). To handle this, whenever a hash match is found, the actual substring is compared to the pattern to confirm the match.
- **Modulo Operation:** A prime number is used in the modulo operation to reduce the likelihood of collisions.

Impact of Hash Collisions

- **Hash Collisions:** When different substrings produce the same hash value, it leads to hash collisions. In the Rabin-Karp algorithm, hash collisions cause unnecessary substring comparisons, impacting performance.
- **Performance:** In the average case, the Rabin-Karp algorithm runs in $O(N + M)$ time. However, if there are many hash collisions, the performance degrades to $O(NM)$ in the worst case.

Handling Hash Collisions

- **Verification Step:** After finding a hash match, perform a direct string comparison to verify the match. This step ensures that false positives due to hash collisions are correctly handled.
- **Choice of q :** Use a large prime number for the modulo operation to minimize the probability of collisions.

Conclusion

The Rabin-Karp algorithm is efficient for multiple pattern searches and large texts due to its hashing and rolling hash mechanism. However, careful handling of hash collisions through verification and appropriate choice of hash function parameters is crucial for maintaining its efficiency.

```
RabinKrap.java ×
1 package com.wipro.patterns;
2
3 import java.util.ArrayList;
4 import java.util.List;
5
6 public class RabinKrap {
7
8     private static final int PRIME = 101;
9
10    public static List<Integer> search(String pattern, String text) {
11        List<Integer> occurrences = new ArrayList<>();
12
13        int patternLength = pattern.length();
14        int textLength = text.length();
15
16        int patternHash = calculateHash(pattern, patternLength);
17        int textHash = calculateHash(text, patternLength);
18
19        for (int i = 0; i <= textLength - patternLength; i++) {
20            if (patternHash == textHash && checkEqual(pattern, text, i)) {
21                occurrences.add(i);
22            }
23
24            if (i < textLength - patternLength) {
25                textHash = recalculateHash(text, i, patternLength, textHash);
26            }
27        }
28
29        return occurrences;
30    }
31
32    private static int calculateHash(String str, int length) {
33        int hash = 0;
34        for (int i = 0; i < length; i++) {
35            hash += str.charAt(i) * Math.pow(PRIME, i);
36        }
37        return hash;
38    }
39
40    private static int recalculateHash(String str, int oldIndex, int patternLength, int oldHash) {
41        int newHash = oldHash - str.charAt(oldIndex);
42        newHash /= PRIME;
43        newHash += str.charAt(oldIndex + patternLength) * Math.pow(PRIME, patternLength - 1);
44        return newHash;
45    }
46
47    private static boolean checkEqual(String pattern, String text, int startIndex) {
48        for (int i = 0; i < pattern.length(); i++) {
49            if (pattern.charAt(i) != text.charAt(startIndex + i)) {
50                return false;
51            }
52        }
53        return true;
54    }
55
56    public static void main(String[] args) {
57        String text = "ABCCDEFGHABC";
58        String pattern = "ABC";
59        List<Integer> occurrences = search(pattern, text);
60        if (occurrences.isEmpty()) {
61            System.out.println("Pattern not found in the text");
62        } else {
63            System.out.println("Pattern found at positions: " + occurrences);
64        }
65    }
66 }
```

```
Console ×
<terminated> RabinKrap [Java Application] C:\Pro...
Pattern found at positions: [0, 9]
```

Task 5: Boyer-Moore Algorithm Application

Use the Boyer-Moore algorithm to write a function that finds the last occurrence of a substring in a given string and returns its index. Explain why this algorithm can outperform others in certain scenarios.

Boyer-Moore Algorithm Explanation

The Boyer-Moore algorithm is efficient because it uses information gained from mismatches to skip sections of the text, potentially moving more than one character ahead with each step. It uses two heuristics:

1. **Bad Character Heuristic:** When a mismatch occurs, the algorithm skips alignments that cannot possibly match based on the character that caused the mismatch.
2. **Good Suffix Heuristic:** When a mismatch occurs, the algorithm uses information about the matched suffix to skip sections of the text where the same suffix would appear.

