

Lecture 6

String Matching

Bo Tang @ SUSTech, Fall 2022

2

Our Roadmap

- ◆ String Concepts
- ◆ String Searching Problem
 - ◆ Brute Force Solution
 - ◆ Rabin-Karp
 - ◆ Finite State Automata
 - ◆ Knuth-Morris-Pratt

String Definition

- ◆ String:
 - ◆ Sequence of characters over some alphabet
 - ◆ Binary {0,1}: S1 = "10000101010101001010101"
 - ◆ DNA {ACGT}: S2 = "ACGTACGTACGTTTCA"
 - ◆ English Characters {a...z, A..Z}: S3 = "Hello World"
- ◆ Applications
 - ◆ Word processors
 - ◆ Virus scanning
 - ◆ Text retrieval
 - ◆ Natural language processing
 - ◆ Web search engine

3

String Operators

- ◆ append: append to string
- ◆ assign: assign content to string
- ◆ insert: insert to string
- ◆ erase: erase characters from string
- ◆ replace: replace portion of string
- ◆ swap: swap string values
- ◆ find: find the specific char in the string
- ◆ Give string s="SUSTechCS203", how many sub string it has?

4

Our Roadmap

- ◆ String Concepts
- ◆ String Searching Problem
 - ◆ Brute Force Solution
 - ◆ Rabin-Karp
 - ◆ Finite State Automata
 - ◆ Knuth-Morris-Pratt



5

Why String Searching?

- ◆ **Applications in Computational Biology**
 - ◆ DNA sequence is a long word (or text) over a 4-letter alphabet
 - ◆ GTTTGAGTGGTCAGTCTTTTCGTTTCGACGGAGCCC.....
 - ◆ Find a Specific pattern W
- ◆ **Finding patterns in documents formed using a large alphabet**
 - ◆ Word processing
 - ◆ Web searching
 - ◆ Desktop search (Google, MSN)
- ◆ **Matching strings of bytes containing**
 - ◆ Graphical data
 - ◆ Machine code
- ◆ **grep in unix**
 - ◆ grep searches for lines matching a pattern.

String Searching

Search Text									
a	s	s	u	s	u	s	t	c	s

Search Pattern				
s	u	s	t	c

Successful Search									
a	s	s	u	s	u	s	t	c	s

Parameter

- ◆ n: # of characters in text
- ◆ m: # of characters in pattern
- ◆ Typically, $n \gg m$
 - e.g., $n = 1 \text{ Billion}$, $m = 100$

7

Our Roadmap



- ◆ String Concepts
- ◆ String Searching Problem
 - ◆ Brute Force Solution
 - ◆ Rabin-Karp
 - ◆ Finite State Automata
 - ◆ Knuth-Morris-Pratt

8

Brute Force

- ◆ Brute force
 - ◆ Check for pattern starting at every text position
- ◆ **Algorithm:** BruteForce(T, P):
 1. $n \leftarrow \text{len}(T)$, $m \leftarrow \text{len}(P)$
 2. **for** $i \leftarrow 0$ to $n-m-1$
 3. **for** $j \leftarrow 0$ to $m-1$
 4. **if** $P[j] \neq T[i+j]$ **then**
 5. **break**;
 6. **if** $j = m-1$
 7. pattern occurs with shift i
- ◆ Time complexity?

9

Analysis of Brute Force

- ◆ Analysis of brute force
 - ◆ Running time depends on pattern and text
 - ◆ Can be slow when strings repeat themselves
 - ◆ Worst case: mn comparisons
 - ◆ Too slow when m and n are large

Search Pattern				
a	a	a	a	b

Search Text														
a	a	a	a	a	a	a	a	a	a	a	a	a	a	b
a	a	a	a	b										
	a	a	a	a	b									
		a	a	a	a	b								
			a	a	a	a	b							

.....

10

Can we do better?

