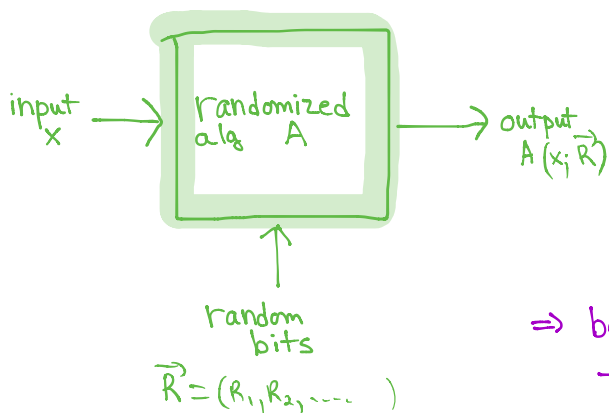
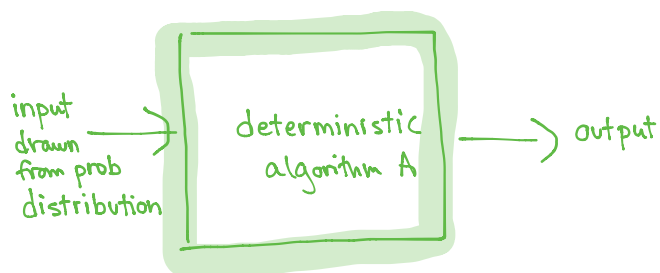


# Randomized Algorithms & Probabilistic Analysis of Algorithms....



Model of computation:  
standard model (TM, RAM)  
with additional input consisting of stream of perfectly random bits.

⇒ behavior can vary on fixed input  
- running time on particular input is a random variable



again, performance of algorithm is a random variable

also other random structures: random graphs, random boolean formulas, etc.

## Example of differences:

quicksort with randomly selected pivots vs QS where input is random  $\Pi$

Why randomized algs?

- often simplest or fastest
- fun!!!

Today: simple examples to illustrate

- searching for a witness
- principle of deferred decisions
- fingerprinting
- probabilistic method
- derandomization
  - method of conditional expts
  - pairwise independence

# Matrix-Product Verification

[MU]1.3 [MR]7.1

Given  $n \times n$  matrices  $A, B, C$  over field  $F$

Told  $AB = C$

Goal: to verify this identity

Obvious method: matrix multiplication

$$O(n^{2.372})$$

Field  $F$ :

Set with 2 operations  
addition, multiplication  
has all the properties of  $\mathbb{R}$   
or rational #'s  
e.g. commutativity,  
associativity  
additive/mult. inverses  
identity elts for  $+, -$

Example:  $GF(2)$

addition XOR (addition mod 2)  
multiplication AND

[Freivalds Alg] simple & elegant

one of first published uses of randomization in algs

Pick random vector  $\vec{r} = (r_1, r_2, \dots, r_n) \in \{0, 1\}^n$

each  $r_i$  indep, equally likely to be 0 or 1

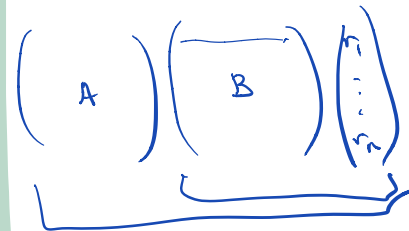
↑ additive identity of field  
↑ multiplicative identity of field

compute  $A(Br) = z$

If  $Cr = z$

then output "yes,  $AB = C$ "

else output "no"



Running Time:  $O(n^2)$

Errors: if  $AB = C$  always output yes

if  $AB \neq C$  may make an error ←

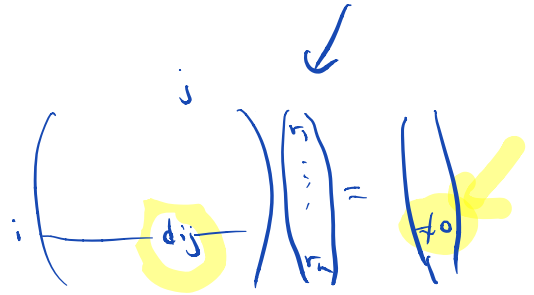
$$AB \neq C$$

Claim:  $\Pr(\text{output an incorrect answer}) \leq \frac{1}{2}$

Proof: Define  $D = AB - C$

Suppose  $D \neq 0$

Then  $\exists$  entry, say  $(i,j)$  s.t.  $d_{ij} \neq 0$



$$\Pr(Dr = 0) \leq \Pr\left(\sum_k d_{ik} r_k = 0\right)$$

$$= \Pr\left(d_{ij} r_j = -\sum_{k \neq j} d_{ik} r_k\right)$$

$$= \Pr\left(r_j = \frac{-\sum_{k \neq j} d_{ik} r_k}{d_{ij}}\right)$$

Example of simple but powerful principle of deferred decisions

multiple r.v.'s - think of setting some of them first  
and deferring setting rest until later  
in analysis

