String Matching II

Chapter 32

Lecture Overview

- String matching with Automata
 - DFAs
 - String-matching automata
 - The transition function
- Knuth-Morris-Pratt
 - The prefix function

(Deterministic) Finite Automata

A (deterministic) finite automaton is a mathematical machine that operates on a string. Formally, it is a 5-tuple $M = (Q, q_0, A, \Sigma, \delta)$ where

• Q is a finite set of states,

- $q_0 \in Q$ is the start state,
- $A \subseteq Q$ is the set of accepting states,
- Σ is a finite input alphabet, and
- $\delta: Q \times \Sigma \rightarrow Q$ is the transition function.



The automaton begins in state q_0 and reads the characters of its input string one at a time. If it is in state q and reads input character a, it moves ("makes a transition") to state $\delta(q, a)$. Whenever its current state q is in A, the machine is said to have accepted the string read so far. We are particularly interested in the state that the machine is in when it has read the entire input string; the machine accepts the string if it is in an accepting state, otherwise it rejects the string.

Final State Function

Given a DFA M = (Q, q₀, A, Σ, δ), we can define its final state function $Φ_M$: Σ* \rightarrow Q recursively as:

 $Φ_M(ε) = q_0$ $Φ_M(wa) = δ(Φ_M(w), a)$ for $w ∈ Σ^*$ and a ∈ Σ

 $\Phi_M(w)$ is simply the state that machine M ends up in when its input string is w.



 $\Phi(\epsilon) = 0$ $\Phi(a) = 1$ $\Phi(aa) = 0$ $\Phi(aab) = 1$ $\Phi(aaba) = 0$ $\Phi(babb) = 2$ $\Phi(abab) = 2$ $\Phi(abba) = 1$ $\Phi(aba) = 1$ $\Phi(abab) = 2$

String-matching Automata

There is a string-matching automaton M_P for every pattern P. It will be constructed from the pattern in a preprocessing step and then be used to search the text string T. We want M_P to have the property that $\Phi_M(w) \in A$ iff $P \sqsupset w$. Then we can find the valid shifts by finding the accepting states in Φ_M for text T and its prefixes.



Suffix Function

Given a pattern P, we define the suffix function $\sigma: \Sigma^* \rightarrow \mathbb{Z}_{m+1}$ such that $\sigma(x)$ is the length of the longest prefix of P that is a suffix of x:

 $\sigma(\mathbf{x}) = \max \{ \mathbf{k} : \mathbf{P}_{\mathbf{k}} \sqsupset \mathbf{x} \}$

This is well-defined since $P_0 = \varepsilon$ is a suffix of x.

For pattern aba:

 $\sigma(baa) = 1$ $\sigma(baab) = 2$ $\sigma(baba) = 3$ $\sigma(babb) = 0$

String-Matching Automaton

Given a pattern P[1..m], $M_p = (Q, q_0, A, \Sigma, \delta)$, where:

- Q = {0, 1, ... m}
- $q_0 = 0$
- A = m
- $\delta(q, a) = \sigma(P_q a)$

Lemma. This automaton maintains the invariant $\Phi(T_i) = \sigma(T_i)$

Proof. Suppose $T_{i+1} = T_i a$. Then because $\Phi(T_i) = \sigma(T_i)$, $\Phi(T_{i+1}) = \delta(\Phi(T_i), a) = \sigma(P_{\Phi(Ti)}, a) = \sigma(T_i a) = \sigma(T_{i+1})$. Also, $\Phi(\varepsilon) = 0 = \sigma(\varepsilon)$.

Finite Automaton Matcher FINITE-AUTOMATON-MATCHER(T, δ , m) 1. n = length(T)O(1)2. q = 0O(1) **3. for** i = 1 to n n iterations 4. $q = \delta(q, T[i])$ O(1)5. **if** q = m O(1)6. output i – m O(1)

matching time: O(n)

preprocessing time: ?