- ◆ How to avoid re-computation?
 - ◆ Pre-analyze search pattern
 - ◆ Example: suppose the first 4 chars of pattern are all a's
 - If $t[0..3]$ matches $p[0..3]$ then $t[1..3]$ matches $p[0..2]$
 - No need to check $i=1, j=0,1,2$
 - Saves 3 comparisons
 - ◆ Need better ideas in general

Search Pattern				
a	a	a	a	b

Search Text														
a	a	a	a	a	a	a	a	a	a	a	a	a	a	b
a	a	a	a	b										
	a	a	a	a	b									

11

Our Roadmap



- ◆ String Concepts
- ◆ String Searching Problem
 - ◆ Brute Force Solution
 - ◆ Rabin-Karp
 - ◆ Finite State Automata
 - ◆ Knuth-Morris-Pratt

12

Rabin-Karp Algorithm

- Given search text T and search pattern P as follows:

Pattern					
1	3	5	9		

Search Text											
2	4	6	8	0	1	2	1	3	5	9	7
							1	3	5	9	

- Any idea?

2468	4680	6801	8012	121	1213	2135	1359	3597	5972
							1359		

13

Rabin-Karp Algorithm

- General idea
 - Convert search pattern to a number p
 - Convert search text to an array of numbers $t[0], \dots, t[n-m-1]$
 - Compare p with $t[i]$, for each i in $[0, n-m-1]$
 - if $p = t[i]$, pattern p occurs
- Example
 - $p = 1359$
 - Array t is:

2468	4680	6801	8012	121	1213	2135	1359	3597	5972
------	------	------	------	-----	------	------	------	------	------
 - $t[7] = p \rightarrow T[7,8,9,10] = P[0,1,2,3]$

14

Rabin-Karp Algorithm

- How to convert size-m characters to a number?
 - E.g., the alphabet $\Sigma = \{a, \dots, z, A, \dots, Z\}$
 - Solution: radix-d ($d = |\Sigma|$) Horner's rule
 - $p = P[m-1] + d(P[m-2] + d(P[m-3] + \dots + d(P[1] + dP[0])))$
- When m is large, p may be too large to work
 - Modulo a proper prime number q
 - $p = (P[m-1] + d(P[m-2] + d(P[m-3] + \dots + d(P[1] + dP[0]))) \bmod q$
- Compute $t[0], t[1], \dots, t[n-m-1]$ in time $O(n-m)$
 - Compute $t[i+1]$ by using $t[i]$ in $O(1)$ time
 - $t[i+1] = d(t[i] - d^{m-1}T[i]) + T[i+m]$
 - $t[i+1] = ((t[i] - hT[i]) + T[i+m]) \bmod q$, where $h \equiv d^{m-1} \pmod{q}$
 - $t[0] \rightarrow t[1] \rightarrow t[2] \rightarrow t[3] \rightarrow \dots \rightarrow t[n-m-1]$ in $O(n-m)$

15

Rabin-Karp Algorithm

- Correctness analysis
 - $p \not\equiv t[i] \pmod{q}$ we have $p \neq t[i]$, thus, $P[0..m-1] \neq T[i..i+m-1]$
 - $p \equiv t[i] \pmod{q}$, it does not imply $p = t[i]$ (**spurious hit**)
- Example: search P:

2	3	1	4	1	5	2	6	7	3	9	9	2

Valid match mod 13 Invalid match
- Additional test to check
 - $P[0..m-1] = T[i..i+m-1]$

16

Rabin-Karp Algorithm

- Algorithm:** Rabin-Karp(T, P, d, q):
 - $n \leftarrow \text{len}(T)$, $m \leftarrow \text{len}(P)$
 - $h \leftarrow d^{m-1} \pmod{q}$, $p \leftarrow 0$, $t_0 \leftarrow 0$
 - for** $j \leftarrow 0$ to $m-1$
 - $p \leftarrow (dp + P[j]) \bmod q$,
 - $t_0 \leftarrow (dt_0 + T[j]) \bmod q$,
 - for** $i \leftarrow 0$ to $n-m$
 - if** $p \neq t_i$ then
 - $t_{i+1} \leftarrow (d(t_i - T[i]h) + T[i+m]) \bmod q$
 - else**
 - If $P[0..m-1] = T[i..i+m-1]$
 - pattern occurs with shift i
 - Else**
 - $t_{i+1} \leftarrow (d(t_i - T[i]h) + T[i+m]) \bmod q$

18

Analysis of Rabin-Karp Alg.

