

ALG 4.1
*Randomized Pattern
Matching:*

Reading Selection:
CLR: Chapter 12

Handout: R. Karp and M. Rabin,
"Efficient Randomized
Pattern-Matching"

Generalized Pattern Matching

input index set R
for each $r \in R$ strings $X(r), Y(r) \in \{0,1\}^m$
problem find $r \in R$ s.t. $X(r) = Y(r)$

Examples:

(1) *ID String Pattern Matching*

input pattern $X = X_1 \dots X_m \in \{0,1\}^m$

text $Y = Y_1 \dots Y_n \in \{0,1\}^n$

index set $R = \{1,2,\dots, n-m+1\}$

$\forall r \in R$ $X(r) = X$
 $Y(r) = Y_r Y_{r+1} \dots Y_{r+(m-1)}$

(2) 2D Array Matching

input pattern

$s \times s$ binary array $X = (X_{ij})$, $m = s^2$

text $b \times b$ binary array $Y = (Y_{ij})$, $n = b^2$

index set $R = \{ \langle i,j \rangle \mid s \leq i,j \leq b \}$

$X(\langle i,j \rangle)$ = string of rows of X

$Y(\langle i,j \rangle)$ = string of rows in $s \times s$ block of Y
with (i,j) in lower right position.

(note Karp & Pratt reverse n,m)

Pattern Matching by Fingerprinting

S is a finite set

$\forall p \in S$ $\Phi_p(\cdot)$ is function $\{0,1\}^m \rightarrow$ small range D_p

$\Phi_p(X)$ is "*fingerprint*" for string X

idea compare $X(r) = Y(r)$ only if
fingerprints agree: $\Phi_p(X(r)) = \Phi_p(Y(r))$

Algorithm

$p \leftarrow$ random element of S

for each $r \in R$ in order *do*

begin

compute $a_p(r) = \Phi_p(X(r))$

compute $b_p(r) = \Phi_p(Y(r))$

if $a_p(r) = b_p(r)$ *then*

if $X(r) = Y(r)$ *then* output "Match at r "

end

fingerprint fn $\Phi_p : \{0,1\}^m \rightarrow D_p$

Requirements

- (1) small domain D_p
- (2) small probability of *false match*
 $\Phi_p(X(r)) = \Phi_p(Y(r))$ but $X(r) \neq Y(r)$
- (3) fingerprints $\Phi_p(X(r))$, $\Phi_p(Y(r))$ are *easily updatable* from previous r

Examples of fingerprints:

- (A) integer modular fns.
- (B) unimodular matrices
- (C) irreducible polynomial modular fns.

represent *binary string* $X = X_1 \dots X_m$

by integer $H(x) = \sum_{i=1}^m X_i 2^{m-i}$

modular fingerprint $\Phi_p(x) = \text{res}(H(x), p)$

modular fingerprint $\Phi_p(X) = \text{res}(H(X), p)$
 $= p \cdot \lfloor \frac{H(X)}{p} \rfloor - H(X)$

note $\Phi_p(x) \equiv H(X) \pmod{p}$

Define $S = \{p \mid p \text{ is prime and } p \leq M\}$

where M is a (suf. large) integer

idea choose *random* $p \in S$

\Rightarrow must prove $\Phi_p(X) = \text{res}(H(X), p)$ is

good fingerprint

Facts about Prime Numbers

let $\Pi(k) = \text{number of primes } \leq k$

FACT 1 If $k \geq 29$ and $a \leq 2^k$, then
 a has $\leq \Pi(k)$ distinct prime divisors

proof follows from Rosser & Schoenfeld bound:

$$\prod_{\substack{p \text{ prime} \\ p \leq k}} p > e^{k - 2.05282k^{\frac{1}{2}}}$$

FACT 2 for all $k \geq 17$

$$\frac{k}{\ln k} \leq \Pi(k) \leq 1.25506 \frac{k}{\ln k}$$

(prime number theorem)

Suppose

Randomized Pattern Match

Algorithm is executed, with

fingerprint $\Phi_p(X) = \text{res}(H(x), p)$

with set $S = \{p \mid p \text{ prime} \leq M\}$

where $M = mn^2$ and $mn \geq 29$

Theorem

for each $r \in R$, probability of
false match $\Phi_p(X(r)) = \Phi_p(Y(r))$ but $X(r) \neq Y(r)$

$$\text{is } \leq \frac{2.511}{n}$$

proof

A false match occurs only if $\exists r \in R$

$X(r) \neq Y(r)$ and $P \mid (H(x(r)) - H(Y(r)))$

iff $p \mid L$ where $L = \prod_{X(r) \neq Y(r)} |H(X(r)) - H(Y(r))|$

But $L \leq 2^{mn}$ so by Fact 1,
 L has at most $\Pi(mn)$ prime divisors

Since p is chosen at *random* from $\Pi(M)$ primes,
and only $\Pi(mn)$ give false matches,

$$\text{Prob}(\text{false match}) \leq \frac{\Pi(mn)}{\Pi(M)} \leq \frac{2.511}{n}$$

by Fact 2.