Formally, use law of total probability; condition on values of vars set 1<sup>st</sup>

$$\Pr \left( r_j = \frac{-\sum_{k \neq j} d_{ik} r_k}{d_{ij}} \right)$$

$$= \sum_{\substack{(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \\ \in \{0,1\}^{n-1}}} \Pr \left( r_j = \frac{-\sum_{k \neq j} d_{ik} r_k}{d_{ij}} \mid \underbrace{\begin{matrix} (r_1, \dots, r_{i-1}, r_{i+1}, \dots, r_n) \\ = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \end{matrix}}_A \right) \Pr(A)$$

$\downarrow$   
 number

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

$$\leq \sum_{\substack{(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \\ \in \{0,1\}^{n-1}}} \frac{1}{2} \Pr(A) = \frac{1}{2}$$

If want to reduce probability of error, can do so at expense of small  $\uparrow$  in running time

- ① Run alg  $k$  times
- ② Output yes if get yes all  $k$  times

$$\Pr(\text{error}) \leq \frac{1}{2^k}$$

by independence of trials.

## Searching for witnesses

- using randomization to check whether  $P(r) = \text{True} \forall r$

- Pick random  $r$  from suitable set
- If  $P(r)$  false  $\rightarrow$  "No"
- else  $\rightarrow$  "Yes"

- works well if density of witnesses that  $P(r)$  false high enough

---

## 2 types of randomized algs

**Monte Carlo algo** - halt in finite time but may output wrong answers

- One-sided error  $\uparrow$  confidence with repetition.

- Two-sided error  $\begin{matrix} \text{if true answer yes} \Rightarrow \Pr(\text{output yes}) \geq \frac{1}{2} + \epsilon \\ \text{no} \Rightarrow \Pr(\text{output no}) \geq \frac{1}{2} + \epsilon \end{matrix}$

$\uparrow$  confidence with repetition & majority vote.

Claim: If  $\checkmark$  <sup>decision</sup> alg correct w.p.  $\geq \frac{1}{2} + \epsilon$  & we run it  $t$  times & output majority answer, probability answer correct  $\geq 1 - e^{-2\epsilon^2 t}$

$\Pr(\text{output correct answer}) \geq \frac{1}{2} + \epsilon$

Proof:

$$\begin{aligned} \Pr(\text{majority wrong}) &\leq \sum_{i=0}^{t/4} \binom{t}{i} \left(\frac{1}{2} + \epsilon\right)^i \left(\frac{1}{2} - \epsilon\right)^{t-i} \\ &\leq \sum_{i=0}^{t/4} \binom{t}{i} \left(\frac{1}{2} + \epsilon\right)^{t/4} \left(\frac{1}{2} - \epsilon\right)^{3t/4} = \left(\frac{1}{4} - \epsilon^2\right)^{t/4} \sum_{i=0}^{t/4} \binom{t}{i} \\ &\leq \left(\frac{1}{4} - \epsilon^2\right)^{t/4} 2^t = \left(\frac{1}{4} - \epsilon^2\right)^{t/4} 4^{t/4} = \left(\frac{1 - 4\epsilon^2}{4}\right)^{t/4} \leq e^{-\frac{4\epsilon^2 t}{4}} = e^{-\epsilon^2 t} \end{aligned}$$

Most useful approx:  $1-x \leq e^{-x}$

Las Vegas algorithms

always output correct answers. runtime is r.v.

Ex: randomized Quicksort

Big open question

Does randomness help in computation?