- Algorithm:** Rabin-Karp(T, P, d, q):

Overall Cost: $O(mn)$

Our Roadmap

- String Concepts
- String Searching Problem
 - Brute Force Solution
 - Rabin-Karp
 - Finite State Automata
 - Knuth-Morris-Pratt



19

Midterm Exam (tentative)

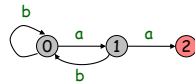
- Time: 12 Nov. 16:30-18:30
- Venue: To be decided
- Scope: Lecture 1 to 6

20

Finite State Automata

- A finite State automaton is defined by:
 - Q , a set of states
 - $q_0 \in Q$, the start state
 - $A \subseteq Q$, the accepting states
 - Σ , the input alphabet
 - δ , the transition function, from $Q \times \Sigma$ to Q

	0	1
a	1	2
b	0	0

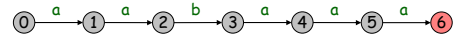


FSA idea for String Matching

- Start in state q_0
- Perform a transition from q_0 to q_1 if next character of $T = P[1]$
- State q_i means first i characters of P match.
- Transition from q_i to q_{i+1} if the next character of $T = P[i+1]$

Search Pattern					
a	a	b	a	a	a

	0	1	2	3	4	5
a	1	2	?	4	5	6
b	?	?	3	?	?	?



- How to fill these ???
 - Reset to q_0 ? Why not?

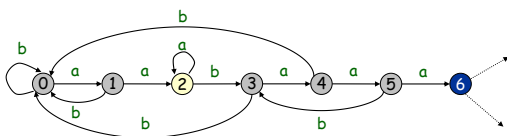
21

22

FSA construction

- FSA construction
 - FSA builds itself
- Example. Build FSA for aabaaabb
 - State 6. $P[0..5] = \text{aabaaa}$
 - assume you know state for $p[1..5] = \text{abaaa}$
 - if next char is b (match): go forward
 - if next char is a (mismatch): go to state for abaaaa
 - update X to state for $p[1..6] = \text{abaaab}$

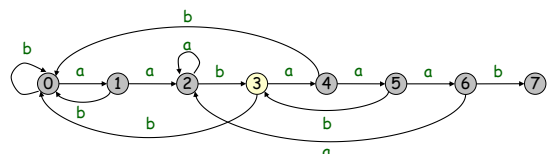
$X = 2$
 $6 + 1 = 7$
 $X + 'a' = 2$
 $X + 'b' = 3$



23

FSA construction

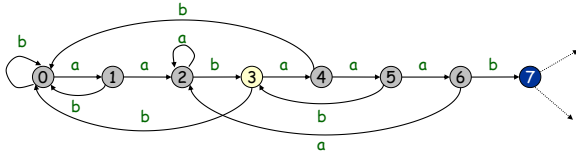
- FSA construction
 - FSA builds itself
- Example. Build FSA for aabaaabb



24

FSA construction

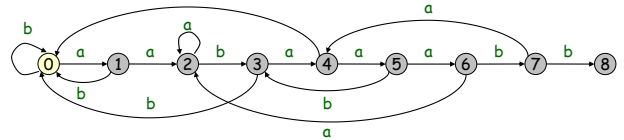
- FSA construction
 - FSA builds itself
- Example. Build FSA for aabaaabb
 - State 7. $p[0..6]=\text{aabaaab}$
 - assume you know state for $p[1..6] = \text{abaaab}$ $X = 3$
 - if next char is b (match): go forward $7 + 1 = 8$
 - if next char is a (mismatch): go to state for abaaaba $X + 'a' = 4$
 - update X to state for $p[1..7] = \text{abaaabb}$ $X + 'b' = 0$