Updating Modular Fingerprints

pattern $X = X_1 \dots X_m$

text $Y = Y_1 \dots Y_n$

$X(r) = X, \quad Y(r) = Y_r Y_{r+1} \dots Y_{r+m-1}$

Since $Y(r+1) = (Y(r) - 2^{m-1} Y_r) \cdot 2 + Y_{r+m}$

update formula:

$a_p(r+1) = (a_p(r) + a_p(r) + \xi Y_r + Y_{r+m}) \bmod p$

where $\xi = -2^m \bmod p$

gives $\Phi_p(Y(r+1)) = a_p(r+1)$ from $\Phi_p(Y(r)) = a_p(r)$
in $O(1)$ arithmetic steps using

$O(\log n)$ bit integers on range $[0, mn^2 - 1]$

Theorem

**Total Exp. Time for finding a match
is $O(n)$**

proof

Expected Time is:

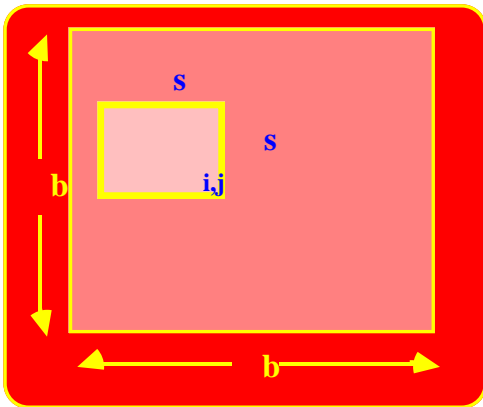
$O(n) + nm \text{ Prob}(\text{false match}) \leq O(n) + nm O\left(\frac{1}{n}\right)$
 $\leq O(n)$

2D Randomized Pattern Matching

input pattern $X = (X_{ij})$ is $s \times s$ boul. array

text $Y = (Y_{ij})$ is $b \times b$ boul. array

text window $Y(\langle i,j \rangle)$ = concatenation of rows of $s \times s$ subarray of Y with $\langle i,j \rangle$ in lower right corner



Index i Update

$$\text{fingerprint } a_p(\langle i+1, j \rangle) = a_p(\langle i, j \rangle) + \left[\left(a_p(\langle i, j \rangle) - \lambda \cdot Y_{i-s+1, j} + Y_{i+1, j} \right) \right] \bmod p$$

where $\lambda = 2^s \bmod p$

Index j Update

$$a_p(\langle i, j+1 \rangle) = a_p(\langle i, j \rangle) \cdot 2^{-s} + \left[Y_{i, j+1} \cdot 2 + Y_{2, j+1} \cdot 2^2 + \dots + Y_{s, j+1} \cdot 2^s \right] \theta - \left[Y_{i, 1} \cdot 2 + Y_{2, 1} \cdot 2^2 + \dots + Y_{s, 1} \cdot 2^s \right] \bmod p$$

where $\theta = -2^{s(s-1)} \bmod p$

Unimodular Matrices as Fingerprints

Definition

homomorphism k from $\{0,1\}^*$ into unimodular matrices with

$$k(\epsilon) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad k(0) = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \quad k(1) = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

$$\begin{array}{ccc} k(X * Y) & = & k(X) \cdot k(Y) \\ \uparrow & & \uparrow \\ \text{concatenation} & & \text{matrix multiplication} \end{array}$$

Fact 1' If $X \in \{0,1\}^m$, then each entry of $k(X)$ is $\leq F_m = m$ th Fibonacci
(where $F_0 = F_1 = 1$, and $F_m = F_{m-1} + F_{m-2}$ for $m \geq 2$)

Fact 2' $\log F_m \sim .694m$

Suppose the unimodular fingerprint fn Φ_p is used with random $p \in S = \{p | p \leq M \text{ is prime}\}$ where $M = mn^2$.

