

CS70 - Spring 2024

Lecture 27 - April 25

Today: Intro. to Randomized Algorithms

- Finding large prime numbers
- Fingerprinting
- Pattern matching

1. Primality testing

Cryptographic applications (e.g. RSA) require very large prime numbers (~100s of digits)

Q: How do we get hold of them?

A: Generate a random number with 100s of digits & test if it is prime!
If not, try another one...

Q: How many numbers do we expect to try?

A: Depends on the density of primes

Prime Number Theorem: Let $\pi(n)$ denote the number of primes $\leq n$. Then

$$\frac{n}{\ln n} \leq \pi(n) \leq 1.26 \frac{n}{\ln n} \quad \forall n \geq 17$$

Corollary: Roughly 1 in every $\ln n$ numbers $\leq n$ is prime. [Eg. $n = 10^{500} \Rightarrow \ln n \approx 1150$]

So no. of trials until we find a prime has Geometric ($1/\ln n$) distribution

$$\Rightarrow E[\# \text{ trials}] \approx \ln n$$

$$P[\text{more than } k \ln n \text{ trials}] = \left(1 - \frac{1}{\ln n}\right)^{k \ln n} \leq e^{-k}$$

$$\left(1 - \frac{1}{m}\right)^m \sim e^{-1}$$

Q: How do we get hold of them?

A: Generate a random number with 100s of digits
& test if it is prime!
If not, try another one...

Bigger Question: How do we test if a very large number n (with 100s of digits) is prime?

Simple algorithm: try all divisors!

for $a = 2, 3, \dots, \sqrt{n}$

if a divides n then halt & output "not prime"

output "prime"

Simple algorithm : try all divisors!

for $a = 2, 3, \dots, \sqrt{n}$

if a divides n then halt & output "not prime"

output "prime"

Is this a good algorithm?

How many divisors will we have to try?

Think : $n \approx 10^{500}$ (500 digits, ≈ 1700 bits)

A simple randomized algorithm

repeat many times

pick $a \in [2, \sqrt{n}]$ uniformly at random

if a divides n then halt & output "not prime"

output "prime?"



witness

Is this a good algorithm?

A better randomized algorithm?

Need a better witness for n being not prime (so that the number of witnesses is large)

Recall: Fermat's Little Theorem

For any prime p and any $a \in [1, \dots, p-1]$:

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

Corollary: If we find an $a \in [1, \dots, n-1]$ s.t.

$$a^{n-1} \neq 1 \pmod{n}$$

then we know for sure that n is not prime!

A better algorithm: Fermat Test

pick $a \in [1, \dots, n-1]$ u.a.r.

if $\gcd(a, n) \neq 1$ then halt & output "not prime"

if $a^{n-1} \neq 1 \pmod{n}$ then _____ .. _____
else output "prime?"

Properties

- Outputs "not prime" \Rightarrow n is definitely not prime
- Outputs "prime?" \Rightarrow either n is prime, or it's not prime but the algorithm picked an a that's not a witness
- Running time: $O(\log n) = O(\# \text{ of } \underline{\text{digits}} \text{ in } n)$

- Running time: $O(\log n) = O(\# \text{ of } \underline{\text{digits in } n})$

Why?

→ $\text{gcd}(a, n)$ runs in $O(\log n)$ steps via Euclid's Algorithm

→ a^{n-1} can be computed in $O(\log n)$ steps by repeated squaring:

e.g. for a^{53}

$$53 = 32 + 16 + 4 + 1 = 110101 \text{ in binary}$$

$$a^{53} = a^{32} \times a^{16} \times a^4 \times a^1$$

→ enough to compute
 $a, a^2, a^4, a^8, a^{16}, a^{32}$

Density of witnesses

Let $\mathbb{Z}_n^* = \{a \in [1, \dots, n-1] : \gcd(a, n) = 1\}$

Defn: A witness for n is a number $a \in \mathbb{Z}_n^*$
s.t. $a^{n-1} \neq 1 \pmod{n}$

Claim: For any n , if \exists a witness a for n
then at least half of all $a \in \mathbb{Z}_n^*$ are witnesses!