25

FSA construction

- FSA construction
 - FSA builds itself
- Example. Build FSA for aabaaabb



26

FSA construction

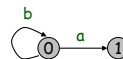
- FSA construction
 - FSA builds itself
- Crucial Insight
 - To compute transitions for state n of FSA, suffices to have:
 - FSA for state 0 to n-1
 - State X that FSA ends up in with input $p[1..n-1]$
 - To compute state X' that FSA ends up in with input $p[1..n]$, it suffices to have
 - FSA for states 0 to n-1
 - State X that FSA ends up in with input $p[1..n-1]$

FSA construction

Search Pattern							
a	a	b	a	a	a	b	b

j	pattern[1..j]	X
---	---------------	---

a
b



27

28

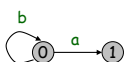
FSA construction

Search Pattern							
a	a	b	a	a	a	b	b



j	pattern[1..j]	X
0		0

0
a 1
b 0



29

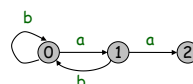
FSA construction

Search Pattern							
a	a	b	a	a	a	b	b



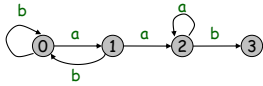
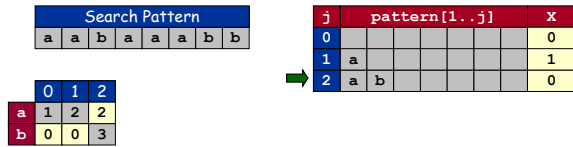
j	pattern[1..j]	X
0		0
1	a	1

0	1
a 1	2
b 0	0



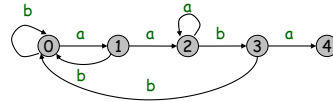
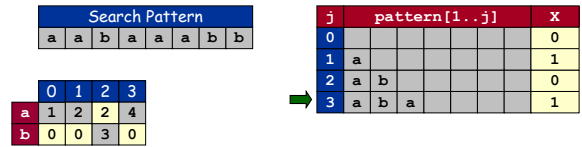
30

FSA construction



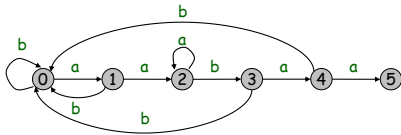
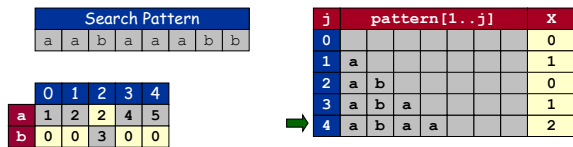
31

FSA construction



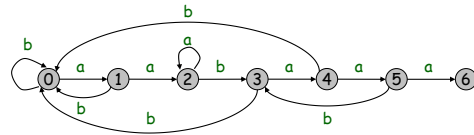
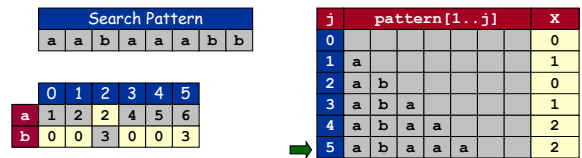
32

FSA construction



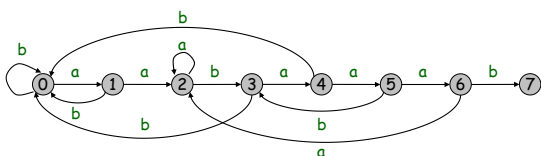
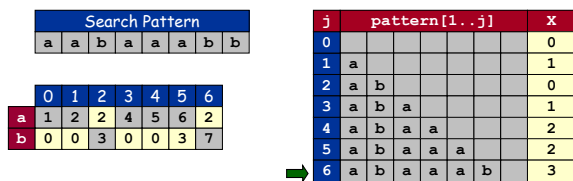
33

