

String Matching

①

→ Why it's imp

↳ Word processors ↳ Crawlers

↳ Search Engine

↳ medical field (DNA matching) / (takes hours/days)

↳ Database Search

↳ Compilers (Syntax/Parsing)

2 Types : Exact String Matching / Approximate String matching

① → Naive String matching (Brute force)

② → Rabin Kert (RK Algo).
 ↓
 defined number of errors allowed

③ → Finite Automata
 ↓
 misspelled / misspelled

④ → KMP → Knuth - Morris - Pratt Algo

⇒ Naive String Matching (Brute force).

Text length = N

Pattern length = m .

for ($i = 0$ to $n - m$) ++)

if ($P(0 + m) = T(i + 1 + n)$)

found,

e.g

$T =$ a b c d a b b c x y z p q

$P =$ a b b c.

$((n-m+1) m)$ time

No preprocessing

Notations

Σ^* = set of all finite strings from alphabet Σ .

ϵ = empty string

$|x|$ = length of string

xy = string ~~concat~~ ~~concat~~ concat length
 $|x| + |y|$

w is prefix of x $w \in x$ if $x = wy$

where $y \in \Sigma^*$

w is suffix of x if $x = yw$ where $y \in \Sigma^*$

ϵ is both suffix & prefix of any string.

(2)

ab \neq ab cca
cca \neq ab cca

if $(m = n/2) \Rightarrow O(n^2)$

Rabin Karp Algo. ^{preparation = $O(m)$}
Basic Idea

\rightarrow a a a a b (q) ^{not}
Algo

\rightarrow a a b. (p)

We change our text in numeric value.

e.g. a = 97 (ascii value).

So assume

$a = 1$
 $b = 2$
 $c = 3$
 $d = 4$
 $e = 5$
 $f = 6$
 $g = 7$
 $h = 8$
 $i = 9$
 $j = 10$

So a a b

$$m = \underbrace{1 + 1 + 2}_{\text{hash code}} = \underbrace{4}_1$$

function to get this value
is hash function

$$|m| = 3$$

Now

$$N = \begin{matrix} & 1 & 1 & 1 & 1 & 1 & 2 \\ & a & a & a & a & a & b \end{matrix}$$

So these character $\neq |m|$

not matching, slide

$$\begin{matrix} a & a & a & a & a & b \\ & \underbrace{\hspace{1cm}} \end{matrix}$$

not matching
again slide

$$\begin{matrix} a & a & a & a & a & b \\ & \underbrace{\hspace{1cm}} \end{matrix}$$

slide = rolling hash function

~~a a a~~

Q Q Q 'a a b

3

as $4 = 4$
 match $4 = 4$ & $Qab = aab$
 So matched

How we roll

e.g. a b c d ~~a b~~ c e
 $1 + 2 + 3 = 6$

for shifting

$$(6 - 1) + 4 = 9$$

$$b \cdot c \cdot e = 10$$

example.

when hash matched, check actual pattern

$$O(n - m + 1) \text{ time (avg)}$$

⇒ Now issue of approach.

⇒ worst case for this

c c a c c a a e d b a

$$3 + 3 + 1 = 7$$

d b a

$$1 + 2 + 1 = 7$$

→ code match.

→ compare alphabet

Next

again 7.

→ compare alphabet

⇒ Every time hash code is same.

⇒ Simple hash function lead to collisions (Spurious hits)

So max time is $(m \cdot n)$

→ So we need a strong hash function 4
by Robin Karp.

c c a c c a q e d b a

a b a

$$4 \times 10^2 + 2 \times 10^1 + 1 \times 10^0$$

$$P[1] \times \underbrace{10^{m-1}}_{\substack{\text{we can take} \\ \text{any value}}} + P[2] \times 10^{m-2} + P[3] \times 10^{m-3}$$

base value = no. of ch taken if we take all english alphabets = 26 = base

if all lower, upper, special = 127

So value of $P = 421$.

⇒ Now we need match

$$c c a = 331$$

So not a superious hit.
Now rolling function? Here

$$\left[(3 \times 10^2 + 3 \times 10^1 + 1 \times 10^0) - 3 \times 10^2 \right] \times 10 + 3 \times 10^0$$

⇒ Robin Finger Print function.

where base = input considered.

$$O(n-m+1),$$

This reduces the chances of spurious hit, but still might release it

So worst (mm)

If we want to avoid overflow (due to big values), we will mod these values.

So what should be mod value depends on data type if $m1 \% 2^31$ could be any value.

↳ will get superior hits
(Total extended ASCII = 256)

Algo

(5)

$n \rightarrow \text{length } T$
 $m \rightarrow \text{length } P$
 $h \rightarrow d^{m-1} \text{ mod } q$
 $p = 0$
 $t_0 = 0$
 $q = \text{prime no.}$
 $\rightarrow \text{Any random prime no.}$
 $\rightarrow \text{usually } 256 \Rightarrow (\text{base})$

for $i = 1$ to m // preprocessing

$p \rightarrow [dp + P[i]] \text{ mod } q$
 $t_s \rightarrow [d \cdot t_r + T[i]] \text{ mod } q$

for $s = 0$ to $n - m$

if $p = T_0$

then

check $P[1 \text{ to } m] = T(s+1 \text{ to } s+m)$

print pattern

if $s < n - m$

$t_{s+1} = (d (t_s - T[s+1])^h + T[s+m+1]) \% d$

if $(t < 0)$
 $t = (t + q)$