Corollary: The Fermat Test is correct with
probability $\geq 1/2$ on all inputs n except
for non-primes n that have no witnesses at all

"Carmichael Numbers"

A better algorithm: Fermat Test

pick $a \in [1, \dots, n-1]$ u.a.r.

if $\gcd(a, n) \neq 1$ then halt & output "not prime"

if $a^{n-1} \neq 1 \pmod{n}$ then _____ .. _____
else output "prime?"

Properties

- Outputs "not prime" \Rightarrow n is definitely not prime
- Outputs "prime?" \Rightarrow either n is a Carmichael Number, or

$$\Pr[n \text{ is prime}] \approx 1/2$$

Claim: For any n , if \exists a witness a for n then at least half of all $a \in \mathbb{Z}_n^*$ are witnesses!

Proof: \mathbb{Z}_n^* is a group under multiplication (mod n). I.e.:

$$\begin{array}{c} \textcircled{0^s} \\ |S| \mid |G| \end{array}$$

- [Identity]: $1 \in \mathbb{Z}_n^*$
- [Inverses]: $a \in \mathbb{Z}_n^* \Rightarrow a^{-1} \in \mathbb{Z}_n^*$
- [Closure]: $a, b \in \mathbb{Z}_n^* \Rightarrow ab \in \mathbb{Z}_n^*$

The set $S = \{a \in \mathbb{Z}_n^* : a^{n-1} = 1 \pmod{n}\}$ is a subgroup (because $a, b \in S \Rightarrow ab \in S$)

Lagrange's Theorem: $|S|$ divides $|\mathbb{Z}_n^*|$

n is not a CN $\Rightarrow |S| < |\mathbb{Z}_n^*| \Rightarrow |S| \leq \frac{1}{2} |\mathbb{Z}_n^*|$ ✓

Two remaining questions:

1. Why is $\Pr[n \text{ is prime}] \approx 1/2$ good enough?
2. What about Carmichael Numbers?

1. Why is $\Pr[n \text{ is prime}] \approx 1/2$ good enough?

A: Just repeat the Fermat Test k times, choosing a independently each time

Suppose n is not prime and not a CN

Then $\Pr[\text{all } k \text{ tests output "prime?"}] \leq 2^{-k}$

If we take $k = 1000$ (say), this is negligible!

2. What about Carmichael Numbers?

Defn: n is a Carmichael Number if it is not prime and $a^{n-1} = 1 \pmod{n} \quad \forall a \in \mathbb{Z}_n^*$

CNs are rare: $255 \text{ CNs} \leq 10^8$
 $\sim 20 \text{ million CNs} \leq 10^{21}$

First few CNs: 561, 1105, 1729, ...

- There are similar randomized algorithms (using more complicated witnesses) that handle CNs
- There are also deterministic algorithms (that are always correct) but they are too inefficient ($\sim O((\log n)^6)$)

Proof that 561 is a Carmichael Number:

$$561 = 3 \times 11 \times 17$$

Claim: $a^{560} = 1 \pmod{561} \quad \forall a \in \mathbb{Z}_{561}^*$

Sufficient to show $a^{560} = 1 \pmod{3}$
 $a^{560} = 1 \pmod{11}$
 $a^{560} = 1 \pmod{17}$ } \Rightarrow by C.R.T.
 $a^{560} = 1 \pmod{561}$

