

Substring Search

Introduction

Find a pattern of length M in a text of length N

Application

- screen scraping
- spam detection

Brute Force

Basic

- check for pattern starting at each text position

Performance

- can be slow if text and pattern are repetitive. $\sim MN$ char compares
- Brute-force needs **backup** for every mismatch---should avoid backup in text stream

Implementation

```
1 public class Brute {
2
3     public static int search1(String pat, String txt) {
4         int m = pat.length();
5         int n = txt.length();
6
7         for (int i = 0; i <= n - m; i++) {
8             int j;
9             for (j = 0; j < m; j++) {
10                 if (txt.charAt(i+j) != pat.charAt(j))
11                     break;
12             }
13             if (j == m) return i;           // found at offset i
14         }
15         return n;                         // not found
16     }
17
18     // return offset of first match or N if no match
19     public static int search2(String pat, String txt) {
20         int m = pat.length();
21         int n = txt.length();
22         int i, j;
23         for (i = 0, j = 0; i < n && j < m; i++) {
24             if (txt.charAt(i) == pat.charAt(j)) j++;
25             else {
```

```

26         i -= j;
27         j = 0;
28     }
29 }
30 if (j == m) return i - m;    // found
31 else      return n;        // not found
32 }
33
34 public static int search1(char[] pattern, char[] text) {
35     int m = pattern.length;
36     int n = text.length;
37
38     for (int i = 0; i <= n - m; i++) {
39         int j;
40         for (j = 0; j < m; j++) {
41             if (text[i+j] != pattern[j])
42                 break;
43         }
44         if (j == m) return i;    // found at offset i
45     }
46     return n;                  // not found
47 }
48
49 // return offset of first match or n if no match
50 public static int search2(char[] pattern, char[] text) {
51     int m = pattern.length;
52     int n = text.length;
53     int i, j;
54     for (i = 0, j = 0; i < n && j < m; i++) {
55         if (text[i] == pattern[j]) j++;
56         else {
57             i -= j;
58             j = 0;
59         }
60     }
61     if (j == m) return i - m;    // found
62     else      return n;        // not found
63 }
64
65 public static void main(String[] args) {
66     String pat = args[0];
67     String txt = args[1];
68     char[] pattern = pat.toCharArray();
69     char[] text    = txt.toCharArray();
70
71     int offset1a = search1(pat, txt);
72     int offset2a = search2(pat, txt);
73     int offset1b = search1(pattern, text);
74     int offset2b = search2(pattern, text);
75
76     // print results
77     StdOut.println("text:    " + txt);

```

```

78
79     // from brute force search method 1a
80     stdout.print("pattern: ");
81     for (int i = 0; i < offset1a; i++)
82         stdout.print(" ");
83     stdout.println(pat);
84
85     // from brute force search method 2a
86     stdout.print("pattern: ");
87     for (int i = 0; i < offset2a; i++)
88         stdout.print(" ");
89     stdout.println(pat);
90
91     // from brute force search method 1b
92     stdout.print("pattern: ");
93     for (int i = 0; i < offset1b; i++)
94         stdout.print(" ");
95     stdout.println(pat);
96
97     // from brute force search method 2b
98     stdout.print("pattern: ");
99     for (int i = 0; i < offset2b; i++)
100         stdout.print(" ");
101     stdout.println(pat);
102 }
103 }

```

Knuth-Morris-Pratt

clever method to always avoid backup

DFA

Deterministic finite state automaton

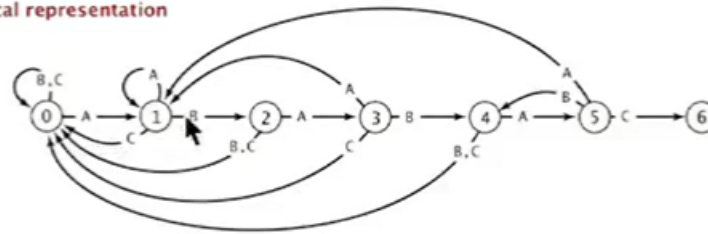
- DFA is abstract string-searching machine
 - finite number of states (including start and halt)
 - exactly one transition for each char in alphabet
 - accept if sequence of transitions leads to halt state

internal representation

j	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[j][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	4

