

Plan:

Frequency Moments F_k $k \in (0, 2]$

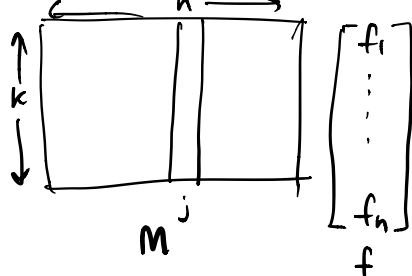
Lower Bound l_2 dimension reduction

Heavy Hitters

F_k $k \in (0, 2]$

p -stable Distribution D_p

$k \times n$ matrix M $M_{ij} \sim D_p$ $k = O(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$



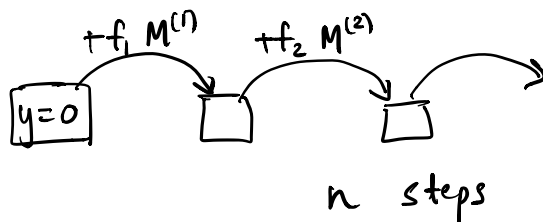
$$y = Mf$$

$$y_i = \sum_{j=1}^n M_{ij} f_j$$

$$y_i \sim (\sum |f_j|^p)^{1/p} D_p$$

$$y := y + M^{(j)} \text{ when } j \text{ appears}$$

$$\text{Space} = \underbrace{O(\frac{1}{\epsilon^2} \log(\frac{1}{\delta}) \log n)}_S + \text{space for random matrix}$$



$$\square \quad y = Mf$$

S random bits in each step $R = n$ steps

$$S = O(\frac{1}{\epsilon^2} \log(\frac{1}{\delta}) \log n)$$

U_t : uniform random string in $\{0, 1\}^t$

Nisan's pseudorandom generator:

$$\exists h: \{0, 1\}^{S \log R} \rightarrow \{0, 1\}^{SR}$$

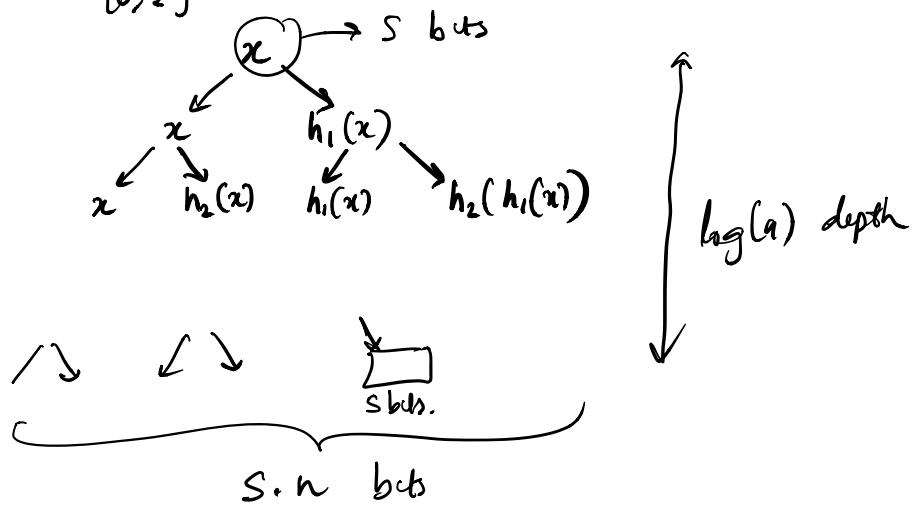
$$|\Pr(f(U_{SR}) = 1) - \Pr(f(h(U_{S \log R})) = 1)| \leq 2^{-\Omega(S)}$$

f : results of executing space S decision algorithm for R steps
with S random bits per step

$$|\text{Seed}| = S \log R = O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \log^2 n\right)$$