FSA construction



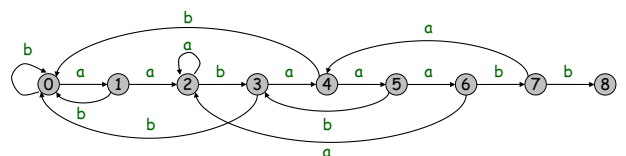
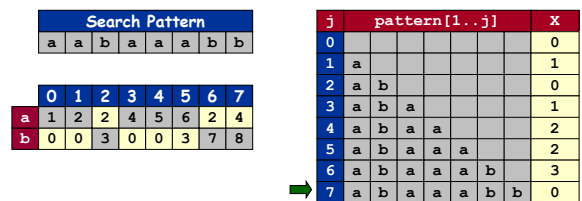
34

FSA construction



35

FSA construction



36

Transition function

Algorithm: Transition(P, Σ):

1. $m \leftarrow \text{len}(P)$
2. $X \leftarrow 0$
3. Initialize $\delta(0, a)$ for each $a \in \Sigma$
4. **for** $j \leftarrow 1$ to $m-1$
5. **for** each character $a \in \Sigma$
6. if $P[j+1] = a$ then // char match
7. $\delta(j, a) \leftarrow j + 1$
8. else // char mismatch
9. $\delta(j, a) \leftarrow \delta(X, a)$
10. $X \leftarrow \delta(X, P[j+1])$
11. **return** δ

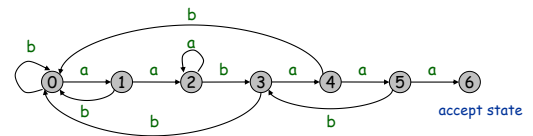
37

Finite State Automata (FSA)

FSA-matching algorithm.

- ◆ Use knowledge of how search pattern repeats itself.
- ➔ ◆ Build FSA from pattern.
- ◆ Run FSA on text.

Search Pattern					
a	a	b	a	a	a



38

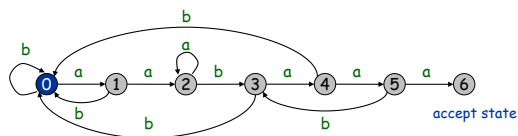
Finite State Automata (FSA)

FSA-matching algorithm.

- ◆ Use knowledge of how search pattern repeats itself.
- ◆ Build FSA from pattern.
- ➔ ◆ Run FSA on text.

Search Pattern					
a	a	b	a	a	a

Search Text											
a	a	a	b	a	a	b	a	a	a	b	



39

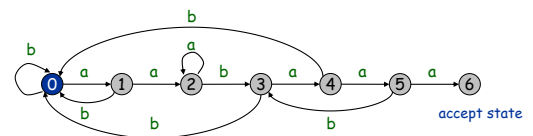
Finite State Automata (FSA)

FSA-matching algorithm

- ◆ Use knowledge of how search pattern repeats itself.
- ◆ Build FSA from pattern.
- ➔ ◆ Run FSA on text.

Search Pattern					
a	a	b	a	a	a

Search Text											
a	a	a	b	a	a	b	a	a	a	b	



40

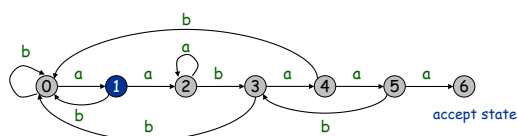
Finite State Automata (FSA)

FSA-matching algorithm

- ◆ Use knowledge of how search pattern repeats itself.
- ◆ Build Finite State Automata (FSA) from pattern.
- ➔ ◆ Run FSA on text.

Search Pattern					
a	a	b	a	a	a

Search Text											
a	a	a	b	a	a	b	a	a	a	b	



41

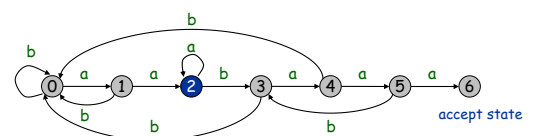
Finite State Automata (FSA)

FSA-matching algorithm.