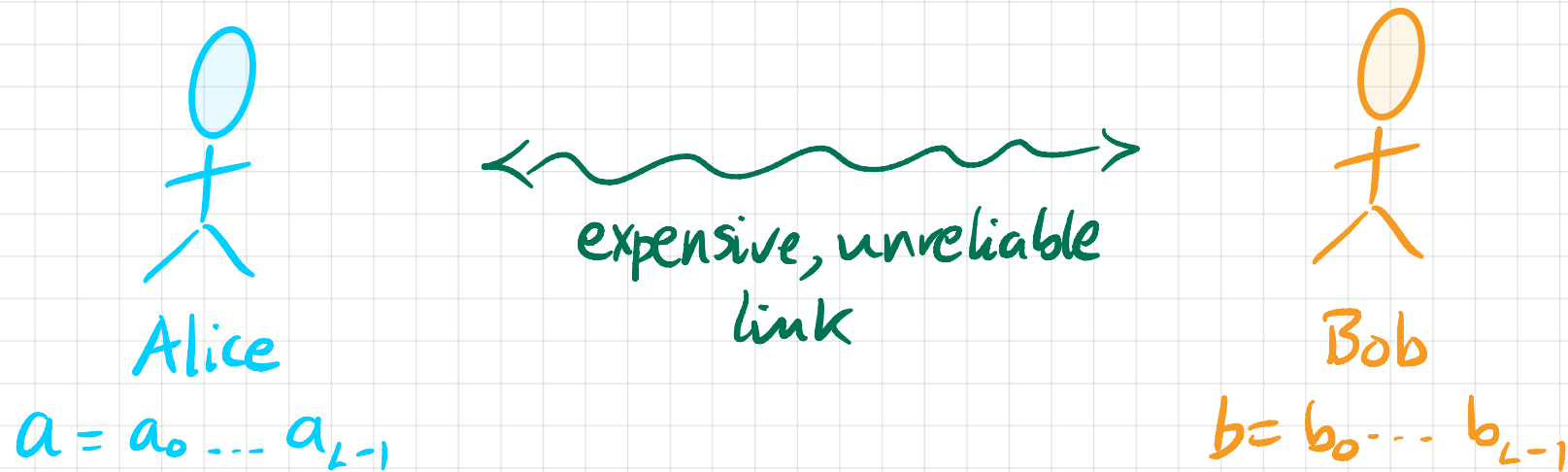
Now note: $a^2 = 1 \pmod{3} \Rightarrow a^{560} = 1 \pmod{3}$

$$a^{10} = 1 \pmod{11} \Rightarrow a^{560} = 1 \pmod{11}$$

$$a^{16} = 1 \pmod{17} \Rightarrow a^{560} = 1 \pmod{17}$$



2. Fingerprinting & Pattern Matching

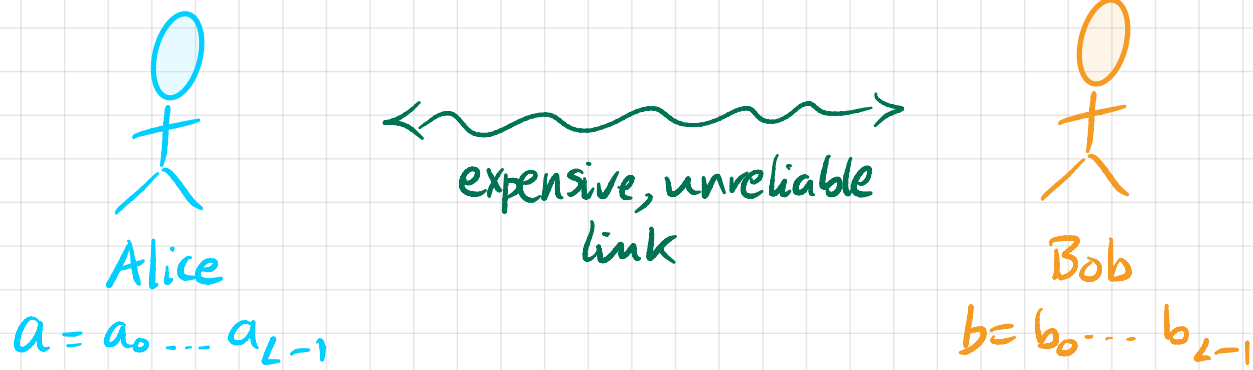


Alice & Bob each have a copy of a large database consisting of L bits (L very large)

They want to check if their copies are identical

But they don't want to send all L bits

Idea: Send a much smaller fingerprint of their data



Idea: Send a much smaller fingerprint of their data

View $a = a_0 \dots a_{L-1}$ & $b = b_0 \dots b_{L-1}$ as L -bit numbers

Alice: picks a random prime $p \in [2 \dots T]$

computes $F_p(a) := a \bmod p$

sends p and $F_p(a)$ to Bob

Bob: computes $F_p(b) := b \bmod p$

if $F_p(a) \neq F_p(b)$ outputs "not identical"
else outputs "identical?"

