

## Frequency Moments

$f_i$ : # times  $i$  occurs in stream

$$F_t = \sum f_i^t$$

$F_0$  = # distinct elements  $0^0 = 1$

$F_1$  = size of stream

$F_2$  = measure of skewness

[Alon, Matias, Szegedy '96]

$$h: U \rightarrow \{\pm 1\}$$

$$Y = \sum_i h(x_i)$$

Estimator:  $Y^2$

$$X_i = h(x_i) \quad Y = \sum_i f_i X_i$$

$$\begin{aligned} E[Y^2] &= E\left[\left(\sum_i f_i X_i\right)^2\right] \\ &= E\left[\sum_{i,j} f_i f_j X_i X_j\right] \\ &= \sum_i f_i^2 + E\left[\sum_{i \neq j} f_i f_j X_i X_j\right] \end{aligned}$$

$$i \neq j \quad E[X_i X_j] = 0$$

$$E[Y^2] = \sum f_i^2$$

$$\begin{aligned} E[Y^4] &= E\left[\left(\sum_i f_i X_i\right)^4\right] \\ &= E\left[\sum_{i_1, i_2, i_3, i_4} f_{i_1} f_{i_2} f_{i_3} f_{i_4} X_{i_1} X_{i_2} X_{i_3} X_{i_4}\right] \\ &= \sum_i f_i^4 + 6 \sum_{i \neq j} f_i^2 f_j^2 \end{aligned}$$

$$\left(\sum f_i X_i\right) \left(\sum f_i X_i\right) \left(\sum f_i X_i\right) \left(\sum f_i X_i\right)$$

$$\binom{4}{2}$$

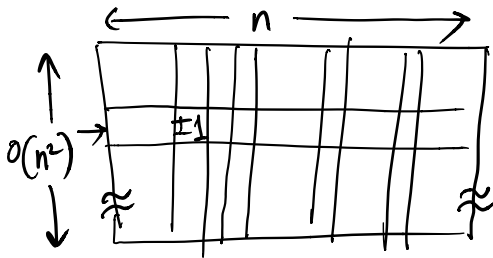
$$\begin{aligned} \text{Var}[Y^2] &= E[Y^4] - (E[Y^2])^2 & E[Y^2] &= \sum_i f_i^2 \\ &= \left( \sum_i f_i^4 + 6 \sum_{i < j} f_i^2 f_j^2 \right) - \left( \sum_i f_i^4 + 2 \sum_{i < j} f_i^2 f_j^2 \right) \\ &= 4 \sum_{i < j} f_i^2 f_j^2 \leq 2 F_2^2 = 2 (E[Y^2])^2 \end{aligned}$$

$O(\frac{1}{\epsilon^2})$  copies and take mean get  $(1 \pm \epsilon)$  approx to  $F_2$   
w. prob.  $\frac{2}{3}$

$O(\log(\frac{1}{\epsilon}))$  groups and compute median of means  
→ prob. of success  $1 - \delta$

$(1 \pm \epsilon)$  approximation to  $F_2$  with prob.  $(1 - \delta)$   
 $O(\frac{1}{\epsilon^2} \log(\frac{1}{\epsilon}))$  space

$$\begin{aligned} E[X_{i_1}] &= 1 & E[X_{i_1} X_{i_2}] &= 0 \\ E[X_{i_1} X_{i_2} X_{i_3} X_{i_4}] &= 0 \end{aligned}$$



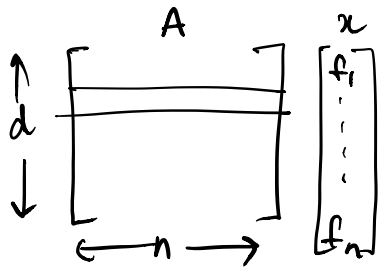
Any 4 columns  $i_1, i_2, i_3, i_4$   
 $\{\pm 1\}^4 \quad +1, -1, -1, +1$

Orthogonal arrays of strength 4  
parity check matrices of BCH codes

seed  $\equiv$  one of  $O(n^2)$  rows  $O(\log n)$

Any entry of matrix can be computed in  $O(\log n)$  space  
&  $O(\log n)$  time

$$\begin{aligned} Y_1, Y_2, \dots, Y_d & \quad d = O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right) \\ x \in \mathbb{R}^n & \quad x = (f_1, f_2, \dots, f_n) \end{aligned}$$



Algorithm stores  $Ax$

$$x + \Delta x \quad \Delta x = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$$

$$Ax + A\Delta x$$

linear sketch  $Ax + Ay = A(x+y)$

furnstake model: input is  $x$

updates to co-ordinates of  $x$  (need not be  $\pm 1$ )

Linear sketches can be used in furnstake model.

$$\underbrace{A(x+\Delta x)}_{\text{new sketch}} = \underbrace{Ax}_{\text{old sketch}} + \underbrace{A\Delta x}_{\text{update}}$$

Johnson-Lindenstrauss Lemma

Lemma:  $G \subset (\mathbb{R}^d, \ell_2)$  be a set of  $n$  points.

