

## Lecture 3 :- Concentration and saving memory (II)

So far :-  $X_i = I[G_i \text{ is "bad"}]$

$M := \text{Median of } (G_1, \dots, G_k)$ ,  $G_i$  independent.

$$E[X_i] = p \leq \frac{1}{8}, \text{ and}$$

$$\begin{aligned} \Pr[|M - F_2| \geq \varepsilon F_2] &\leq \Pr\left[\sum_{i=1}^k X_i \geq \frac{k}{2}\right] \\ &\leq \Pr\left[\sum_{i=1}^k (X_i - p) \geq k\left(\frac{1}{2} - p\right)\right] \\ &\leq \exp(-2k \cdot (\frac{1}{2} - p)^2) \end{aligned}$$

$$\leq \exp\left(-\frac{9k}{32}\right) \quad \text{From the Chernoff-Hoeffding bound discussed above.}$$

We wanted  $\Pr[|M - F_2| > \varepsilon F_2] \leq \delta$ .

From  $\times$  this requires  $k \geq \lceil \frac{32}{9} \log\left(\frac{1}{\delta}\right) \rceil$ .

What we gained :-

Only  $O(\log(1/\delta))$  samples of the  $G_i$  are needed to drive the error probability to below  $\delta$ !

Summing up :- Fix  $\varepsilon > 0$  (accuracy parameter)  
 $\delta > 0$  (error probability tolerance).

1. If our stream is over  $D$  dimensions:

$$Y_{ij} = \left( \sum_{l=1}^D f_l u_{l,(i,j)} \right)^2 \quad [f_l = \text{Frequency of the } l\text{-th element in the stream.}]$$

Requirement:- For a fixed  $(i, j)$  pair.

the collection  $(u_{l,(i,j)})_{l=1}^D$  of  $l \in \{+1, -1\}$  r.v.s with expectation 0 is 4-wise independent.

2.

$$G_i := \frac{1}{s} \sum_{j=1}^s Y_{i,j} \quad s \geq \left\lceil \frac{16}{\varepsilon^2} \right\rceil.$$

We needed  $(Y_{i,j})_{j=1}^s$  (fixed  $i$ ) to be at least pairwise independent. — ①

3.  $M = \text{Median}(G_1, \dots, G_k), k \geq \lceil \frac{32}{\varepsilon^2} \log(\frac{1}{\delta}) \rceil$

We need  $G_1, \dots, G_k$  to be independent. — ②

Both ① and ② are satisfied if we assume the vectors  $(u_{l,(i,j)})_{l=1}^D$ , as  $(i, j)$  vary, are independent  $[O(\frac{1}{\varepsilon^2} \log(\frac{1}{\delta}))$  vectors]

Thus:-

We need storage for  $sk = O(\frac{1}{\varepsilon^2} \log(\frac{1}{\delta}))$  counters.  
 $Y_{ij}$ .

But we make all the  $u_{l,(i,j)}$  independent we will

have to store  $\Omega(D \cdot \frac{1}{\epsilon^2} \log(\frac{1}{\delta}))$  random bits.

But the whole point of this exercise was to not need storage that grows linearly in terms of the alphabet size of the stream.

- With  $\Omega(D)$  storage, we might just use an associative array (HashMap / Balanced binary tree map).

So, how do we get rid of this linear dependence on D?

We need to use that each of the  $s_k$  arrays only needs to have  $k$ -wise independent (and not necessarily fully independent!!) signs.

~~X      X      X      X~~  
Limited independence:- We only want to store a small number of truly random bits. (called the 'seed') such that a deterministic algorithm working with this seed can produce a much larger sequence of bits, which look (say) pairwise independent to an "outside observer" who has not seen the seed.

A construction:- Store as seed  $s = s_1, \dots, s_n$  sampled uniformly at random from  $\{0, 1\}^{n+1}$ .