Theorem

The Random Pattern Matching Algorithm

has probability $\leq \frac{6.971}{n}$ of false match.

proof

A false match occurs if $\exists i,j \in \{1,2\}$

$\Phi_p(\mathbf{X}(r))_{i,j} = \Phi_p(\mathbf{Y}(r))_{i,j}$ but $k(\mathbf{X}(r))_{i,j} \neq k(\mathbf{Y}(r))_{i,j}$

iff $p \mid L'$ where $L' = \prod_{\substack{r \in \mathbf{R} \\ i,j \in \{1,2\}}} |k(\mathbf{X}(r))_{i,j} - k(\mathbf{Y}(r))_{i,j}|$

But $L' \leq (F_m)^{4n} \leq 2^{\lceil 4n \log F_m \rceil}$

By fact 1, the number of primes that

divide L' is at most $\Pi\left(\lceil 4n \log F_m \rceil\right)$

The probability of a false match is

$$\frac{\Pi\left(\lceil 4n \log F_m \rceil\right)}{\text{LCM}} \leq \frac{6.971}{n} \text{ since by Fact 2' } \log(F_m) \sim .694n$$

Updating Using Unimodular Matrices

ID string matching: $a_p(r) = \Phi_p(\mathbf{Y}(r))$

$$a_p(r+1) = \mathbf{K}_p(\mathbf{Y}_r)^{-1} a_p(r) \mathbf{H}(\mathbf{Y}_{r+m}) \pmod p$$

$$\text{where } \mathbf{K}_p(0)^{-1} = \begin{pmatrix} 1 & 0 \\ p-1 & 1 \end{pmatrix}$$

$$\text{and } \mathbf{K}_p(1)^{-1} = \begin{pmatrix} 1 & p-1 \\ 0 & 1 \end{pmatrix}$$

(can also extend to 2D string matching)

Simplifications of Modular Fingerprinting

$S = \{p \mid p \leq M \text{ and } p \text{ prime or pseudoprime}\}$

p is *pseudoprime* if $2^{p-1} \equiv 1 \pmod p$ but p not prime

pseudoprimes is $\leq Me^{-c(\log m \log \log m)^{-1}}$

19

(2) $S = \{p \mid p \leq M \text{ and } p \text{ is } M\text{-fat}\}$
 p is *M-fat* if $p \leq M$ and p has
 prime divisor $> \sqrt{M}$

fact

$|S| \sim M(\ln 2 + o(1))$

A false match occurs when $p|L$

where $L = \prod_{X(r) \neq Y(r)} (H(X(r)) - H(Y(r)))$

Bound $L < 2^{nm}$

Let N be the number of M -fat integers dividing L

Then $(\sqrt{M})^{\frac{N}{\sqrt{M}}} < L < 2^{nm}$ so

$$N \leq \frac{\ln 2}{2} \frac{n^4}{(\log n)^2}$$

The *prob of false match* is:

$$\frac{N}{|S|} < .5 \quad \text{if} \quad M = \frac{n^4}{(\ln n)^2}$$

20

(3) $S = \{p \mid p \leq M\}$ with some M
idea use *new* p when get false match

expected time

$$cn \left(.5 + (.5)^2 + (.5)^3 + \dots \right) \leq O(n)$$

Fingerprinting by Random Polynomials

Galois Field

$$\text{GF}(2^k) = \{b_1 \dots b_k \mid b_1, \dots, b_k \in \{0,1\}\}$$

$\mathbb{Z}_2[t]$ = polynomials of form

$$p(t) = t^k + a_{k-1} \cdot t^{k-1} + \dots + a_0 \quad \text{where}$$

$$a_{k-1}, a_{k-2}, \dots, a_0 \in \{0,1\}$$

$p(x)$ irreducible

if can't be factored

Lemma

If k is prime, the number of irreducible polynomials of degree k in $\mathbb{Z}_2[t]$

$$\text{is } \frac{(2^k - 2)}{k}$$

Fingerprint fn $\Phi_p(x) = x_1 t^{m-1} + \dots + x_m \text{ mod } p(t)$
(residue comp can be done efficiently)

Theorem

If use random Fingerprint fn
with degree $k > \log(nm \epsilon^{-1})$, then prob of
false match is $< \epsilon$.

proof

use usual argument and above Lemma

Open Problems

- (1) Are there deterministic methods
for Fingerprinting?
- (2) What are optimal trade offs for
prob of error and size of S
for randomized Fingerprinting?

*Application of Fingerprinting
to Computer Security*

for fixed random p ,

idea store $\Phi_p(F_1), \dots, \Phi_p(F_k)$

fingerprints of files F_1, \dots, F_k

Security:

only operator knows p , so if

any file F_i modified, to F_i'

then with high likelihood

$$\Phi_p(F_i') \neq \Phi_p(F_i)$$

⇒ can build a secure operating system

from this idea!