If in state j reading char C :
if j is 6 halt and accept
else move to state $dfa[c][j]$

graphical representation



- Mismatch transition

performance

- KMP substring search accesses no more than $M + N$ chars to search for a pattern of length M in a text of length N
- IKMP constructs $dfa[][]$ and space proportional to RM

Implementation

```

1 public class KMP {
2     private final int R;          // the radix
3     private final int m;          // length of pattern
4     private int[][] dfa;          // the KMP automaton
5
6     public KMP(String pat) {
7         this.R = 256;
8         this.m = pat.length();
9
10        // build DFA from pattern
11        dfa = new int[R][m];
12        dfa[pat.charAt(0)][0] = 1;
13        for (int x = 0, j = 1; j < m; j++) {
14            for (int c = 0; c < R; c++)
15                dfa[c][j] = dfa[c][x];    // Copy mismatch cases.
16            dfa[pat.charAt(j)][j] = j+1;    // Set match case.
17            x = dfa[pat.charAt(j)][x];    // Update restart state.
18        }
19    }
20
21    public KMP(char[] pattern, int R) {
22        this.R = R;
23        this.m = pattern.length;
24
25        // build DFA from pattern
26        int m = pattern.length;
27        dfa = new int[R][m];
28        dfa[pattern[0]][0] = 1;

```

```

29     for (int x = 0, j = 1; j < m; j++) {
30         for (int c = 0; c < R; c++)
31             dfa[c][j] = dfa[c][x];    // Copy mismatch cases.
32         dfa[pattern[j]][j] = j+1;    // Set match case.
33         x = dfa[pattern[j]][x];    // Update restart state.
34     }
35 }
36
37 public int search(String txt) {
38
39     // simulate operation of DFA on text
40     int n = txt.length();
41     int i, j;
42     for (i = 0, j = 0; i < n && j < m; i++) {
43         j = dfa[txt.charAt(i)][j];
44     }
45     if (j == m) return i - m;    // found
46     return n;    // not found
47 }
48
49 public int search(char[] text) {
50
51     // simulate operation of DFA on text
52     int n = text.length;
53     int i, j;
54     for (i = 0, j = 0; i < n && j < m; i++) {
55         j = dfa[text[i]][j];
56     }
57     if (j == m) return i - m;    // found
58     return n;    // not found
59 }
60
61 public static void main(String[] args) {
62     String pat = args[0];
63     String txt = args[1];
64     char[] pattern = pat.toCharArray();
65     char[] text = txt.toCharArray();
66
67     KMP kmp1 = new KMP(pat);
68     int offset1 = kmp1.search(txt);
69
70     KMP kmp2 = new KMP(pattern, 256);
71     int offset2 = kmp2.search(text);
72 }
73 }

```

Boyer-Moore

Basic

Mismatch character heuristic

- scan characters in pattern from right to left
- can skip as many as M text chars when finding one not in the pattern

Mismatch cases

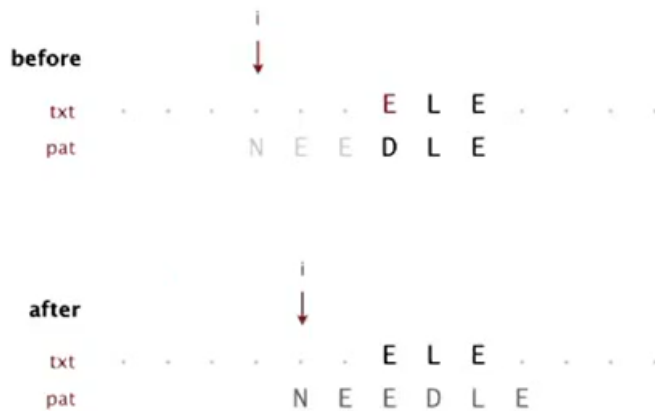
- mismatch character not in pattern



- mismatch character in pattern



Case 2b. Mismatch character in pattern (but heuristic no help).



Performance

Property. Substring search with the Boyer-Moore mismatched character heuristic takes about $\sim N/M$ character compares to search for a pattern of length M in a text of length N . sublinear!