I do

Algorithm :- When the observer asks for bit indexed  $i$  where  $0 \leq i \leq 2^n - 1$ , the algorithm outputs

$$A(s, i) := s + \bigoplus_{j=1}^n (s_j \cdot I[j\text{th bit of } i \text{ is } 1])$$

Claim :-  $\forall i, j, i \neq j, 0 \leq i, j \leq 2^n - 1, \forall \alpha, \beta \in \{0, 1\}$ ,

$$\begin{aligned} & \Pr_{\substack{s \sim \{0,1\}^n \\ \text{unif.}}} [A(s, i) = \alpha \text{ and } A(s, j) = \beta] \\ &= \Pr_{\substack{s \sim \{0,1\}^n \\ \text{unif.}}} [A(s, i) = \alpha] \cdot \Pr_{\substack{s \sim \{0,1\}^n \\ \text{unif.}}} [A(s, j) = \beta] \\ &= \frac{1}{4}. \quad [\text{Chor-Goldreich}] \end{aligned}$$

Second construction:  $\xrightarrow{\text{k-wise}} \xrightarrow{\text{independence}} \xrightarrow{[\text{BCH, Reed-Solomon}]}$

Fix a finite field  $\mathbb{F}$  :- A field is just a set with a 0 element, a 1 element (different from the 0 element) commutative addition and multiplication, with multiplication distributing over addition, in which every element (except 0) has a unique multiplicative inverse and every element has an additive inverse.

Examples :-  $\mathbb{C}, \mathbb{R}, \mathbb{Q}$ .

Non-examples :-  $\mathbb{Z}$ .

A field with finitely many elements is a finite field. [Any finite field has size

$p^d$  for some  $p$  prime (including 2) and  $d$  a positive integer. Further, there is exactly one finite field of each such size.]

One can define polynomials over a finite field  $\mathbb{F}$

$$p(x) = \sum_{i=0}^d \alpha_i x^i, \quad \alpha_i \in \mathbb{F}.$$

Fact:-

- If  $p, q$  are polynomials of degree  $d$  over  $\mathbb{F}$ , and  $p(\beta_i) = q(\beta_i)$  for at least  $d+1$  distinct  $\beta_i \in \mathbb{F}$  then  $p = q$ .

- If  $r_0, \dots, r_d \in \mathbb{F}$  are distinct values and  $\beta_0, \dots, \beta_d \in \mathbb{F}$  are arbitrary values then there always exist a polynomial  $p$  s.t.  $p(r_i) = \beta_i, 0 \leq i \leq d$ .

This polynomial is unique by the previous fact.

Construction  $\xrightarrow{\quad} \xrightarrow{\quad} \xrightarrow{\quad} \xrightarrow{\quad} \xrightarrow{\quad}$

$k$ -wise independent r.v. taking values in  $\mathbb{F}$ .  $k \leq \log |\mathbb{F}|$

- Pick  $\alpha_0, \dots, \alpha_{k-1}$  u.a.r. from  $\mathbb{F}$  ["Seed"].

- Outside observers inputs  $y \in \mathbb{F}$ . The algorithm outputs:

$$A(\vec{x}, y) = \sum_{i=0}^{k-1} \alpha_i y^i$$

Claim :- Let  $r_1, r_2, \dots, r_k \in \mathbb{F}$  be distinct, and  $\beta_1, \beta_2, \dots, \beta_k \in \mathbb{F}$  arbitrary. Then.

$$\Pr_{\vec{r} \in \mathbb{F}^k} \left[ \bigwedge_{j=1}^k [A(\vec{r}, r_j) = \beta_j] \right] = \prod_{j=1}^k \Pr_{r \in \mathbb{F}} [A(\vec{r}, r_j) = \beta_j] \\ = \frac{1}{|\mathbb{F}|^k}.$$

- Just the previous fact in disguise !!

(There is a unique  $\vec{r} \in \mathbb{F}^k$  for which

$$A(\vec{r}, r_j) = \beta_j \text{ for } k \text{ distinct } r_j.)$$

So when  $\vec{r} \in \mathbb{F}^k$  is chosen w.r.t. the random variables  $(A(\vec{r}, r))_{r \in \mathbb{F}}$  are  $\mathbb{F}$ -valued  $k$ -wise independent random variable.

Requirement :-

D random variables

4-wise independent

$\{+1, -1\}$ -valued, with uniform distribution

An option : choose  $\mathbb{F} = \mathbb{F}_p$  where p is the smallest prime satisfying  $p \geq D$ .

Next option :- choose d smallest s.t.  $2^d \geq D$  and take  $\mathbb{F} = \mathbb{F}_{2^d}$ .