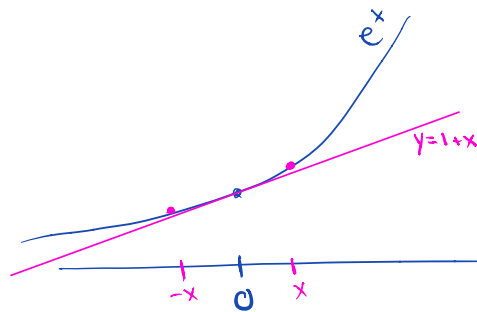
Can every poly time randomized alg be "de-randomized" with at most polynomial loss in efficiency?

Most useful approx:  $1-x \leq e^{-x}$

$e^x$  convex everywhere

$\Rightarrow$  tangent at  $x=0$  lies below curve everywhere

$\Rightarrow e^{-x} \geq 1-x$   
 $e^x \geq 1+x$



# Fingerprinting

[MR] 7.4 [CG] 2.2.1

A & B each have large DB, separated by long distance

↓ ↓  
a b both n-bit strings

want to check if  $a=b$ ?

Deterministically n bits of communication necessary

Next: randomized protocol that uses  $O(\log n)$  bits of communication

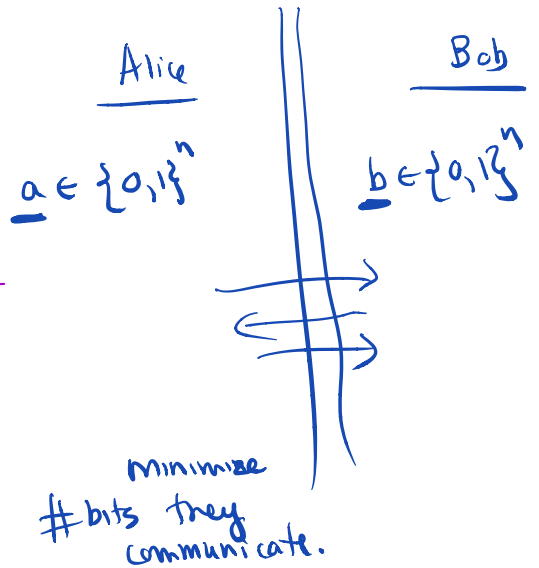
A picks prime  $p \in [2..x]$  u.a.r.  
 ↘ to be determined

A sends  $(p, a \bmod p)$  to B

B computes  $b \bmod p$

If  $a \bmod p = b \bmod p$ , B sends back "yes", else "no"

u.a.r.  
≡ uniformly at random.



Notes need random prime ...

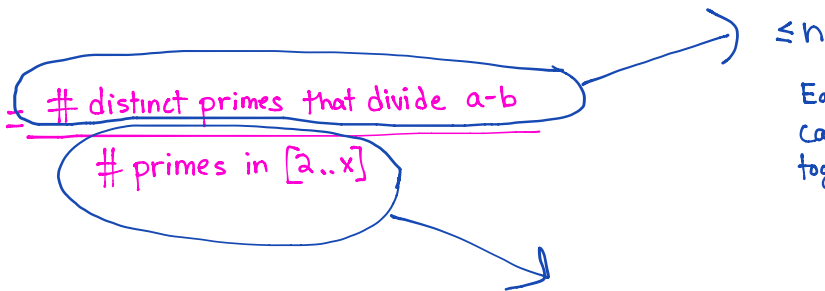
Always gives right answer if  $a=b$ .  
may give wrong answer if  $a \neq b$

Suppose  $a \neq b$

$$\Pr(a \bmod p = b \bmod p) = \Pr(a-b \text{ is multiple of } p)$$

p is one of prime factors of  $a-b$

$$= \frac{\# \text{ distinct primes that divide } a-b}{\# \text{ primes in } [2..x]}$$



Each prime  $\geq 2$   
 Can't multiply  $> n$   
 together before get  $> 2^n$

Prime # Thm:

$$\# \text{ primes } \leq x \geq \frac{x}{\ln x} \quad \forall x \geq 17$$

$$p \leq x$$

$a \bmod p$

$$\leq \frac{n \ln x}{x}$$

choosing  $x = c n \ln n$

$$\frac{x \ln x}{c x \ln n}$$

$$\leq \frac{1}{c} \frac{\ln x}{\ln n} = \frac{1}{c} + o(1)$$

# bits transmitted  
 $= 2 \log x = O(\log n)$

Example:  $n = 2^{33} \approx 1 \text{ gigabyte}$

$x = 2^{64}$  (fingerprints are 64 bit words)

