

Today: Fingerprinting for Pattern Matching

- Find occurrences of pattern P in long document D
- $O(|P||D|)$ obvious algorithm
- Lots of methods

Approach 1: Hash Pattern to x -bit string

Compare hash value to hash value

Hash Function: mod prime

$$P = 17935 \quad p = 251$$

$$P \bmod p = 114$$

$$0386179357342$$

$$03861 \bmod p = 107$$

$$38617 \bmod p = 214$$

$$86179 \bmod p = 86 \dots$$

$$17935 \bmod p = 114$$

$$57342 \bmod p = 114$$

$$N \bmod p \rightarrow N' \bmod p \quad \text{(first digit } a, \text{ last digit } b)$$

$$N' = (N - 10^{|P|-1} \times a) \times 10 + b$$

$$N' \bmod p = ((N - 10^{|P|-1} \times a) \times 10 + b) \bmod p$$

Assuming $\bmod p$ operations are $O(1)$ → not constant but only calculate once.

$$\pi(x) = \# \text{ of primes } \leq x$$

$$\lim_{x \rightarrow \infty} \pi(x) \text{ grows } x / \ln(x)$$

$$\frac{x}{\ln(x)} \leq \pi(x) \leq 1.26 x / \ln x$$

$$\rightarrow \text{for } 10^{100}, 10^{100} / 100 \times \ln(10) \text{ primes}$$

Random Prime

FP: two substrings P and A have $(P-A) \bmod p = 0$
 $\Pr[\text{any false positives in document}] =$

$$\frac{\# \text{ primes } p \text{ s.t. } p \mid P-A}{\pi(z)}$$

$$\text{pick a prime} \\ 1 \leq p \leq z$$

$$|P-A| \leq 10^{|P|}$$

How many primes are a factor of n , where $n \leq 10^{|P|}$

Claim. n has at most $\log_2(n)$ distinct factors

$$n \geq \pi \text{ prime factors}$$

$$\geq 2^{\# \text{ prime factors}}$$

$$\# \text{ prime factors} \leq \log_2 n$$

$$\Pr[\text{FP}] \leq \frac{\log_2 10^{|P|}}{\pi(z)}$$

$$\text{Expected \# of FP} = |D| \times \Pr[\text{FP}] = |D| \times \log_2 10^{|P|} / \pi(z)$$

* can decrease by using multiple hash functions (primes)

$$\Pr[\text{FP}] = \left(\frac{\log_2 10^{|P|}}{\pi(z)} \right)^k \text{ where } k \text{ primes used}$$

Primality Test

To pick random prime:

→ pick random #
 test if prime
 repeat if necessary

Test if prime:

is 2 a factor
 3 a factor
 :

√n a factor

$n \rightarrow \log_2 n$ digits
 want polynomial in $\log_2 n$

Fermat's Little Theorem:

If p is a prime, $1 \leq a < p$, $a^{p-1} = 1 \pmod{p}$

Pf. $\{1, \dots, p-1\}$ and $\{a \cdot 1, a \cdot 2, \dots, a(p-1)\}$
 \pmod{p} \pmod{p}

In the second set, all #s between 1, $p-1$

are different because if $a \cdot i$ and $a \cdot j$ were same:

$$\begin{array}{|l}
 a \cdot i = a \cdot j \pmod{p} \\
 \hookrightarrow i = j \pmod{p} \rightarrow \text{contradiction}
 \end{array}$$

\pmod{p} there are multiplicative inverses

$$1 \times 2 \times 3 \dots \times (p-1) =$$

$$(a \times 1)(a \times 2) \dots (a \times (p-1)) \pmod{p}$$

$$= a^{p-1} \times 1 \times 2 \times \dots \times (p-1) \pmod{p}$$

$$1 = a^{p-1} \pmod{p}$$

FLT Primality Test

Given n , is n prime?

Test: compute 2^{n-1} square using repeated squaring
 if $2^{n-1} \equiv 1 \pmod n$ say prime
 $2^{n-1} \not\equiv 1 \pmod n$ say composite

Correctness:

only one direction, 1 means not prime

FP possible

2-pseudoprime numbers (satisfies Fermat's Little Thm but not prime)

Carmichael numbers: composites s.t. $a^{n-1} \equiv 1 \pmod n$ for any a that does not share a factor with n .
↑ infinite #!

Improvement: Rabin-Miller Primality Test

Composites have nontrivial square roots of 1

trivial: $(-b)^2 \equiv 1 \pmod 7 \rightarrow$ trivial

$$x^2 \equiv 1 \pmod n$$

$$x^2 - 1 \equiv 0 \pmod n$$

nontrivial: $(4)^2 \equiv 1 \pmod{15} = (11)^2$

$$(x+1)(x-1) \equiv 0 \pmod n$$

compute a^{n-1}

$$a^n, a^{2n}, a^{4n}, a^{8n}, \dots, a^{2^v n} = a^{n-1}$$

$n-1 = 2^t u$ where u is odd

if $a^{n-1} \not\equiv 1$, composite!

elif $a^{n-1} \equiv 1$, go back through the chain as long as you see a 1. if see -1 at end, prime. if see anything else, composite.

For any composite odd n , for at least $3/4$ of the vals, $1 \leq a < n$, the test will return composite

$$\leq \left(\frac{1}{4}\right)^{50 \text{ times}} = 2^{-100}$$