Finite Automaton Matcher



(a)





2

2

3

3

— a b a b a b a

4

(c)

5

5

6

4

7

5

8

С

6 7

i

state $\phi(T_i) = 0 - 1$

T[i]

ababaca

9

а

10 11

а

3

b

2

Transition Function COMPUTE-TRANSITION-FUNCTION(P, Σ) 1. m = length(P)O(1)**2. for** q = 0 to m O(m) iterations 3. **for** each character $a \in \Sigma$ $O(|\Sigma|)$ iterations 4. k = min(m+1, q+2)O(1) 5. repeat k = k - 1O(m) iterations 6. **until** $P_k \sqsupset P_a a$ **O(m)** 7. $\delta(q, a) = k$ O(1)8. return δ O(1)

total time: $O(m^3|\Sigma|)$

String Matching with Automata

So string matching with automata takes $O(m^3|\Sigma|)$ preprocessing and O(n) matching time, for a total of $O(n + m^3|\Sigma|)$ time.

The automaton construction was not optimal, and in fact it can be done in optimal $O(m|\Sigma|)$ time. This is optimal because $m|\Sigma|$ is the size of the transition table δ . This gives a total of $O(n + m|\Sigma|)$ time.

The next algorithm, Knuth-Morris-Pratt, achieves O(n) matching time with only O(m) preprocessing time. It doesn't compute the transition function at all.

The Prefix Function

The Knuth-Morris-Pratt (KMP) algorithm uses a prefix function π that allows a transition to be computed in O(1) amortized time. One nice thing about π is that it is independent of the input character and only contains m elements.

Given a pattern P[1..m], the prefix function π for P maps {1..m} to {0..m-1} such that:

 $\pi[q] = \max \{k : k < q \text{ and } P_k \sqsupset P_a\}$

It's the length of the maximum proper prefix of P that is a suffix of P_{q} .

The Prefix Function

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

- for q=3, P₁ = a is the longest proper prefix of P₃ = aba that is also a suffix.
- for q=4, $P_2 = ab$ is the longest proper prefix of $P_4 = abab$ that is also a suffix.
- for q=5, $P_3 = aba$ is the longest proper prefix of $P_5 = ababa$ that is also a suffix.
 - for q = 6, $P_0 = \varepsilon$ is the longest proper prefix of P_6 = ababac that is also a suffix.

Using The Prefix Function



KMP Matcher

```
KMP-MATCHER(T, P)
1. n = length(T)
2. m = length(P)
3. \pi = COMPUTE-PREFIX-FUNCTION(P)
4. q = 0
5. for i = 1 to n
6. while q > 0 and P[q+1] \neq T[i]
7.
     q = \pi [q]
8. if P[q+1] = T[i]
9.
        q = q + 1
10. if q = m
11.
         output i – m
12.
         \mathbf{q} = \pi \left[ \mathbf{q} \right]
```

KMP Matcher

KMP-MATCHER(T, P) 1. n = length(T)2. m = length(P)3. π = COMPUTE-PREFIX-FUNCTION(P) 4. q = 0**5. for** i = 1 to n O(n) iterations 6. while q > 0 and $P[q+1] \neq T[i]$??? 7. $q = \pi [q]$ 8. **if** P[q+1] = T[i]9. q = q + 1**10. if** q = m O(1)11. output i – m 12. $q = \pi [q]$

KMP Matcher

KMP-MATCHER(T, P) 1. n = length(T)2. m = length(P)3. π = COMPUTE-PREFIX-FUNCTION(P) 4. q = 0// invariant: coin on every state < q</pre> **5. for** i = 1 to n O(n) iterations 6. while q > 0 and $P[q+1] \neq T[i]$ O(1) 7. $q = \pi [q]$ // paid for by coin on $\pi [q]$ 8. **if** P[q+1] = T[i] q = q + 1 // add a coin to q first 9. 10. **if** q = m O(1)11. output i – m 12. $q = \pi [q]$

Computing the Prefix Function

COMPUTE-PREFIX-FUNCTION(P) 1. m = length(P)2. $\pi[1] = 0$ 3. k = 0**4. for** q = 2 to m 5. while k > 0 and $P[k+1] \neq P[q]$ 6. $k = \pi[k]$ 7. **if** P[k+1] = P[q]8. k = k + 19. $\pi[q] = k$ b b а a а C а **10. return** π 7 1 2 3 5 6 4 q 3 k 0 1 2 1 0 0

π

0

1

2

1

0

3

0

String Matching

Algorithm	Preprocessing time	Matching time			
Naive	0	O((n-m+1)m)			
Rabin-Karp	$\Theta(m)$	$O((n - m + 1)m)$ $\Theta(n)$ $\Theta(n)$			
Finite automaton	$O(m \Sigma)$				
Knuth-Morris-Pratt	$\Theta(m)$				