

String Matching

- text editor
- dna sequencing $3 \cdot 10^9$
- web search

T: a c a b a a b a c c a a b a c c

P: a b a a

$$\Sigma = \{a, b, c\}$$

$$\Sigma^* \text{ all strings over } \Sigma$$

$w \sqsubset x$ = w is a prefix of x

$w \sqsupset x$ = w is a suffix of x

$ab \sqsubset abcb$

$ba \sqsupset caba$

$$P[j:k] = P[j] P[j+1] \dots P[k]$$

$$P[:k] = P[1:k] \quad P[:0] = \epsilon$$

Lemng 32.1

Naive string matcher (t, p, n, m)

$n-m$ $\left\{ \begin{array}{l} \text{for } s = 0 \text{ to } n-m \\ \text{if } p[1:m] = t[s+1:s+m] \\ \text{print } s \end{array} \right.$

$$\Theta((n-m+1) \cdot m) = \Theta(n \cdot m)$$

Rabin-Karp

$$\Sigma = \{a, b, c, d, e, f, g, h, i, j\}$$

0 1 2 3 4 5 6 7 8 9

T: c c a d b e b f c b

2 2 0 3 1 4 1 4 2 1

└──────────┘

P: d b e b f

3 1 4 1 5

3 1 4 1 5

$$\Theta(m) + \Theta(n - m + 1)$$

Trick:

$$\begin{array}{r} 22031 \cdot 10^4 = \\ - 20000 \\ \hline 20 \times 2031 \\ \hline 20310 \\ + \quad 4 \\ \hline 20314 \end{array}$$

4 ops

$$10^{m-1}$$

Heuristic

$$t_{s+1} = (t_s - T[s+1]) \times d^{m-1} + T[s+m+1] \pmod q$$

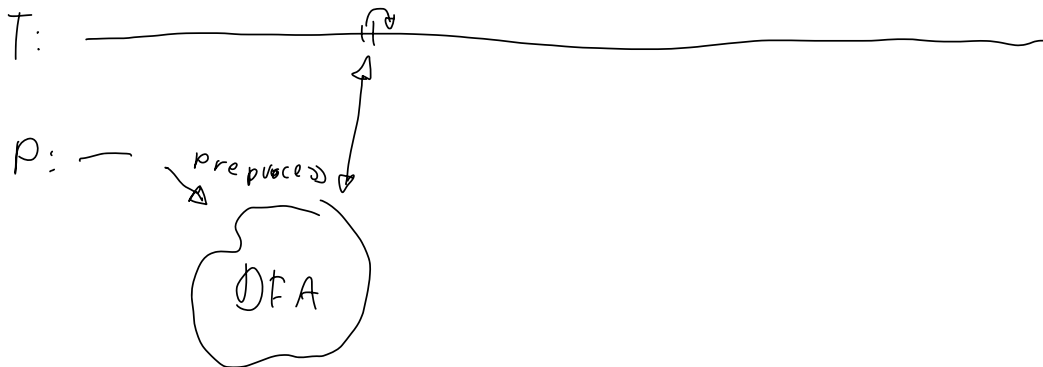
$d = |\Sigma|$

$$t_s \xrightarrow{d} t_{s+1}$$

$T[s+1: s+m]$

choose prime q as large as possible
so that $q \cdot d$ fits in 1 word

String Matching with a finite automata



DEF A $M = (Q, q_0, A, \Sigma, \delta)$

Q finite set of states

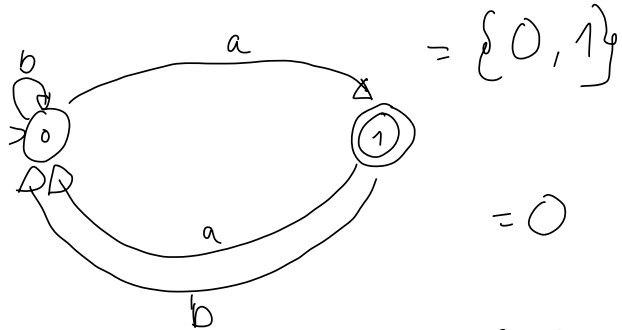
$q_0 \in Q$ start state

$A \subseteq Q$ accepting states

Σ finite alphabet

$\delta: Q \times \Sigma \rightarrow Q$ transition function

δ	Σ	
	a	b
0	1	0
1	0	0



$= \{0, 1\}$

$= 0$

$\{1\}$

$\Sigma = \{a, b\}$

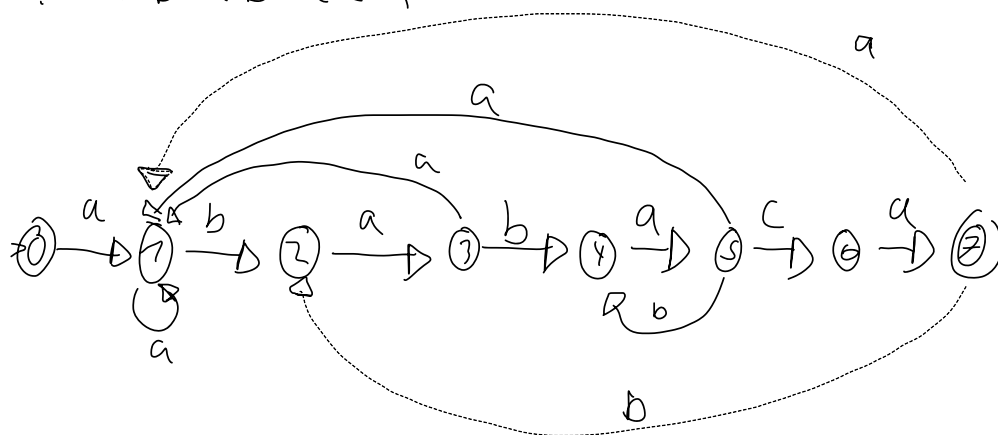
Suffix function

$$\sigma : \Sigma^* \rightarrow \{0, 1, \dots, m\}$$

$$\sigma(x) = \max \{ k \mid p[:k] \sqsupseteq x \}$$

$$\delta(q, a) = \sigma(p[:q] a)$$

P: a b a b a c a



missing arrows
are
going to

zero
all states
need an
arrow

$$\nabla(x) = \max \{k \mid P[:k] \supset x\}$$

$$\delta(q, a) = \nabla(P[:q]a)$$

$$\delta(3, b) = \nabla(\underline{a b a} b) = 4$$

$$\delta(3, a) = \nabla(a b a \underline{a}) = 1$$

$$\delta(3, c) = \nabla(a b a \underline{c}) = 0$$

$$\delta(5, b) = \nabla(a b a b a \underline{b}) = 4$$

$$P[:k] = \epsilon$$

T: bacbabababcbab

P: |abab|ca
 |ab|ab|ca
 |a|ba|...

$$\pi[q] = \max \{k \mid k < q \text{ and } P[:k] = P[:q]\}$$

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

Kmp - Matcher (T, P, n, m) put code

Compute - Prefix - Function (P, m) put code

DFA vs. kmp

$\rightarrow 0 \dots 0 \xrightarrow{c} 0 \dots 0 \xrightarrow{b} 0 \dots 0 \xrightarrow{a} a \dots 0$