

①

Random variables X_1, \dots, X_n over Ω .

Mutually independent if:

$$\forall i_1, \dots, i_n \in \Omega, \Pr\left(\bigwedge_{j=1}^n X_j = i_j\right) = \prod_{j=1}^n \Pr(X_j = i_j)$$

k-wise independent if:

$\forall S \subset \{1, \dots, n\}$ where $|S| \leq k$, $S = \{s_1, \dots, s_k\}$

$$\forall i_1, \dots, i_k \in \Omega, \Pr\left(\bigwedge_{j=1}^k X_{s_j} = i_j\right) = \prod_{j=1}^k \Pr(X_{s_j} = i_j)$$

Pairwise independent is $k=2$, i.e.,

$\forall j, l \in \{1, \dots, n\}$, $a, b \in \Omega$,

$$\Pr(X_j = a, X_l = b) = \Pr(X_j = a) \times \Pr(X_l = b).$$

Lemma: Let X_1, \dots, X_n be pairwise independent,
 & let $X = \sum_{i=1}^n X_i$

Then,
$$\text{Var}(X) = \sum_{i=1}^n \text{Var}(X_i)$$

& if X_i 's are binary/indicator random variables,
 0-1

then

$$\text{Var}(X) \leq \sum_{i=1}^n E[X_i^2] = \sum_i E[X_i] = E[X].$$

Simple construction of pairwise independent random variables:

from ~~a~~ b mutually indep't. random bits,
generate $m = 2^b - 1$ pairwise indep't. random bits.

Let $X_1, \dots, X_b \in \{0, 1\}$ be uniform, mutually indep't. random bits

Let S_1, \dots, S_{2^b-1} be the nonempty subsets of $\{1, \dots, b\}$.

Set $Y_j = \bigoplus_{i \in S_j} X_i = \sum_{i \in S_j} X_i \pmod 2$

Note, $Y_j \in \{0, 1\}$

Lemma: The Y_j 's are pairwise indep't.

Proof: Claim 1: $\Pr(Y_j = 1) = \Pr(Y_j = 0) = \frac{1}{2}$

Why? Let $S_j = \{z_1, \dots, z_{\ell}\} \subseteq \{1, \dots, b\}$

So, $Y_j = \left(\sum_{i=1}^{\ell-1} X_{z_i} \pmod 2 \right) + X_{z_\ell} \pmod 2$

Reveal $X_{z_1}, \dots, X_{z_{\ell-1}}$. Whatever this is,

with prob. $\frac{1}{2}$ $X_{z_\ell} = 1$ & Y_j is opposite

& w.p. $\frac{1}{2}$ $X_{z_\ell} = 0$ & Y_j is the same.

This is the principle of deferred decisions.

Now to see pairwise independence:

Fix S_j & S_l and consider some $z \in S_j \setminus S_l$.

$$\begin{aligned} \Pr(Y_j=a, Y_l=b) &= \Pr(Y_j=a | Y_l=b) \Pr(Y_l=b) \\ &= \Pr(Y_j=a | Y_l=b) \cancel{\Pr(Y_l=b)} \times \frac{1}{2} \end{aligned}$$

we just showed this.

Reveal all X_i 's but X_z

Then with prob. $\frac{1}{2}$ $X_z=1$ & Y_j flips
& w.p. $\frac{1}{2}$ $X_z=0$ & Y_j is the same

$$\begin{aligned} \text{Therefore, } \Pr(Y_j=a | Y_l=b) &= \Pr(Y_j=a | X_1, \dots, X_b \setminus X_z) \\ &= \frac{1}{2} \end{aligned}$$

& thus, $\Pr(Y_j=a, Y_l=b) = \frac{1}{4}$.

More sophisticated construction:

For prime P , given a, b which are indep't. & uniform over $\{0, 1, \dots, P-1\}$

then we construct Y_1, \dots, Y_{P-1} which are pairwise indep't. & uniform over $\{0, 1, \dots, P-1\}$.

Namely, let $Y_i = a + ib \pmod P$ for $i=0, \dots, P-1$.

Lemma: The Y_i 's are pairwise indep't.

Proof: First, Y_i is uniform over $\{0, 1, \dots, P-1\}$.

Why? By principle of deferred decisions again.