Worst-case. Can be as bad as $\sim MN$.

i skip		0	1	2	3	4	5	6	7	8	9
txt →		B	B	B	B	B	B	B	B	B	B
0	0	A	B	B	B	B	← pat				
1	1		A	B	B	B	B				
2	1			A	B	B	B	B			
3	1				A	B	B	B	B		
4	1					A	B	B	B	B	
5	1						A	B	B	B	B

Implementation

```
1 public class BoyerMoore {
2     private final int R;        // the radix
3     private int[] right;        // the bad-character skip array
4
5     private char[] pattern;     // store the pattern as a character array
6     private String pat;         // or as a string
7
8     public BoyerMoore(String pat) {
9         this.R = 256;
10        this.pat = pat;
11
12        // position of rightmost occurrence of c in the pattern
13        right = new int[R];
14        for (int c = 0; c < R; c++)
15            right[c] = -1;
16        for (int j = 0; j < pat.length(); j++)
17            right[pat.charAt(j)] = j;
18    }
19
20    public BoyerMoore(char[] pattern, int R) {
21        this.R = R;
22        this.pattern = new char[pattern.length];
23        for (int j = 0; j < pattern.length; j++)
24            this.pattern[j] = pattern[j];
25
26        // position of rightmost occurrence of c in the pattern
27        right = new int[R];
28        for (int c = 0; c < R; c++)
29            right[c] = -1;
30        for (int j = 0; j < pattern.length; j++)
31            right[pattern[j]] = j;
32    }
33 }
```

```

34     public int search(String txt) {
35         int m = pat.length();
36         int n = txt.length();
37         int skip;
38         for (int i = 0; i <= n - m; i += skip) {
39             skip = 0;
40             for (int j = m-1; j >= 0; j--) {
41                 if (pat.charAt(j) != txt.charAt(i+j)) {
42                     skip = Math.max(1, j - right[txt.charAt(i+j)]);
43                     break;
44                 }
45             }
46             if (skip == 0) return i;    // found
47         }
48         return n;                    // not found
49     }
50
51     public int search(char[] text) {
52         int m = pattern.length;
53         int n = text.length;
54         int skip;
55         for (int i = 0; i <= n - m; i += skip) {
56             skip = 0;
57             for (int j = m-1; j >= 0; j--) {
58                 if (pattern[j] != text[i+j]) {
59                     skip = Math.max(1, j - right[text[i+j]]);
60                     break;
61                 }
62             }
63             if (skip == 0) return i;    // found
64         }
65         return n;                    // not found
66     }
67
68     public static void main(String[] args) {
69         String pat = args[0];
70         String txt = args[1];
71         char[] pattern = pat.toCharArray();
72         char[] text = txt.toCharArray();
73
74         BoyerMoore boyermoore1 = new BoyerMoore(pat);
75         BoyerMoore boyermoore2 = new BoyerMoore(pattern, 256);
76         int offset1 = boyermoore1.search(txt);
77         int offset2 = boyermoore2.search(text);
78
79         // print results
80         StdOut.println("text:    " + txt);
81
82         StdOut.print("pattern: ");
83         for (int i = 0; i < offset1; i++)
84             StdOut.print(" ");
85         StdOut.println(pat);

```



```

86
87     stdout.print("pattern: ");
88     for (int i = 0; i < offset2; i++)
89         stdout.print(" ");
90     stdout.println(pat);
91 }
92 }

```

Rabin-Karp

Basic

modular hashing

- compute a hash of pattern characters 0 to $M - 1$
- for each i , compute a hash of text characters i to $M + i - 1$
- if pattern hash = text substring hash, check for a match

Two Versions

- Monte Carlo version---return match if hash matches
- Las Vegas version---check for substring match if hash matches; continue search if false collision

Theory. If Q is a sufficiently large random prime (about MN^2), then the probability of a false collision is about $1/N$.

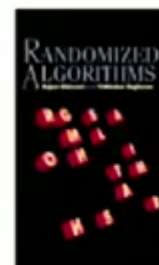
Practice. Choose Q to be a large prime (but not so large to cause overflow). Under reasonable assumptions, probability of a collision is about $1/Q$.

Monte Carlo version.

- Always runs in linear time.
- Extremely likely to return correct answer (but not always!).

Las Vegas version.