Alice: picks a random prime $p \in [2 \dots T]$

computes $F_p(a) := a \bmod p$

sends p and $F_p(a)$ to Bob

Bob: computes $F_p(b) := b \bmod p$

if $F_p(a) \neq F_p(b)$ outputs "not identical"

else outputs "identical?"

Properties

Outputs "not identical" \Rightarrow definitely $a \neq b$

Outputs "identical?" \Rightarrow either $a = b$ or

$a \neq b$ but $a = b \bmod p$

\swarrow error

Outputs "identical?" \Rightarrow either $a = b$ or
 $a \neq b$ but $a = b \pmod p$

Claim: If $a \neq b$ and p is a random prime in $[2..T]$
then $\Pr[a = b \pmod p] \leq \frac{L \ln T}{T}$

Proof: $a = b \pmod p \iff a - b = 0 \pmod p$
 $\iff p \mid |a - b|$

But $|a - b|$ is an L -bit number, so it has
at most L prime factors!

Hence $\Pr[a = b \pmod p] \leq \frac{L}{\pi(T)} \leq \frac{L \ln T}{T}$ ✓
prime number thm.

Corollary: $\Pr[\text{error}] \leq \frac{L \ln T}{T}$

Corollary : $\Pr[\text{error}] \leq \frac{L \ln T}{T}$

How should we set T ?

If we set $T = 4L \ln L$ then

$$\Pr[\text{error}] \leq L \cdot \frac{\ln L + \ln \ln L + \ln 4}{4L \ln L}$$

$$= \frac{1}{4} \left[1 + \frac{\ln \ln L}{\ln L} + \frac{\ln 4}{\ln L} \right]$$

$$\leq \boxed{\frac{1}{2}}$$

→ Can boost to $\leq 2^{-k}$ with k independent trials

Also: $T = 4L \ln L \Rightarrow p$ has $O(\log L)$ bits

So fingerprint $F_p(a)$ sent by Alice is exponentially smaller than the database itself!

Example :

$$\text{Suppose } L = 2^{33} \quad (\approx 1 \text{ GB})$$

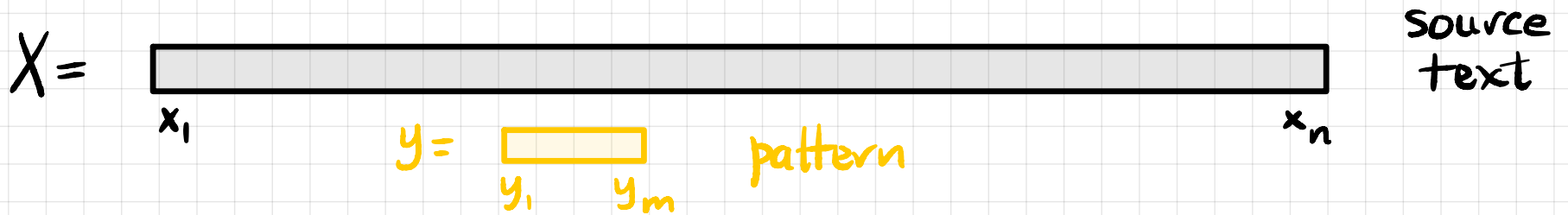
$$T = 2^{64} \quad (\Rightarrow \text{fingerprints are 64-bit words})$$

$$\text{Then } \Pr[\text{error}] \leq \frac{L \ln T}{T} \leq 2^{33} \times \frac{64}{2^{64}}$$

$$= 2^{-25}$$

$$\approx \boxed{3.4 \times 10^{-7}}$$

3. Pattern Matching



Question : Does Y occur as a contiguous substring in X ?
I.e., is $Y = X(i) := X_i X_{i+1} \dots X_{i+m-1}$ for some i ?