For any b & i & α in $\{0, \dots, P-1\}$,

$$\Pr(Y_i = \alpha) = \Pr(a + bi \equiv \alpha \pmod P)$$

$$= \Pr(a \equiv \alpha - bi \pmod P)$$

$$= \frac{1}{P} \quad \text{Since there is a unique such } a \text{ in } \{0, 1, \dots, P-1\}.$$

Now consider $i, j \in \{0, \dots, P-1\}$ and $\alpha, \beta \in \{0, \dots, P-1\}$
we'll show: $\Pr(Y_i = \alpha, Y_j = \beta) = \frac{1}{P^2}$ & we're done.

$$Y_i = \alpha \Leftrightarrow a + ib \equiv \alpha \pmod{P}$$

$$Y_j = \beta \Leftrightarrow a + jb \equiv \beta \pmod{P}$$

$$\text{Thus, } \alpha - \beta \equiv b(i - j) \pmod{P}$$

$$b \equiv \frac{\alpha - \beta}{i - j} \pmod{P}$$

which is valid since $i - j \neq 0$
& P is prime

$$\& a \equiv \alpha - bi \pmod{P}$$

So there is a unique (a, b) pair
so that $Y_i = \alpha, Y_j = \beta$

$$\text{Therefore, } \Pr(Y_i = \alpha, Y_j = \beta)$$

$$= \Pr\left(b \equiv \frac{\alpha - \beta}{i - j} \pmod{P}, a \equiv \alpha - i \left(\frac{\alpha - \beta}{i - j}\right) \pmod{P}\right)$$

$$= \frac{1}{P^2}$$

□

For n which is not prime, can choose $P > n$
where P is prime & $P < 2n$.

Note, the random variables take $O(\log n)$ bits to represent a, b .

Streaming model:

Stream $S = \{s_1, \dots, s_m\}$ for HUGE m .

where each $s_i \in \{0, 1, \dots, n-1\}$

Let $f = (f_1, \dots, f_n)$ where

$$f_i = |\{j : s_j = i\}| = \text{frequency \# of occurrences of value } i \text{ in } S.$$

Let $Q = F_0 = |\{i : f_i > 0\}| = \# \text{ of distinct values in } S$

Goal: find Q with $O(\log n)$ space.

aim for (ϵ, δ) -approx of Q :

$$\text{output } \hat{Q} \text{ where } \Pr((1-\epsilon)\hat{Q} \leq Q \leq (1+\epsilon)\hat{Q}) \geq 1-\delta$$

in ~~time~~ poly space $\text{Poly}(\log n, \frac{1}{\epsilon}, \log(\frac{1}{\delta}))$

For integer $k > 0$,

let $\text{zeros}(k) = \#$ of 0's at end of binary representation of k .

$$= \max \{ l : 2^l \text{ divides } k \}$$

e.g., if k is odd then $\text{zeros}(k) = 0$
 even then $\text{zeros}(k) \geq 1$ since it ends in 0.

AMS algorithm:

(Diff. alg. than last class but same Paper)

1. Choose a random hash function $h: [n] \rightarrow [n]$
 $\{0, 1, \dots, n-1\} \rightarrow \{0, 1, \dots, n-1\}$
 which is pairwise independent

(if n is not prime, choose prime p where $n \leq p \leq 2n$ & think of n as p)

2. Set $z = 0$

3. Go through stream & at element k :
 if $\text{zeros}(h(k)) > z$
 then $z = \text{zeros}(h(k))$

4. output $(2^{z + \frac{1}{2}})$

Intuition for alg.:

9

$h(k)$ is uniformly random bitstring so prob. it has $\text{zeros}(h(k)) = l$ is prob. of last l bits all = 0 which is prob. 2^{-l} .

Thus, prob. that $\text{zeros}(h(k)) = \log Q$ is $2^{-\log Q} = \frac{1}{Q}$

So if Q distinct items then expect 1 out of these Q to have $\text{zeros}(h(k)) = \log Q$.

& for $l \gg \log Q$ it's unlikely that $\text{zeros}(h(k)) \gg \log Q$.

Thus, the $\max_k \text{zeros}(h(k))$ is a good approx. for $\log Q$.

Analysis:

For $k \in \{0, \dots, n-1\}$, & integer $l \geq 0$,

let $X_{l,k} = \begin{cases} 1 & \text{if } \text{zeros}(h(k)) \geq l \\ 0 & \text{o/w} \end{cases}$

& let $Y_l = \sum_{k: f_k > 0} X_{l,k}$

Let t be the final value of z at end of algorithm.

