

①

PLAN

1. ROLLING HASH (RECAP)

2. PROBABILISTIC TOOLS

- UNION BOUND
- MARKOV'S INEQUALITY
- LINEARITY OF EXPECTATION
- AMPLIFICATION

3. BIRTHDAY PARADOX AND THE REVERSE PROBLEM

4. DNA PROBLEM

5. PATTERN MATCHING

ROLLING HASH

$s[1 \dots n]$

$$h(s[i:i+l]) = (s[i] \cdot b^{l-1} + s[i+1] \cdot b^{l-2} + \dots + s[i+l-1]) \pmod p$$

$$h(s[i+1:i+l+1]) = \left(\left(\dots \right) \cdot b + s[i+l] \right) \pmod p$$

Can compute $\left(b^{l-1} \pmod p \right)$ in $O(l)$ time

Then $h(s[i+1:i+l+1])$ from $h(s[i:i+l])$ in $O(1)$ time

②

TOOLS

UNION BOUND

B_1, \dots, B_k - bad events

$$\Pr[\text{any bad event}] \leq \sum_{i=1}^k \Pr[B_i]$$

MARKOV'S INEQUALITY

$$X \geq 0$$

↑
random variable

a - any positive real

$$E[X] \geq a \cdot \Pr[X \geq a]$$

that is

$$\Pr[X \geq a] \leq \frac{E[X]}{a}$$

Example:

expected running time: 10s

probability runs longer than 100s: $\leq \frac{10}{100} = \frac{1}{10}$

LINEARITY OF EXPECTATION

$$E[aX + bY] = a \cdot E[X] + b \cdot E[Y]$$

Example:

run 5 times procedure A, expected time of A: 3s

run 2 times B : 2s

total expected time: $5 \cdot 3s + 2 \cdot 2s = 19s$

③

AMPLIFICATION

A = ALGORITHM THAT OUTPUTS YES/NO
CORRECT WITH PROBABILITY $\geq 2/3$

RUN k TIMES AND RETURN MAJORITY

CAN SHOW: PROBABILITY OF ERROR $2^{-\Omega(k)}$

BIRTHDAY PARADOX

WE'RE SELECTING RANDOM NUMBERS FROM $\{1, \dots, N\}$
WITH REPLACEMENT

HOW MANY SAMPLES TO DRAW THE SAME
NUMBER TWICE WITH CONSTANT PROBABILITY?

ANSWER: $\Theta(\sqrt{N})$

WE'LL ONLY SHOW LOWER BOUND

SAY, WE PICK k NUMBERS: a_1, a_2, \dots, a_k

$$\Pr[a_i = a_j] = 1/N$$

$\swarrow \searrow$
 $i \neq j$

$$\Pr[\text{collision}] \leq \sum_{i \neq j} \Pr[a_i = a_j] = \frac{k(k-1)}{2} \cdot \frac{1}{N}$$

UNION BOUND

④ EXAMPLE: WANT

$$\Pr[\text{collision}] \gg \frac{1}{2}$$

THEN

$$\frac{1}{2} \leq \frac{k(k-1)}{2} \cdot \frac{1}{N}$$

$$N \leq k(k-1)$$

SO

$$k^2 \gg N$$

$$k \gg \sqrt{N}$$

EXAMPLE 2: n elements hashed into $\{1 \dots 100n^2\}$
uniformly

$$\Pr[\text{collision}] \leq \frac{n(n-1)}{2} \cdot \frac{1}{100n^2} = \frac{1}{200}$$

DNA PROBLEM

(uniform hashing
assumption
but we use rolling
hash)

ATTEMPT 1: algorithm from the last lecture: create $N = \Theta(n^2)$ -size array
expected number of collisions in each iteration =

$$\leq \cancel{\binom{n}{2}} \cdot \frac{1}{N} = O(1)$$

rolling hash: can compute all hash values in $O(n)$ time

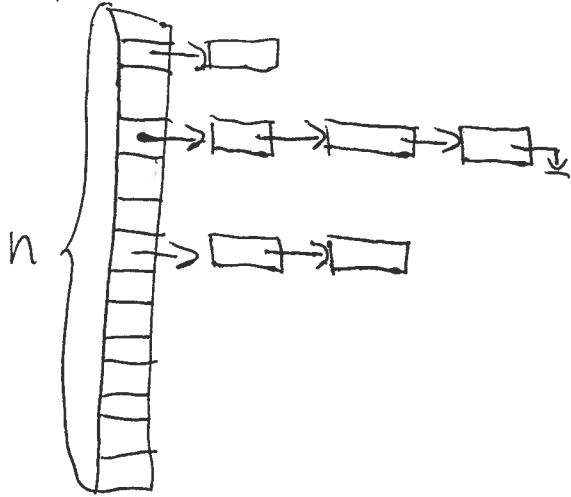
expected ~~time~~ time for comparing ~~the~~ strings: $O(1) \cdot O(n) = O(n)$

number of iterations: $O(\log n)$

total expected time: $O(n \log n) + O(n^2)$
initialization \uparrow

5

ATTEMPT 2: HASH INTO A TABLE OF SIZE $\Theta(n)$ BUT COMPARE FIRST SIGNATURES FROM $\{1, \dots, n^2\}$



WHEN INSERTING NEW SUBSTRING; EXPECTED NUMBER OF SIGNATURE COMPARISONS = $O(1)$

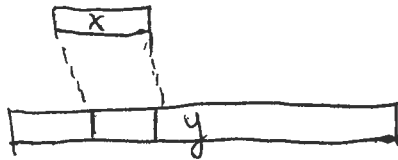
EXPECTED TOTAL NUMBER OF STRING COMPARISONS STILL $O(1)$

EXPECTED RUNNING TIME: $O(n) + O(n) \cdot O(1) + O(1) \cdot O(n) = O(n)$
PER ITERATION \uparrow ROLLING HASH

ALL $O(\log n)$ ITERATIONS: $O(n \log n)$

PATTERN MATCHING

FIND AN OCCURENCE OF PATTERN x IN TEXT ~~in~~ y



NAIVE SOLUTION: $O(|x| \cdot |y|)$
WANT $O(|y|)$

(6)

1. Compute rolling hash for all subwords of y of length $|x|$
2. Run naive check for subwords y' such that $h(y') = h(x)$ until you find occurrence of x

If probability that $h(y') = h(x)$ is $O(1/|x|)$ for any y' , then expected running time:

$$\cancel{O(|y|)} \quad O(|y|) + O(|y|) \cdot O(1/|x|) \cdot O(|x|) = O(|y|)$$

↑
rolling hash

This algorithm known as Rabin-Karp algorithm

Our rolling hash invented for this algorithm