- Always returns correct answer.
- Extremely likely to run in linear time (but worst case is MN).



Modular Hash Function

Modular hash function. Using the notation t_i for `txt.charAt(i)`, we wish to compute

$$x_i = t_i R^{M-1} + t_{i+1} R^{M-2} + \dots + t_{i+M-1} R^0 \pmod{Q}$$

Intuition. M -digit, base- R integer, modulo Q .

Horner's method. Linear-time method to evaluate degree- M polynomial.

Implementation

```
1  import java.math.BigInteger;
2  import java.util.Random;
3
4  public class Rabinkarp {
5      private String pat;          // the pattern // needed only for Las Vegas
6      private long pathHash;       // pattern hash value
7      private int m;               // pattern length
8      private long q;              // a large prime, small enough to avoid long
9      // overflow
10     private int R;                // radix
11     private long RM;              // R^(M-1) % Q
12
13     public Rabinkarp(char[] pattern, int R) {
14         this.pat = String.valueOf(pattern);
15         this.R = R;
16         throw new UnsupportedOperationException("Operation not supported yet");
17     }
18
19     public Rabinkarp(String pat) {
20         this.pat = pat;           // save pattern (needed only for Las Vegas)
21         R = 256;
22         m = pat.length();
23         q = longRandomPrime();
24
25         // precompute R^(m-1) % q for use in removing leading digit
26         RM = 1;
27         for (int i = 1; i <= m-1; i++)
28             RM = (R * RM) % q;
29         pathHash = hash(pat, m);
30     }
31
32     // Compute hash for key[0..m-1].
33     private long hash(String key, int m) {
34         long h = 0;
35         for (int j = 0; j < m; j++)
36             h = (R * h + key.charAt(j)) % q;
37         return h;
38     }
39
40     // Las Vegas version: does pat[] match txt[i..i-m+1] ?
```

```

40     private boolean check(String txt, int i) {
41         for (int j = 0; j < m; j++)
42             if (pat.charAt(j) != txt.charAt(i + j))
43                 return false;
44         return true;
45     }
46
47     public int search(String txt) {
48         int n = txt.length();
49         if (n < m) return n;
50         long txtHash = hash(txt, m);
51
52         // check for match at offset 0
53         if ((patHash == txtHash) && check(txt, 0))
54             return 0;
55
56         // check for hash match; if hash match, check for exact match
57         for (int i = m; i < n; i++) {
58             // Remove leading digit, add trailing digit, check for match.
59             txtHash = (txtHash + q - RM*txt.charAt(i-m) % q) % q;
60             txtHash = (txtHash*R + txt.charAt(i)) % q;
61
62             // match
63             int offset = i - m + 1;
64             if ((patHash == txtHash) && check(txt, offset))
65                 return offset;
66         }
67         // no match
68         return n;
69     }
70
71     // a random 31-bit prime
72     private static long longRandomPrime() {
73         BigInteger prime = BigInteger.probablePrime(31, new Random());
74         return prime.longValue();
75     }
76
77     public static void main(String[] args) {
78         String pat = args[0];
79         String txt = args[1];
80
81         Rabinkarp searcher = new Rabinkarp(pat);
82         int offset = searcher.search(txt);
83
84         // print results
85         StdOut.println("text:    " + txt);
86
87         // from brute force search method 1
88         StdOut.print("pattern: ");
89         for (int i = 0; i < offset; i++)
90             StdOut.print(" ");
91     }

```

```

92     stdout.println(pat);
93 }
94 }

```

Comparison

algorithm	version	operation count		backup in input?	correct?	extra space
		guarantee	typical			
brute force	—	MN	$1.1N$	yes	yes	1
Knuth-Morris-Pratt	full DFA (Algorithm 5.6)	$2N$	$1.1N$	no	yes	MR
	mismatch transitions only	$3N$	$1.1N$	no	yes	M
Boyer-Moore	full algorithm	$3N$	N/M	yes	yes	R
	mismatched char heuristic only (Algorithm 5.7)	MN	N/M	yes	yes	R
Rabin-Karp [†]	Monte Carlo (Algorithm 5.8)	$7N$	$7N$	no	yes [†]	1
	Las Vegas	$7N^{\dagger}$	$7N$	yes	yes	1