$$t \geq l \iff Y_l > 0$$

$t = \max_{\#} \{l : Y_l > 0\} = \max \{l : \exists k, f_k > 0, \text{zeros}(h(k)) \geq l\}$
 this is same as:

$$t < l \iff Y_l = 0.$$

" "

$$t \leq l - 1$$

Note, $h(k)$ is uniform, i.e., $\Pr(h(k)=j) = \frac{1}{n}$ for all $j \in \{0, \dots, n-1\}$

$$\begin{aligned}
 \text{Thus, } E[X_{l,k}] &= \Pr(\text{zeros}(h(k)) \geq l) \\
 &= \Pr(\text{last } l \text{ bits of } h(k) \text{ are all } 0) \\
 &= \frac{1}{2^l}
 \end{aligned}$$

$$\text{Var}(Y_\ell) = \sum_{k: f_k > 0} \text{Var}(X_{\ell,k}) \quad \text{since } X_{\ell,k} \text{ \& } X_{\ell,k'} \text{ are pairwise indep.}$$

$$\leq \sum_{k: f_k > 0} E[X_{\ell,k}^2]$$

$$= \sum_{k: f_k > 0} E[X_{\ell,k}]$$

$$= \frac{d}{2^\ell}$$

By Markov's ineq.,

$$\begin{aligned} \Pr(Y_\ell > 0) &= \Pr(Y_\ell \geq 1) \\ &\leq \frac{E[Y_\ell]}{1} = \frac{d}{2^\ell} \end{aligned}$$

By Chebyshev's (can apply since pairwise indep't.)

$$\begin{aligned} \Pr(Y_\ell = 0) &\leq \Pr(|Y_\ell - E[Y_\ell]| \geq \frac{d}{2^\ell}) \\ &\leq \frac{\text{Var}(Y_\ell)}{\left(\frac{d}{2^\ell}\right)^2} \leq \frac{d}{2^\ell} \end{aligned}$$

Our goal is to output Q .

Let's aim for \hat{Q} where $\frac{\hat{Q}}{3} \leq Q \leq 3\hat{Q}$

So \hat{Q} is a 3-approx. of Q .

Let a be the smallest integer s.t. $2^{a+\frac{1}{2}} \geq 3Q$.

We want to show that the prob. that $t = \text{final value of } z$ is unlikely to be as large as a .

$$\Pr(\hat{Q} \geq 3Q) = \Pr(t \geq a)$$

$$= \Pr(Y_a > 0)$$

$$\leq \frac{Q}{2^a} = \frac{\sqrt{2}}{3} < \cancel{.48} < \frac{1}{2}$$

Note, $2^{a+\frac{1}{2}} \geq 3Q$

$$\frac{Q}{2^a} \leq \frac{\sqrt{2}}{3}$$

Thus, $\Pr(\hat{Q} < 3Q) > .51$

On the other side, we want $\hat{Q} \geq \frac{Q}{3}$

Let b be the largest integer s.t. $2^{b+\frac{1}{2}} \leq \frac{Q}{3}$

$$\Pr(\hat{Q} \leq \frac{Q}{3}) = \Pr(t \leq b)$$

$$= \Pr(Y_{b+1} = 0)$$

$$\leq \frac{2^{b+1}}{Q} = \left(\frac{2^{b+\frac{1}{2}}}{Q}\right) \sqrt{2} \leq \frac{\sqrt{2}}{3} < .48$$

Therefore, $\Pr\left(\frac{Q}{3} < Q < 3\hat{Q}\right) \geq .04$

Since $\Pr(\hat{Q} \leq \frac{Q}{3} \text{ or } \hat{Q} \geq 3Q) \leq 2 \times .48 = .96$

How to boost this prob. to $\geq 1-\delta$?

D. $k = O(\log(1/\delta))$ trials,
get outputs D_1, \dots, D_k

Output $D = \text{Median}(D_1, \dots, D_k)$

$$\text{Let } Z_i = \begin{cases} 1 & \text{if } D_i < 3Q \\ 0 & \text{o/w} \end{cases}$$

If D is $\geq 3Q$ then $\geq \frac{k}{2}$ of the trials
exceed $\geq 3Q$.

$$\text{Let } Z = \sum_{i=1}^k Z_i$$

$$E[Z] \geq .52k$$

$$\Pr(D \geq 3Q) \leq \Pr(Z < \frac{k}{2})$$

$$= \Pr(Z \leq \frac{k}{2} - .02k)$$

$$\leq \Pr(Z \leq (1 - .02)E[Z])$$

$$\leq e^{-\frac{.02 \cdot .52k}{3}} \leq \frac{\delta}{2}$$

for $k = c \log(1/\delta)$ with c
big enough
constant.

This gives a $(3, \delta)$ -approx.

13

Using $O(\log n)$ bits per hash function

& $O(\log \log n)$ bits for z

$\Rightarrow O(\log(\frac{1}{\delta}) \log n)$ total bits.

Next class: $(1+\epsilon, \delta)$ -approx. for all $\epsilon > 0$.