

Hash Functions and String Matching

Sriram Sankaranarayanan

Data Structures and Algorithms

Rabin - Karp

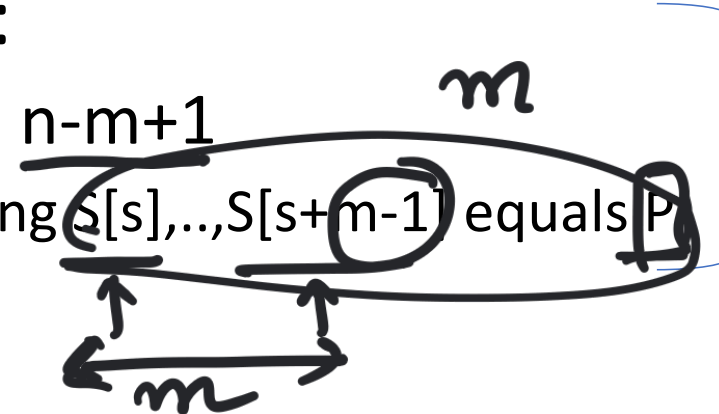
Very important problem with lots of applications!

String Matching

- Does a given pattern P of size m occur in a string S of size n?
- Example:
 - P = "GATTACA" of size m = 7
 - S = "GTAGATAGATTTAATTACGATTACATGATGTTGATTAGGATGATTACATATATGAATA
ATAGCGCCGATATAGAT" $\uparrow \dots$
 - Answer: P occurs in S at the position shown in red.

- Simple algorithm:

For each $s = 1, \dots, n-m+1$
• Check if substring $S[s], \dots, S[s+m-1]$ equals P



$$\begin{array}{c} 1000 \quad 10^6 \\ \swarrow \quad \searrow \\ \Theta(\underbrace{m}_{1000} \times (n - m + 1)) \\ \approx 10^9 \text{ ops} \end{array}$$

Speeding up Matching Using Hash-Functions

Idea: Hash the pattern P using hash function: $h(P)$ *← hash of this pattern*

compute pattern hash: $r = h(P)$

for each $s = 1, \dots, n-m+1$

compute: $q = h(S[s] \dots S[s+m-1])$ *← $\Theta(m)$ filtering*

if $r = q$:
 compare $S[s] \dots S[s+m-1]$ with P.

Diagram: $S[1] \dots S[s] \dots S[s+m-1]$ with a double-headed arrow under $S[s] \dots S[s+m-1]$ labeled m .

Running time: $\Theta(\underline{m} \times (n - m + 1))$

Rolling Hash Functions

$s_1 \dots s_m$ ~~$s_2 \dots s_m$~~ s_{m+1}

p prime number

$$h(s) = (p^{m-1}s_m + p^{m-2}s_{m-1} + \dots + p^1s_2 + s_1) \bmod M$$

$$h(s') = (p^{m-1}s_{m-1} + p^{m-2}s_{m-2} + \dots + p^1s_1 + \hat{s}) \bmod M$$

$$= (\bar{p} \times h(s) - p^{m-1}s_m + \hat{s}) \bmod M$$

Updated

Precompute $p^{m-1} \bmod M$ and perform two multiplications + one subtraction + one addition

Using Rolling Hash Function

$S[1] \dots S[m-1]$
←→

compute pattern hash: $r = h(P)$

for each $s = 1, \dots, n-m+1$

compute: $q = h(S[s] \dots S[s+m-1])$

← $\Theta(1)$ use rolling hash function

if $r = q$:

← compare $S[s] \dots S[s+m-1]$ with P .
collision

Running time: Worst case
continues to be

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

$m = 10^3$
 $n = 10^6$
 10^9

Assuming low probability of
hash collision: we can
improve the running time to

$\times 1000$

$\Theta(m + n)$

10^6

Problem # 2 : Check if two strings have a common substring of size m .

- Inputs: Two strings $S1$ of size $n1$, and $S2$ of size $n2$.
- Output: True if $S1$ and $S2$ have a common substring of size m , FALSE otherwise.
- Example:
 - $S1 = \text{"GATATATACAGACAATAGATAGACACACGTAGGTGCACAGT"}$
 - $S2 = \text{"AGGATTAGGTGGAACCCAGAGAGTTTAGGACCAGATTAGAT"}$
 - $m = 5$
 - Answer: True

Simple Algorithm

- for $i = 1$ to $n_1 - m + 1$
 - $P = S1[i] \dots S1[i+m-1]$ ← m . hash
 - Use previous problem to search for pattern in S2
 - If pattern P found, then return True.
 - Else, continue.
 - Return False

Rabin Karp

Assuming good hash function:

$$\Theta(n_1 \times (m + n_2))$$

Worst Case:

$$\Theta(n_1 \times ((n_2 - m + 1) m))$$

Idea: Use a hash table and hash functions

However, we will need extra space $\Theta(n_1)$ for the hashtable.

Improved Algorithm

For $i = 1$ to $n_1 - m + 1$

- Compute rolling hash $h_i = h(S1[i], \dots, S1[i + m - 1])$ $\Theta(n_1)$

 Insert $\{(h_1, 1), (h_2, 2), \dots, (h_{n_1 - m + 1}, n_1 - m + 1)\}$ into perfect hash table H . $\Theta(n_1 - m + 1)$

For $j = 1$ to $n_2 - m + 1$

- Compute rolling hash $r_j = h(S2[j] \dots S2[j + m - 1])$ $\Theta(n_2)$
- Is key r_j in hashtable H ?
- If yes, let k be the associated value with the key r_j
 - Compare $S1[k] \dots, S1[k + m - 1]$ with $S2[j] \dots S2[j + m - 1]$ $\Theta(m)$

$2 \times 10^6 \neq 10^3$

If there are no spurious collisions: $\Theta(n_1 + n_2 + m)$ Otherwise: $\Theta(n_1 + (n_2 - m + 1) \times m) \sim 10^9$