$$\Pr(\text{error}) < 10^{-9}$$



MaxCut [MU]6.2.1 [CG]1.4.1

simple randomized alg

illustration of **probabilistic method**

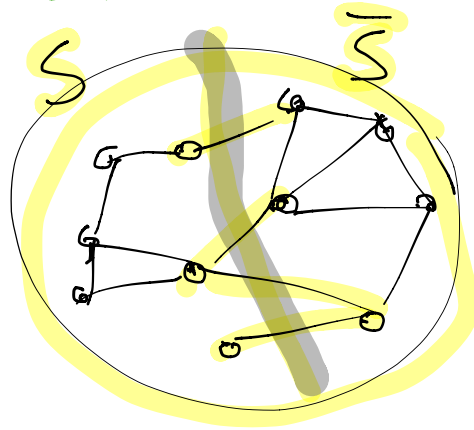
use probabilistic argument to prove  
non-probabilistic mathematical thm

MAXCUT  
Problem: Given  $G=(V,E)$   
find  $S$  that  
maximizes  
 $\text{cut}(S, \bar{S})$   
#edges

Defn cut in graph: partition of nodes into 2 sets  $S$  and  $\bar{S}$   
An edge crosses cut if it has one endpoint in  $S$  & one in  $\bar{S}$

**Thm:**

In any graph  $G=(V,E)$ ,  $\exists$  cut  
s.t. at least  $\frac{1}{2}$  edges cross  
cut.



Proof technique: show that if we pick a random cut, the exp #  
of edges that cross cut is  $\geq \frac{1}{2} |E|$

Pick cut u.a.r.  $\forall v \in V$ , flip fair coin  $\begin{cases} H \rightarrow v \in S \\ T \rightarrow v \in \bar{S} \end{cases}$

Let  $X_e = \begin{cases} 1 & e \text{ crosses cut} \\ 0 & \text{o.w.} \end{cases}$

$$E(X_e) = \Pr(e \text{ crosses cut}) = \frac{1}{2}$$

$X = \sum_{e \in E} X_e$  # edges crossing cut

$$E(X) = ?$$

$$E(X) = E\left(\sum_{e \in E} X_e\right) = \sum_{e \in E} E(X_e) = \frac{1}{2} |E|$$

$\Rightarrow$  sample space must contain at least one cut

in which  $\geq \frac{1}{2}$  edges cross cut. O.W.  $E(X) < \frac{1}{2} |E|$

Typical example of prob method:

- Not everybody can be below (or above) average

- Collection of objects  $\Pr(\exists \text{ object with property } P) > 0$

$\Rightarrow \exists$  object in collection with property P

**MAXCUT**

for  $i := 1$  to  $n$

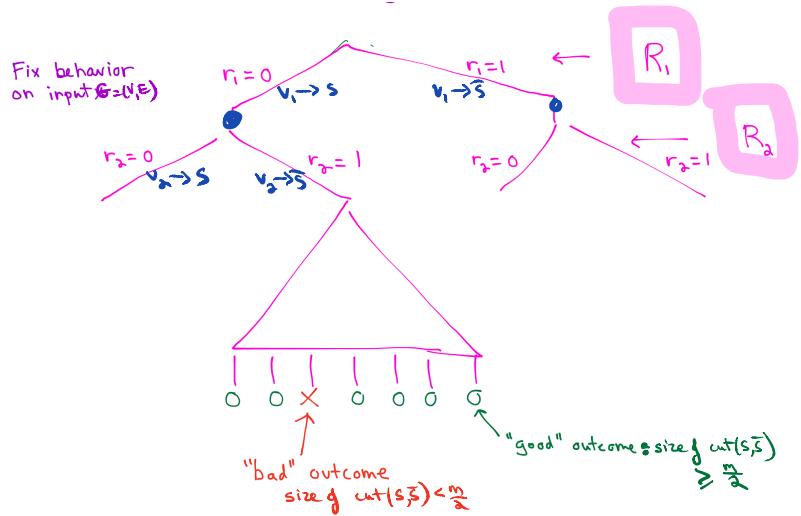
  Let  $R_i \in_{\text{random}} \{0,1\}$