$h_1 \dots h_{\log n}$ pairwise independent hash fns.
 $h_i: [2^S] \rightarrow [2^S]$

choose $x \in \{0, 1\}^S$



$$F_i \quad f_i \quad g_i \quad \sum |f_i - g_i|$$

$$Mf \quad Mg \quad M(f-g) \rightarrow \text{estimate } \|f-g\|_1$$

lower bound for dimension red^n in l_1 $n^{1/\epsilon}$ dimensions
 $(1+\epsilon)$ approxⁿ in space $O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$ with prob $1-\delta$
 $O\left(\frac{1}{\epsilon^2} \log n\right)$

$\text{Dim}^n \text{red}^n$: treat sketch as mapped into l_1 small dimension
compute l_1 norm of sketch

Lower Bound l_2 dimension reduction

Thm: [Alon '00] Let $v_1 \dots v_{n+1} \in \mathbb{R}^d$ $\frac{1}{\sqrt{n}} \leq \epsilon \leq \frac{1}{3}$

$$1 \leq \|v_i - v_j\| \leq 1 + \epsilon \quad \forall i \neq j \in [n+1]$$

Then subspace spanned by $v_1 \dots v_{n+1}$ has dimension

$$d = \Omega\left(\frac{\log n}{\epsilon^2 \log(1/\epsilon)}\right)$$

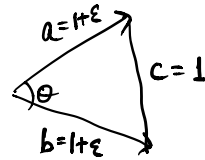
n dimensional simplex cannot be embedded with distortion $1 + \epsilon$ in fewer than $\Omega(\log n)$ dimensions

Assume $v_{n+1} = 0$. So $1 \leq \|v_i\| \leq 1 + \epsilon$

$$\text{Set } v_i' = \frac{v_i}{\|v_i\|}$$

$$\begin{aligned} \langle v_i', v_j' \rangle &= \cos \angle(v_i, v_j) \\ &\leq \frac{1}{2} + \epsilon + \frac{\epsilon^2}{2} \end{aligned}$$

$$\left| \langle v_i', v_j' \rangle - \frac{1}{2} \right| = O(\epsilon)$$



$$c^2 = a^2 + b^2 - 2ab \cos \theta$$

$$\cos \theta = \frac{a^2 + b^2 - c^2}{2ab} = \frac{(1+\epsilon)^2 + (1+\epsilon)^2 - 1}{2(1+\epsilon)^2}$$

Define $n \times n$ matrix B $B_{ij} = \langle v_i', v_j' \rangle$

$$\begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & & \frac{1}{2} + O(\epsilon) & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix}$$

If $\epsilon = 0$, B has rank n

$\Rightarrow v_i$'s span subspace of dim. n

No distortion \Rightarrow dim. $\geq n$

$d = \text{rank}(B)$ Define $C = 2B - J$, $J = ee^T$ (all ones)

$$|\text{rank}(C) - \text{rank}(B)| \leq 1$$

$$?? \leq \text{rank}(C) \leq d + 1$$

$$C = \begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & & 0(\epsilon) & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix}$$

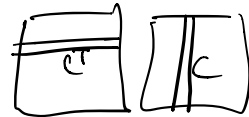
Lemma Consider symmetric matrix C $C_{ii} = 1$
 $|C_{ij}| \leq \frac{1}{\sqrt{n}} \quad i \neq j$
 then $\text{rank}(C) \geq \frac{n}{2}$

Proof: C symmetric \Rightarrow all eigenvalues real
 $d = \text{rank}(C) \quad \lambda_1 \dots \lambda_d \in \mathbb{R}$ non-zero eigenvalues

$$\text{Tr}(C) = \sum_{i \in [n]} C_{ii} = n = \sum_{i \in [d]} \lambda_i$$

Non-zero eigenvalues of $C^2 = C^T C$ are $\lambda_1^2 \dots \lambda_d^2$

$$\text{Tr}(C^2) = \sum_{i \in [d]} \lambda_i^2$$



$$\text{Tr}(C^2) = \sum_i \sum_j C_{ij}^2 \leq n + n(n-1) \frac{1}{n} = 2n - 1 < 2n$$

$$\frac{2n}{d} > \frac{\sum \lambda_i^2}{d} \geq \left(\frac{\sum \lambda_i}{d} \right)^2 = \left(\frac{n}{d} \right)^2$$

$$\Rightarrow d > \frac{n}{2}$$

Lemma: Suppose $n \times n$ matrix A has rank d
 then $F = (A_{ij}^k)$ has rank at most $\begin{pmatrix} d+k-1 \\ d-1 \end{pmatrix}$