Why Boyer-Moore Can Outperform Other Algorithms

1. **Large Alphabet Advantage:** Boyer-Moore is particularly efficient for large alphabets. The bad character heuristic can skip large sections of the text, especially when the alphabet size is large.
2. **Mismatches:** The algorithm skips sections of the text more aggressively on mismatches compared to algorithms like KMP. This is especially beneficial when the pattern length is long and mismatches are frequent.
3. **Performance in Practice:** While the worst-case time complexity is $O(MN)$ (same as the naive algorithm), in practice, Boyer-Moore often performs much better than KMP and Rabin-Karp due to fewer comparisons. It is closer to $O(N/M)$ in many practical scenarios.
4. **Flexibility with Different Patterns:** The algorithm performs well even with varied and complex patterns due to the dual heuristics that complement each other.

Conclusion

The Boyer-Moore algorithm's use of the bad character and good suffix heuristics allows it to skip large parts of the text, making it more efficient in practice compared to other algorithms. This efficiency is especially notable with large alphabets and longer patterns, making it a powerful tool for string searching tasks.

```
BoyerMooreAlgorithm.java x
1 package com.wipro.patterns;
2
3 public class BoyerMooreAlgorithm {
4
5     // Function to find the last occurrence of a pattern in a text
6     public static int lastOccurrence(String txt, String pat) {
7         int[] badChar = badCharHeuristic(pat);
8         int[] goodSuffix = goodSuffixHeuristic(pat);
9
10        int m = pat.length();
11        int n = txt.length();
12
13        int s = 0; // s is the shift of the pattern with respect to text
14        int lastOccurrence = -1;
15
16        while (s <= n - m) {
17            int j = m - 1;
18
19            // Keep reducing index j of pattern while characters of pattern and text are matching
20            while (j >= 0 && pat.charAt(j) == txt.charAt(s + j)) {
21                j--;
22            }
23
24            // If the pattern is present at the current shift, update lastOccurrence
25            if (j < 0) {
26                lastOccurrence = s;
27                s += (s + m < n) ? m - badChar[txt.charAt(s + m)] : 1;
28            } else {
29                s += Math.max(1, j - badChar[txt.charAt(s + j)]);
30            }
31        }
32
33        return lastOccurrence;
34    }
35}
```

<terminated> BoyerMooreAlgorithm [Java Application] CA
Last occurrence of pattern is at index 10

```
36 private static int[] badCharHeuristic(String pat) {
37     int[] badChar = new int[256]; // Assuming ASCII character set
38     int m = pat.length();
39
40     for (int i = 0; i < 256; i++) {
41         badChar[i] = -1; // Initialize all occurrences as -1
42     }
43
44     for (int i = 0; i < m; i++) {
45         badChar[pat.charAt(i)] = i; // Fill the actual value of last occurrence of a character
46     }
47
48     return badChar;
49 }
50
51 private static int[] goodSuffixHeuristic(String pat) {
52     int m = pat.length();
53     int[] suffix = new int[m];
54     int[] shift = new int[m];
55
56     for (int i = 0; i < m; i++) {
57         suffix[i] = -1;
58     }
59
60     int lastPrefixPosition = m;
61     for (int i = m - 1; i >= 0; i--) {
62         if (isPrefix(pat, i + 1)) {
63             lastPrefixPosition = i + 1;
64         }
65         shift[m - 1 - i] = lastPrefixPosition - i + m - 1;
66     }
67
68     for (int i = 0; i < m - 1; i++) {
69         int slen = suffixLength(pat, i);
70         shift[slen] = m - 1 - i + slen;
71     }
72
73     return shift;
74 }
75}
```

```

76 private static boolean isPrefix(String pat, int p) {
77     int m = pat.length();
78     for (int i = p, j = 0; i < m; i++, j++) {
79         if (pat.charAt(i) != pat.charAt(j)) {
80             return false;
81         }
82     }
83     return true;
84 }
85
86 private static int suffixLength(String pat, int p) {
87     int m = pat.length();
88     int len = 0;
89     for (int i = p, j = m - 1; i >= 0 && pat.charAt(i) == pat.charAt(j); i--, j--) {
90         len++;
91     }
92     return len;
93 }
94
95 // Driver code
96 public static void main(String[] args) {
97     String txt = "abacaabadcabacabaabb";
98     String pat = "abacab";
99     int result = lastOccurrence(txt, pat);
100     if (result == -1) {
101         System.out.println("Pattern not found");
102     } else {
103         System.out.println("Last occurrence of pattern is at index " + result);
104     }
105 }
106 }
107

```