  If  $R_i = 0$ , put  $v_i$  in  $S$ , else put  $v_i \in \bar{S}$

$G = (V, E)$      $|V| = n$      $|E| = m$

$E[\text{Cut}(S, \bar{S})] \geq \frac{m}{2}$

$\uparrow$   
 $R_1, R_2, \dots, R_n$



Idea: walk down tree, making good choice at each step

Observation:  $E(\text{Cut}(S, \bar{S}) \mid R_1=r_1, R_2=r_2, \dots, R_i=r_i) \stackrel{= E[S]}{\geq \frac{m}{2}}$

$= \frac{1}{2} E(\text{Cut}(S, \bar{S}) \mid R_1=r_1, R_2=r_2, \dots, R_i=r_i, R_{i+1}=0)$

$+ \frac{1}{2} E(\text{Cut}(S, \bar{S}) \mid R_1=r_1, R_2=r_2, \dots, R_i=r_i, R_{i+1}=1)$

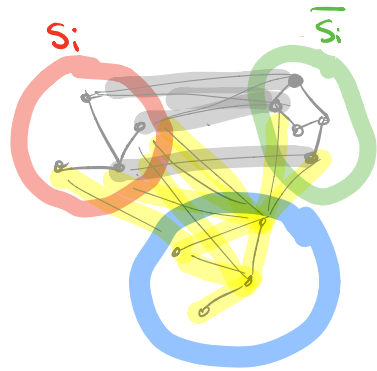
one of these is at least  $\frac{m}{2}$

$v_1, \dots, v_i \rightarrow S, \bar{S}$

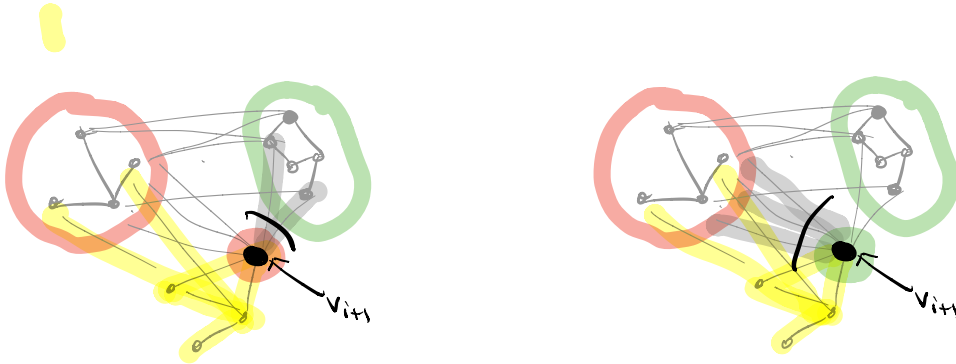
$S_i = \{v_j \mid j \leq i, R_j = 0\}$

$\bar{S}_i = \{v_j \mid j \leq i, R_j = 1\}$

$U_i = \{v_{i+1}, \dots, v_n\}$



$$E(\text{cut}(S, \bar{S}) \mid R_1=r_1, \dots, R_i=r_i) = \underbrace{|\text{cut}(S_i, \bar{S}_i)|}_{\text{red-green}} + \frac{1}{2} \underbrace{|\# \text{ edges with at least one endpoint in } U_i|}_{\text{yellow edges}}$$



$$E(\text{cut}(S, \bar{S}) \mid R_1=r_1, R_2=r_2, \dots, R_{i+1}=r_{i+1}) = |\text{cut}(S_{i+1}, \bar{S}_{i+1})| + \frac{1}{2} \underbrace{|\# \text{ edges with one endpoint in } U_{i+1}|}_{\text{edges of } r_{i+1}}$$

$\Rightarrow$  suffices to set  $r_{i+1}$  to maximize  $|\text{cut}(S_{i+1}, \bar{S}_{i+1})|$

To maximize, pick bigger one  $\Rightarrow$  The greedy algorithm!