- ◆ Use knowledge of how search pattern repeats itself.
- ◆ Build FSA from pattern.
- ➔ ◆ Run FSA on text.

Search Pattern					
a	a	b	a	a	a

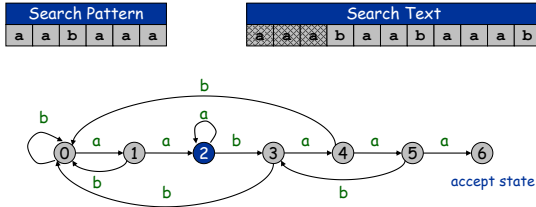
Search Text											
a	a	a	b	a	a	b	a	a	a	b	



42

Finite State Automata (FSA)

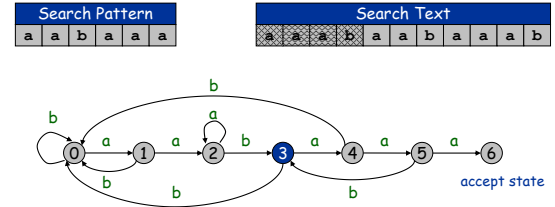
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



43

Finite State Automata (FSA)

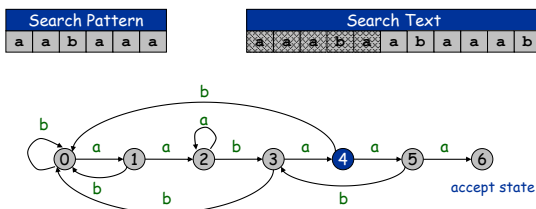
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



44

Finite State Automata (FSA)

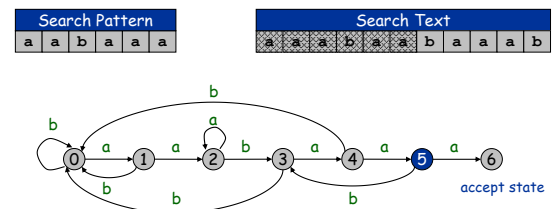
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



45

Finite State Automata (FSA)

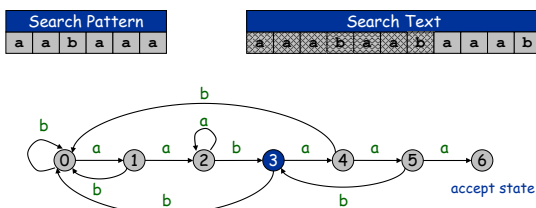
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



46

Finite State Automata (FSA)

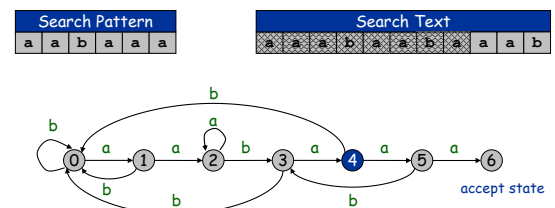
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



47

Finite State Automata (FSA)

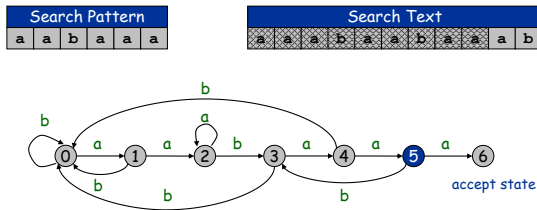
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



48

Finite State Automata (FSA)

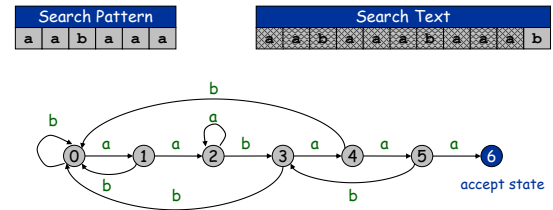
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



