

Strings

In Java, (and other languages - Python, Scala?, C#) strings are immutable.

Algorithms such as:

reverse (String s)

X $s \rightarrow \text{Prefix} + c$

 return c + reverse (Prefix)

} $RT = O(n^2)$

↳ $O(n)$ algorithm: $\left[\begin{array}{l} \text{StringBuilder class} \\ \text{toCharArray()} - \text{convert to} \\ \text{character array.} \end{array} \right.$

In general:

do not write: if (s.equals("Find")) ...

(see LP4 Driver).

- ok only for during testing phase
- should not be used in production code.

Avoid: HashMap < string, x > .

String Matching Input: Text $T[1..n]$, Pattern $P[1..m]$

Output: All occurrences of P in T - valid shifts

shifts s_1, s_2, \dots, s_k : $T[s_i+1 \dots s_i+m] = P[1..m]$.

Ex: $T = \boxed{a} \boxed{b} \boxed{a} \boxed{b} \boxed{a} ab$

$P = aba$: shifts = 0, 2

Naive Algorithm

for $s \leftarrow 0$ to $n-m$ do

if $P.equals(T.substring(s, s+m))$

$P[1..m] = T[s+1..s+m]$

s is a valid shift

then } $\left. \begin{array}{l} RT = \\ O(mn) \end{array} \right\}$

Better algorithms for string matching

1. Rabin-Karp algorithm: use ideas from hashing to reduce RT.

RT = $O(mn)$
worst case,

$O(m+n)$
most of
the times

$$\left. \begin{aligned} h(P[1..m]) &= x \\ h(T[s+1..s+m]) &= y \end{aligned} \right\} \begin{aligned} &\text{if } x \neq y \text{ then} \\ &\quad \text{it is not a valid shift} \\ &\text{else} \end{aligned}$$

Use naive algorithm
on $P[1..m], T[s+1..s+m]$

From $h(T[s+1..s+m])$,

calculate $h(T[s+2..s+m+1])$ in $O(1)$ time.

h. Consider $T[s+1..s+m]$ as a polynomial in some base b

$$s \leftarrow 0 \quad a \leftarrow h(P), \quad t \leftarrow h(T[1..m]) \pmod p \quad (\text{for some random, large prime } p)$$

$$\text{mult} \leftarrow b^{m-1} \pmod p$$

while $s < n-m$ do

if $a = t$ then Naive matcher ($P[1..m], T[s+1..s+m]$)

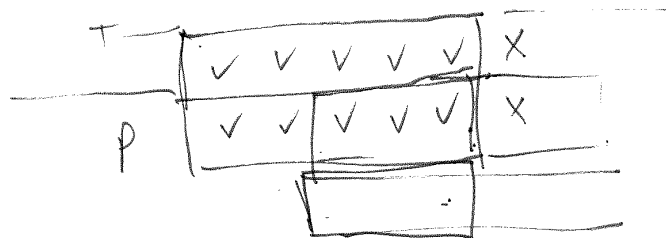
$$t \leftarrow (t - T[s] \times \text{mult}) \times b + T[s+m+1] \pmod p$$

2. Finite Automata matcher - next class. RT = $O(m+n)$.

$$\text{space: } O(m \cdot \Sigma)$$

3. Knuth-Morris-Pratt (KMP) Algorithm

Σ - size of alphabet.



Prefix function π : $\pi[q]$ = Length of
Longest prefix of $P[1..q]$
that is also a suffix of $P[1..q]$

KMP($T[1..n], P[1..m], \pi[1..m]$) // Given prefix function π

$q \leftarrow 0$ // Length of current match.

for $i \leftarrow 1$ to n do // Looking at $T[i]$.

while $q > 0$ and $P[q+1] \neq T[i]$ do

$$q \leftarrow \pi[q]$$

if $P[q+1] = T[i]$ then $q++$

if $q = m$ then $s = i - m$ is a valid shift