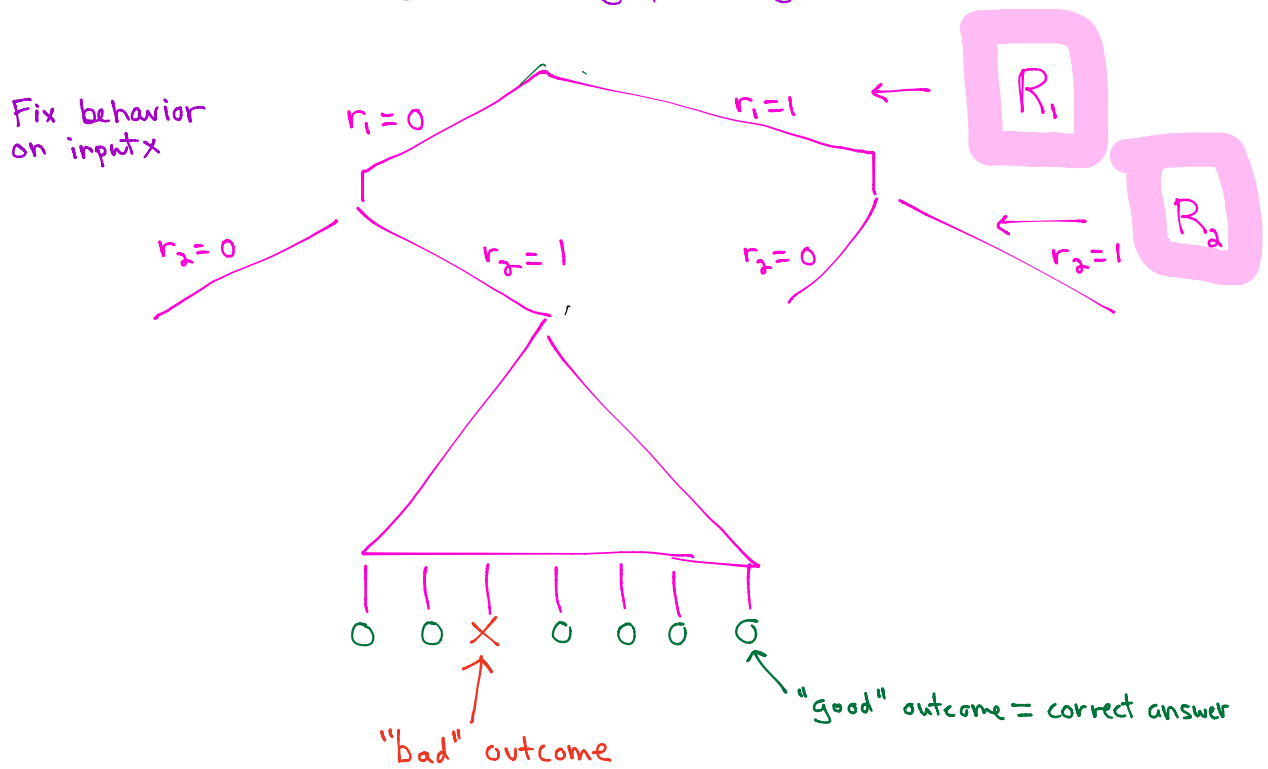
Corollary The greedy alg is guaranteed to find a cut of size  $\geq \frac{|E|}{2}$

# Method of Conditional Expectation

Consider randomized alg  $A$  that uses  $m$  random bits.  $\Pr(A(x; R_1, \dots, R_m) \text{ good})$   
 at least, say,  $\frac{2}{3}$

sequences of coin tosses  $\Leftrightarrow$  binary tree

"Good" randomized alg  $\Rightarrow$  many paths good



$R_1, \dots, R_m$  seq of unif, indep random bits

Define  $P(r_1, \dots, r_i) =$  fraction of continuations that are good.

$$= \Pr(A(x; R_1, \dots, R_m) \text{ good} \mid R_1=r_1, \dots, R_i=r_i)$$

$$= \frac{1}{2} P(r_1, \dots, r_i, 0) + \frac{1}{2} P(r_1, \dots, r_i, 1)$$

$$\Rightarrow \exists r_{i+1} \in \{0,1\} \text{ s.t. } P(r_1, \dots, r_{i+1}) \geq P(r_1, \dots, r_i)$$

To find good path, just walk down tree & pick

$$r_i \in \{0,1\} \text{ for } i=1..m \text{ s.t. } P(r_1, \dots, r_{i+1}) \geq P(r_1, \dots, r_i)$$

At end:

$$P(r_1, \dots, r_m) \geq P(r_1, \dots, r_{m-1}) \geq \dots \geq P(r_1) \geq P(A(x; R_1, \dots, R_m)) \geq \frac{2}{3}$$