49

Finite State Automata (FSA)

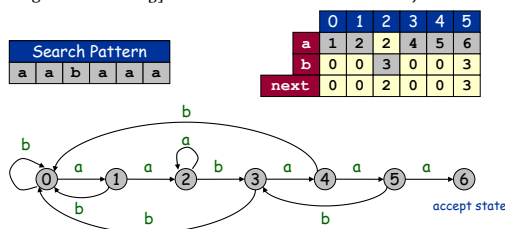
- ◆ FSA-matching algorithm.
 - ◆ Use knowledge of how search pattern repeats itself.
 - ◆ Build FSA from pattern.
- ➡ ◆ Run FSA on text.



50

Finite State Automata (FSA)

- ◆ FSA used in KMP has special property
 - ◆ If match, go to next state
 - ◆ Only need to keep track of where to go upon character mismatch.
 - ◆ go to state next[j] if character mismatches in state j



51

FSA algorithm

- ◆ **Algorithm:** FSA(T, P):
 1. $n \leftarrow \text{len}(T)$, $m \leftarrow \text{len}(P)$
 2. $\delta \leftarrow \text{Transition}(P, \Sigma)$
 3. $q \leftarrow 0$ // q is the state of the FSA.
 4. **for** $i \leftarrow 1$ to n
 5. $q \leftarrow \delta(q, T[i])$
 6. **if** $q = m$
 7. pattern occurs with shift $i - m$

52

Analysis of FSA

- ◆ **Algorithm:** FSA(T, P):

Overall Cost: $O(|\Sigma| m + n)$

53

Our Roadmap

- ◆ String Concepts
- ◆ String Searching Problem
 - ◆ Brute Force Solution
 - ◆ Rabin-Karp
 - ◆ Finite State Automata
 - ◆ Knuth-Morris-Pratt



54

History of KMP

- ◆ Inspired by the theorem of Cook that says $O(m+n)$ algorithm should be possible
- ◆ Discovered in 1976 independently by two groups
- ◆ Knuth-Pratt
- ◆ Morris was hacker trying to build an editor
- ◆ Resolved theoretical and practical problem
 - ◆ Surprise when it was discovered
 - ◆ In hindsight, seems like right algorithm

55

String

- ◆ **String:** "HelloCS203"
- ◆ **Substring:** a substring of s string S is a string S' that occurs in S, e.g., $P[2,...,4] = \text{"ell"}$
- ◆ **Prefix ($P[1,...]$):** a prefix of a string S is a substring of S that occurs at the beginning of S, e.g., $P[1,...,1] = \text{"H"}$ (note that $P[1]=\text{"H"}$), $P[1,...,2] = \text{"He"}$, $P[1,...,5] = \text{"Hello"}$, we denote prefix as: **$P[1,...]$**
- ◆ **Suffix:** a suffix of a string S is a substring of S that occurs at the end of S, e.g., $P[10,...,10] = \text{"3"}$, $P[8,...,10] = \text{"203"}$, $P[6,...,10] = \text{"CS203"}$, we denote suffix as: **$P[...m]$**

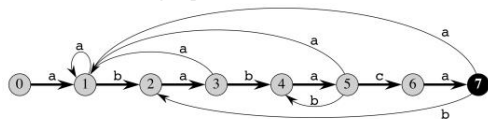
56

Finite State Automata

- ◆ $P = \text{"ababaca"}$
- ◆ Transition function table

State	0	1	2	3	4	5	6	7
a	1	1	3	1	5	1	7	1
b	0	2	0	4	0	4	0	2
c	0	0	0	0	0	6	0	0
P	a	b	a	b	a	c	a	

- ◆ State transition graph



57

Finite State Automata

- ◆ $P = \text{"ababaca"}$ and $T = \text{"abababacaba"}$

i	1	2	3	4	5	6	7	8	9	10	11
T	a	b	a	b	a	b	a	c	a	b	a
1	a	b	a	b	a	c	a				
2			a	b	a	b	a	c	a		
3									a	b	