Proof Let $v_1 \dots v_d \in \mathbb{R}^n$ be basis for row space of A

$$i\text{th row } A_i = \sum_{l \in [d]} \lambda_l v_l \quad \text{for some coeff } \lambda_l$$

$$A_{ij} = \sum_{l \in [d]} \lambda_l v_{l,j}$$

$$F_{ij} = A_{ij}^k = \left(\sum_{l \in [d]} \lambda_l v_{l,j} \right)^k$$

$$= \sum_{k_1 + \dots + k_d = k} \underbrace{\binom{k}{k_1 \dots k_d}}_{\text{coefficient}} \underbrace{\left(\prod_{l \in [d]} \lambda_l^{k_l} \right)}_{\text{coeff}} \underbrace{\left(\prod_{l \in [d]} v_{l,j}^{k_l} \right)}_{j\text{th coord of basis}}$$

row-space of F spanned by

$$(W_{k_1 \dots k_d})_j = \prod_{\ell \in [d]} U_{\ell, j}^{k_\ell}$$

one vector for every choice of k_1, \dots, k_d $\sum k_\ell = K$

$$\# \text{ basis vectors} = \# \text{ partitions} = \binom{k+d-1}{d-1} = \binom{k+d-1}{k}$$

Proof of Thm ?? $\leq \text{rank}(C) \leq d+1$ $|C_{ij}| \leq O(\epsilon)$

$$k \text{ integer } \epsilon^k \leq \frac{1}{\sqrt{n}}$$

$$\text{Consider } F = (C_{ij}^k)_{1 \leq i, j \leq n}$$

$$|F_{ij}| \leq \frac{1}{\sqrt{n}}$$

$$\text{rank}(F) \leq \binom{k+d}{d} \quad (\text{lemma 2})$$

$$\text{rank}(F) \geq \frac{n}{2}$$

$$\frac{n}{2} \leq \binom{k+d}{d} = \binom{k+d}{k} = \frac{(k+d)!}{d! \cdot k!} \leq \binom{k+d}{k} \left(\frac{e}{k}\right)^k$$

$$\text{take logs of both sides } k = \frac{\ln n}{2 \ln(\frac{1}{\epsilon})}$$

$$k \ln\left(\frac{e(d+k)}{k}\right) \geq \ln\left(\frac{n}{2}\right)$$

$$\frac{\ln n}{2 \ln(\frac{1}{\epsilon})} \ln\left(\frac{e(d+k)}{k}\right) \geq \ln\left(\frac{n}{2}\right)$$

$$\ln\left(\frac{e(d+k)}{k}\right) \geq (1-o(1)) 2 \ln\left(\frac{1}{\epsilon}\right)$$

$$\frac{e(d+k)}{k} \geq (1-o(1)) \frac{1}{\epsilon^2}$$

$$\frac{d}{k} \geq (1-o(1)) \frac{1}{\epsilon^2} - 1$$

$$d \gg \Omega\left(\frac{\ln n}{\epsilon^2 \ln(1/\epsilon)}\right)$$

$$\frac{\ln(n/2)}{\ln n} = \frac{\ln(n) - c}{\ln(n)} = 1 - \frac{c}{\ln(n)} = 1 - o(1)$$

Heavy Hitters:

length m element appears $> \frac{m}{2}$ times

Misra-Gries '82

initialize k bins each with elt (initially null)
and a counter (initially 0)

For each element e in stream

If e is in a bin b then
increment b 's counter

else if find a bin whose counter = 0
set its element = e , counter to 1

else
decrement counter of every bin

For each bin b do

$i \leftarrow$ element in bin b ,

return $\hat{f}_i = b$'s counter

If f_i is true frequency of element i

\hat{f}_i : frequency returned by algo

$$f_i - \frac{m}{k} \leq \hat{f}_i \leq f_i$$