↑

0 or 1  $\Rightarrow$  must be 1

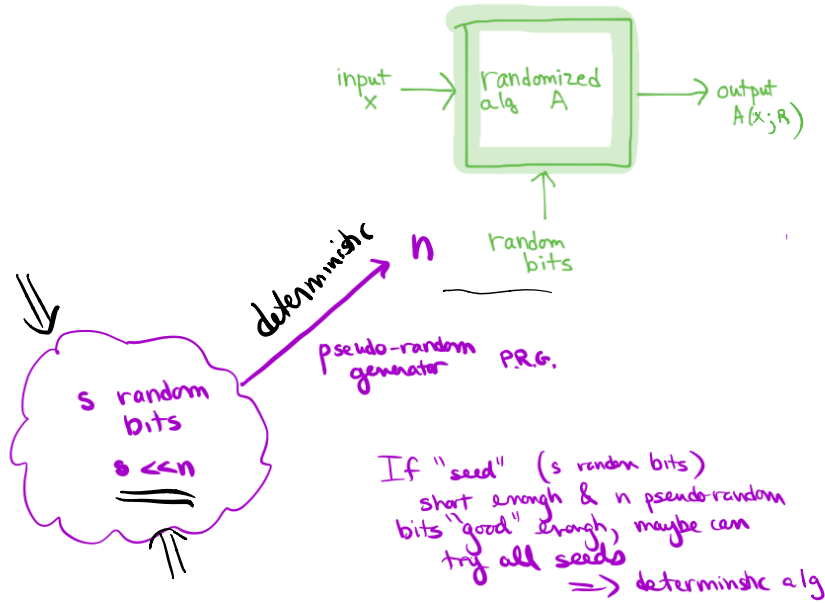
Issue: to do this: need to be able to deterministically

compute  $P(r_1, \dots, r_i)$ ; may be infeasible

but sometimes works

worked for MAX CUT

Another approach - PRGs



Back to MAXCUT

Recall  $E(|cut(S)|) = \sum_{(i,j) \in E} \Pr(R_i \neq R_j) = \frac{|E|}{2}$

random partition  
 used this =  $\frac{1}{2}$ .

$R_1, \dots, R_n$  are random bits used by algorithm

Don't need full independence of  $R_i$ 's

Pairwise independence suffices!

$\Pr(R_i \neq R_j) = \frac{1}{2}$  ✓

Observation: Suppose  $B_1, B_2, \dots, B_k$  are  $k$  indep unbiased random bits

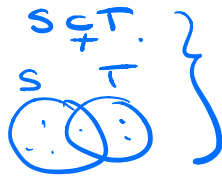
Then  $\forall S \subseteq [k]$  ( $s \neq \emptyset$ ), the  $2^s - 1$  random variables  $R_S = \bigoplus_{i \in S} B_i$  are pairwise indep unbiased random bits

$\oplus = \text{XOR}$

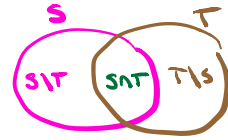
Proof: Unbiased ✓ ...

Pairwise-indep: Consider  $S \neq T \subseteq [k]$  nonempty

either  $S, T$  disjoint ✓



$$\Pr(R_T = 1 \mid R_S) = \frac{1}{2}$$



$\Rightarrow$  Given  $\lceil \lg_2(n+1) \rceil$  indep random bits  $\Rightarrow$   $n$  pairwise indep random bits

Another deterministic MAXCUT Alg

$\forall$  sequences of bits  $b_1, b_2, \dots, b_k$  where  $k = \lceil \lg_2(n+1) \rceil$

run randomized MAXCUT Alg using coin tosses  $(r_S = \bigoplus_{i \in S} b_i)_{S \neq \emptyset}$

choose largest cut obtained

Correctness:  $E(\text{cut}) = \frac{|E|}{2}$

$\Rightarrow \exists b_1, \dots, b_k$  s.t. cut has size  $\geq \frac{|E|}{2}$

Running time:



A family  $\mathcal{H}$  of fns:  $\mathcal{H} = \{h: [n] \rightarrow [m]\}$   
is pairwise indep if, when  $h$  is chosen  
v.o.a.r. from  $\mathcal{H}$  the following conditions hold:

- (1)  $\forall x \in [n]$ ,  $h(x)$  is uniform on  $[m]$
- (2)  $\forall x_1 \neq x_2 \in [n]$ ,  $h(x_1)$  and  $h(x_2)$  are independent

Super important: hashing & well beyond

Often model hash fns as truly random  
infeasible to implement

domain often exponentially large,  
can't even write it down

Want explicit family  
efficiently computable.