Simple algorithm :

for $i := 1$ to $n - m + 1$

if $Y = X(i)$ output "match found" & halt

output "no match"

Running time : $O(nm)$

Clever randomized algorithm: use fingerprints!

pick a prime $p \in [2, \dots, T]$ u.a.r.

compute $F_p(y) := y \bmod p$

for $i := 1$ to $n - m + 1$

compute $F_p(X(i)) := X(i) \bmod p$

if $F_p(y) = F_p(X(i))$ output "match found" & halt

output "no match"

Clever randomized algorithm: use fingerprints!

pick a prime $p \in [2, \dots, T]$ u.a.r.

compute $F_p(y) := y \bmod p$

for $i := 1$ to $n - m + 1$

compute $F_p(x(i)) := x(i) \bmod p$

if $F_p(y) = F_p(x(i))$ output "match found" & halt

output "no match"

$$\Pr[\text{error}] = \Pr[\exists i : y \neq x(i) \ \& \ y - x(i) = 0 \pmod{p}]$$

union bound \rightarrow

$$\leq \sum_{i=1}^{n-m+1} \Pr[y \neq x(i) \ \& \ y - x(i) = 0 \pmod{p}]$$

$$\leq n \times \frac{m \ln T}{T}$$

[$y - x(i)$ has m bits
 $\Rightarrow \leq m$ prime factors]

So if we set $T = 4nm \ln(nm)$ then $\Pr[\text{error}] \leq \frac{1}{2}$

So if we set $T = 4nm \ln(nm)$ then $\Pr[\text{error}] \leq \frac{1}{2}$

With this choice of T , number of bits in p is $O(\log(nm)) = O(\log n)$.

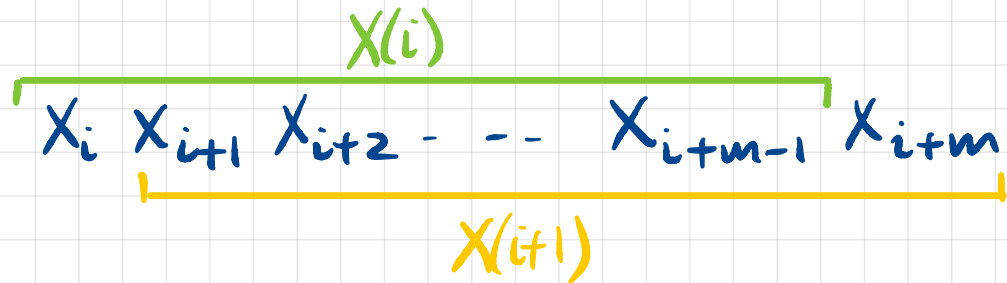
So can assume arithmetic mod p is fast.

Running time of algorithm?

- compute $F_p(y)$ $O(m)$
- compute $F_p(x(1))$ $O(m)$
- n iterations:
 - compute $F_p(x(i))$ $O(1)$
 - compare $F_p(y) = F_p(x(i))$ $O(1)$ } $O(n)$

Total: $O(n)$ [much faster than $O(nm)$]

Computing $F_p(X(i+1))$ from $F_p(X(i))$



$$X(i+1) = 2(X(i) - 2^{m-1}X_i) + X_{i+m}$$

$$\Rightarrow F_p(X(i+1)) = 2(F_p(X(i)) - 2^{m-1}X_i) + X_{i+m} \pmod{p}$$

requires four arithmetic ops. mod $p \rightarrow O(1)$ time

Example :

Searching for a pattern in a DNA sequence

Say $n = 2^{28}$, $m = 2^{11}$ (search for 1000-bp pattern in 100M-bp chromosome)

Take $T = 2^{64}$ (64-bit words)

$$\Pr[\text{error}] \leq nm \frac{\ln T}{T} \leq 2^{39} \cdot \frac{64}{2^{64}}$$

$$= 2^{-19}$$

$$\approx 0.000005$$

Note: Deterministic $O(n)$ algorithms do exist, but they're much harder to implement and have overheads

Life After CS70

CS 170 — Algorithms

CS 172 — Complexity & Computability

CS 174 — Randomized Algorithms

CS 171 — Cryptography

CS 176 — Computational Biology

EECS 126 — Probability & Random Processes

EECS 127 — Optimization



Course Evaluations : course-evaluations.berkeley.edu