

Karp-Rabin fingerprint

- How to check equality in files?
- How to check for computer viruses?

It can be done deterministically, but it is slow!

Idea: $\Sigma = \{0,1\}$ binary sequences

sequence $S = 011101110111011$, $w = |S|$

S can be seen as a huge number $\in [0..2^n - 1]$

another sequence $S' = \dots \dots \dots$

Check whether $S = S'$? It can be done

same file but different names, same byte but ^{in O(n) time} different file

$F(s) = \frac{s}{q} \% p$ for a randomly chosen prime (see pseudocode)

both & sequence
and a number

Issue : collision $k_1, k_2 \in \Sigma^*$

$$F(k_1) \neq F(k_2) \Rightarrow k_1 \neq k_2$$

$F(k_1) = F(k_2) \xrightarrow{k_1 = k_2}$
 $\xleftarrow{k_1 \neq k_2}$ 1-sided error
error \rightarrow probability

Let n be the number of bits representing k_1, k_2
 $\Rightarrow k_1, k_2 \in [0..2^n - 1]$

$$\text{error} = k_1 \neq k_2 \wedge F(k_1) = F(k_2) \quad \text{1-side error}$$

Idea:

- choose σ parameter γ sufficiently large &
see later
- choose uniformly and randomly
 σ prime in $[2 \dots \gamma]$ & BAD prime ?

Before starting the computation

$$F(k_1) = F(k_2) \Rightarrow k_1 \% p = k_2 \% p \quad (k_1 \neq k_2)$$

We call such a prime p **BAD**

$$\Pr(\text{error}) = \Pr_{\substack{\text{primes} \\ p \leq \tilde{\gamma}}} (k_1 \neq k_2 \wedge F(k_1) = F(k_2))$$

$$= \frac{\# \text{ BAD primes} \leq \tilde{\gamma}}{\# \text{ primes } p \leq \tilde{\gamma}}$$

$$= \frac{?}{\tilde{\gamma} / \ln \tilde{\gamma}} \quad (*)$$

Obj

both k_1 and k_2 require n bits of representation

Fact Given an integer k on n bits, there are at most $\leq n$ primes in its prime factorization

$$k = p_1^{e_1} \cdot p_2^{e_2} \cdots p_r^{e_r}, \quad r \leq n$$

Proof $p_i \geq 2 \Rightarrow e_1 + e_2 + \cdots + e_r \leq n$ as otherwise

$k \geq 2^{e_1 + e_2 + \cdots + e_r}$ would be greater than 2^n
which is impossible as k requires n bits

CLAIM: there are $\leq n$ BAD primes

- Common primes for k_1 and k_2 are $\leq n$ in number

$$\Pr(\text{error}) = (*) \leq \frac{n}{\gamma / \ln \gamma} = \frac{1}{n^c} \quad \gamma \sim n^{c+1} \ln n$$

Probabilistic algorithms : w.h.p. with high probability

$$= \frac{1}{\text{poly}(n)}$$

Las Vegas	Monte Carlo
no errors	1-side 2-side errors
running time is expected (random var.)	running time is worst-case time

answer

$$\frac{\gamma}{\ln \gamma} = n^{c+1}$$

Pattern matching :

pattern $P[0..m-1]$

text $T[0..n-1]$

$m \leq n$

P occurs at position i in T if $P = T[i..i+m-1]$

GOAL : find whether P occurs in T (virus scan has several simultaneous patterns)
(See pseudo code)

MonteCarlo : if $F(P) = F(T[i..i+m-1])$ then YES

Pr.error : UNION BOUND $\Pr(A \cup B) \leq \Pr(A) + \Pr(B)$

1-side error : no occurrence is missed
false positive

$$\Pr(\text{error}) \leq \frac{m}{\epsilon/\ln \epsilon} \times (n-m+1) \leq \frac{nm}{\epsilon/\ln \epsilon} \leq \frac{n^2}{\epsilon/\ln \epsilon} \leq \frac{1}{n^c}$$

$$\epsilon \sim n^{c+2} \ln n$$

Las Vegas: if $F(P) = F(T[i..i+m-1])$ then

$\underbrace{\text{if } P = T[i..i+m-1] \text{ then YES}}$
 $\underbrace{\text{O}(m) \text{ time}}$

Small change: when $F(P) = F(T[i..i+m-1])$ and $P \neq T[i..i+m-1]$

we "switch" to the trivial $O(nm)$ time alg.
conceptually

Expected cost:

$m \leq n$

$$O(nm) \cdot \frac{1}{n^c} + O(n+m) \cdot \underbrace{\left(1 - \frac{1}{n^c}\right)}_{\leq 1} = O(n)$$

$c > 2$

$$X_i = \begin{cases} 1 & \text{if } F(P) = F(T[i..i+m-1]) \wedge P \notin T[i..i+m-1] \\ 0 & \text{o.w.} \end{cases}$$

$$X = \sum_{i=0}^{n-m} X_i$$

Real w.r.t $O(n + X_m)$ time

Expect cost $O(n + E[X]m)$ time (Δ)

$$E[X] = (n-m+1) \cdot E[X_c] \leq n/m^c$$

$\Pr[X_i=1] \leq \frac{1}{m^c}$

$$(A) \quad O\left(n + \frac{n}{m^c} \cdot m\right) = O(n) \text{ time}$$

$$\frac{n}{m^{c-1}} \leq n$$

$c > 1$

HANDS-ON

- Set of sequences s_1, s_2, \dots, s_k , $m > 0$
- change a sequence by adding/removing characters
- issue an alert when any two sequence s_i, s_j ($i \neq j$) share a common substring of length m