For any  $0 < \epsilon < \frac{1}{2}$   $k = O(\log n / \epsilon^2)$

$\exists$  mapping  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$

such that for all  $\sigma_i, \sigma_j \in G$

$$(1-\epsilon) \|\sigma_i - \sigma_j\| \leq \|f(\sigma_i) - f(\sigma_j)\|_2 \leq (1+\epsilon) \|\sigma_i - \sigma_j\|_2$$

linear mapping:  $f(v) = \frac{Mv}{\sqrt{k}}$

$M \in \mathbb{R}^{k \times d}$   $M_{ij} = \mathcal{N}(0, 1)$  data-oblivious

Use tail bound for chi-squared distributions.

Lemma: Let  $Z_1 \dots Z_k$  be iid. unit normal random variables.

$$\text{Let } Y = \sum_i Z_i^2$$

Then  $\Pr[(1-\epsilon)^2 k \leq Y \leq (1+\epsilon)^2 k] \geq 1 - 2e^{-c\epsilon^2 k}$   
for some suitable constant  $c$

[Pf idea]  $\Pr\left[\frac{Y}{k} \geq 1+\epsilon\right] \leq \Pr\left[e^{tY} \geq e^{tk(1+\epsilon)}\right]$

$$\leq \frac{\mathbb{E}[e^{tY}]}{e^{tk(1+\epsilon)}}$$

$$= \prod_i \frac{\mathbb{E}[e^{tZ_i^2}]}{e^{t(1+\epsilon)}}$$

(Crux): Bounding  $\mathbb{E}[e^{tZ_i^2}]$  where  $Z_i$  is unit normal.

$$f(\sigma) = \frac{M\sigma}{\sqrt{k}} \quad M \text{ is } k \times d \text{ matrix}$$

$$(M\sigma)_i = M_i^T \sigma = \sum_{j=1}^d M_{ij} \sigma_j = \left(\sum_{j=1}^d \sigma_j^2\right)^{1/2} Y = Y$$

unit vector  $\sigma \in \mathbb{R}^d$  ↑ unit normal

$$\alpha_1 X_1 + \alpha_2 X_2 = (\alpha_1^2 + \alpha_2^2)^{1/2} X$$

$$\Pr\left[k(1-\epsilon)^2 \leq \sum_{i=1}^k (M\sigma)_i^2 \leq k(1+\epsilon)^2\right] \geq 1 - 2e^{-c\epsilon^2 k}$$

$$\Pr\left[(1-\epsilon) \leq \left\|\frac{M\sigma}{\sqrt{k}}\right\|_2 \leq (1+\epsilon)\right] \geq 1 - 2e^{-c\epsilon^2 k}$$

$$\sigma \quad \tilde{\sigma} = \frac{\sigma}{\|\sigma\|_2}$$

$$\Pr\left[(1-\epsilon) \leq \left\|\frac{M\tilde{\sigma}}{\sqrt{k}}\right\|_2 \leq (1+\epsilon)\right] \geq 1 - 2e^{-c\epsilon^2 k}$$

$$\Pr \left[ (1-\epsilon) \|v\|_2 \leq \left\| \frac{Mv}{\sqrt{K}} \right\|_2 \leq (1+\epsilon) \|v\|_2 \right] \geq 1 - 2e^{-c\epsilon^2 K}$$

$$v = v_i - v_j \quad \binom{n}{2} \text{ such pairs.}$$

$$\Pr \left[ (1-\epsilon) \|v_i - v_j\|_2 \leq \underbrace{\left\| \frac{M(v_i - v_j)}{\sqrt{K}} \right\|_2}_{\|f(v_i) - f(v_j)\|_2} \leq (1+\epsilon) \|v_i - v_j\|_2 \right] \geq 1 - 2e^{-c\epsilon^2 K}$$

$$f(v_i) = \frac{Mv_i}{\sqrt{K}} \quad \|f(v_i) - f(v_j)\|_2 = \left\| \frac{M(v_i - v_j)}{\sqrt{K}} \right\|_2$$

$$\text{Choosing } K = \frac{c' \log(n)}{\epsilon^2}$$

$$\text{Failure prob: } 2e^{-cc' \log n} < \frac{1}{n^3}$$

$$\text{Overall failure prob.} \leq \frac{1}{n} \quad (\text{union bound})$$

$$[\text{Alon '00}] \quad \Omega\left(\frac{\log n}{\epsilon^2 \log(\frac{1}{\epsilon})}\right) \text{ dimensions}$$

$$[\text{Larsen, Nelson '17}] \quad \Omega\left(\frac{\log n}{\epsilon^2}\right) \text{ dimensions}$$

$$\text{Single pair: } O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\epsilon}\right)\right) \text{ dimensions}$$

$l_2$  is special

$l_p$   $p = \infty$  dimension red<sup>n</sup> not possible.

roughly, to preserve distances to factor  $\alpha$   
need  $n^{1/\alpha}$  dimensions

$p=1$  dimension red<sup>n</sup> not possible.

to achieve factor  $\alpha$ , need  $n^{1/\alpha^2}$  dimensions  
[Brinkman, C, '03]

$\frac{n}{\alpha}$  upper bound