- ◆ After **failure**: at $i=6$, 'c' was expected, but not found in $T[6]$, FSA transition to state $\delta(5,b)=4$, it means pattern prefix $P[1..4] = \text{"abab"}$ has matched the text suffix $T[2..6] = \text{"abab"}$

	0	1	2	3	4	5	6	7
a	1	1	3	1	5	1	7	1
b	0	2	0	4	0	4	0	2
c	0	0	0	0	0	6	0	0

58

Finite State Automata

- ◆ In general, the FSA is constructed so that the state number tells us how much of a prefix of P has been matched.
- ◆ FSA transition function:
 - 1) Find the longest prefix of P is also a suffix of $T[...i]$, denote as k , i.e., $P[1,...,k] = T[i-k+1,...,i]$
 - 2) Read the next character at " $k+1$ " (i.e., $T[i+1]$), there are two kinds of transitions:
 - ◆ $P[k+1] = T[i+1]$, it is matched, continues.
 - ◆ Otherwise, it is mismatched, go to $\delta(k, T[i+1])$

59

Prefix Function

- ◆ Consider the first step of FSA transition function:
 - ◆ Find the longest prefix of P is also a suffix of $T[...i]$, denote as k , i.e., $P[1,...,k] = T[i-k+1,...,i]$
- ◆ Suppose it is mismatched at " $P[k+1]$ ", it means:
 - ◆ $P[k+1] \neq T[i+1]$ then,
 - ◆ we should find the longest prefix of $P[1,...,k]$ is also a suffix of $T[i-k+1, ..., i]$.
- ◆ **Prefix function (next array in general)**, given $P[1..m]$, the prefix function π for P is $\pi : \{1, 2, ..., m\} \rightarrow \{0, 1, ..., m-1\}$ such that:

$$\pi[i] = \max\{k, k < i \text{ and } P[1..k] = P[i-k+1..i]\}$$

60

Prefix Function

- ◆ **Prefix function**, given P , the prefix function π for P is $\pi : \{1, 2, \dots, m\} \rightarrow \{0, 1, \dots, m-1\}$ such that:

$$\pi[q] = \max\{k, k < q \text{ and } P[1..k] = P[q-k+1..q]\}$$

- ◆ Example: $P = \text{"ababaca"}$

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

61

Compute next array

- ◆ **Algorithm: NextArray(P):**
 1. $m \leftarrow \text{len}(P)$
 2. Let $\pi[1..m]$ be a new array
 3. $\pi[1] = 0, k \leftarrow 0$
 4. **for** $q = 2$ to m
 5. **while** $k > 0$ and $P[k+1] \neq P[q]$
 6. $k \leftarrow \pi[k]$
 7. **if** $P[k+1] = P[q]$
 8. $k \leftarrow k + 1$
 9. $\pi[q] \leftarrow k$
 10. **return** π

62

KMP algorithm

- ◆ **Algorithm: KMP(T, P):**
 1. $n \leftarrow \text{len}(T), m \leftarrow \text{len}(P)$
 2. $\pi \leftarrow \text{NextArray}(P)$
 3. $q \leftarrow 0$
 4. **for** $i = 1$ to n
 5. **while** $q > 0$ and $P[q+1] \neq T[i]$
 6. $q \leftarrow \pi[q]$
 7. **if** $(P[q+1] = T[i])$
 8. $q \leftarrow q + 1$
 9. **if** $q == m$
 10. **print** "Pattern occurs with shift" $i-m$
 11. $q \leftarrow \pi[q]$

63

Analysis of KMP

- ◆ **Algorithm: KMP(T, P):**

Overall Cost: $O(m+n)$

64

Our Roadmap

- ◆ String Concepts
- ◆ String Searching Problem
 - ◆ Brute Force Solution
 - ◆ Rabin-Karp
 - ◆ Finite State Automata
 - ◆ Knuth-Morris-Pratt



Thank You!

65

66