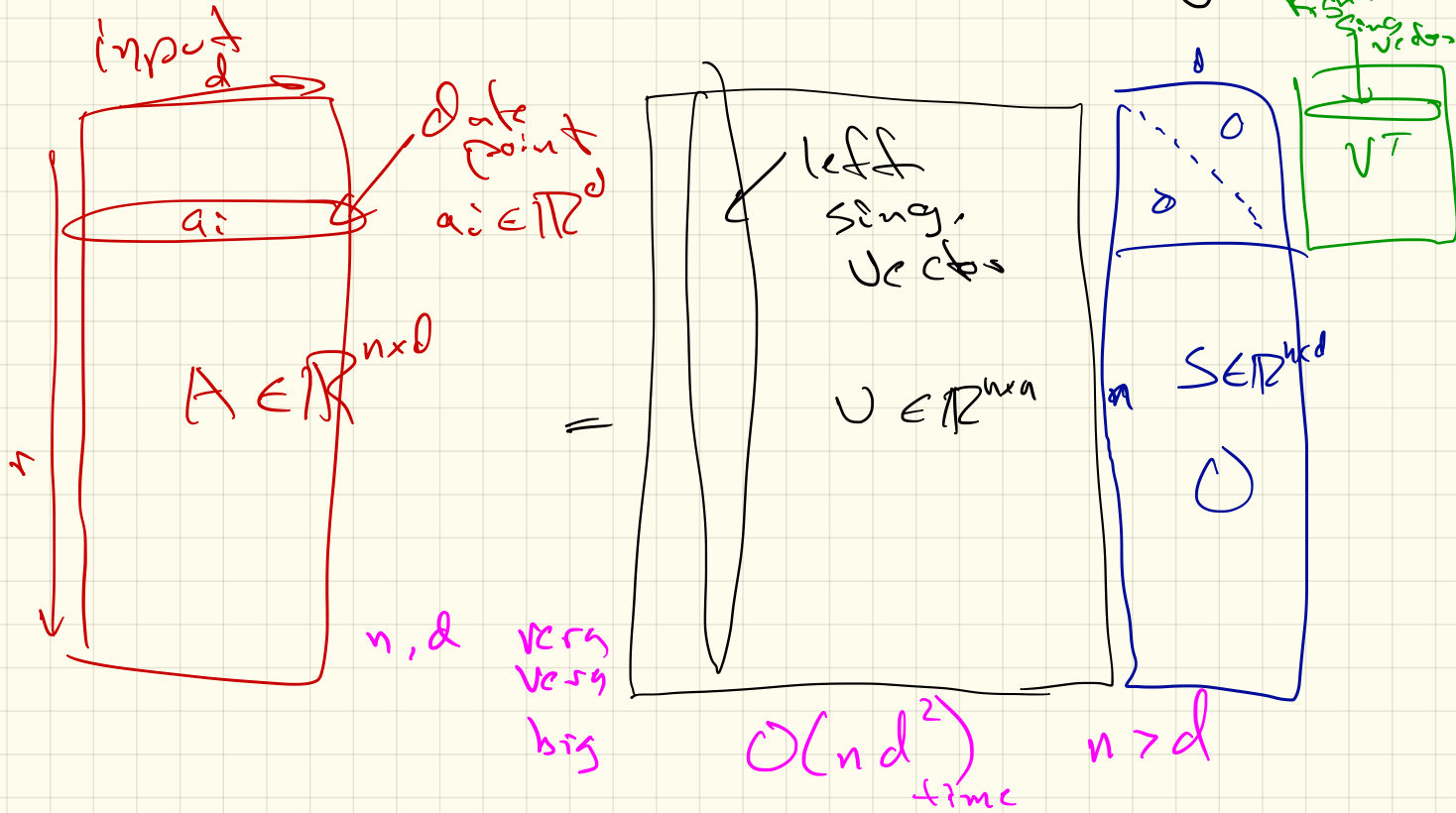


Matrix Sketching



Streaming Setting

$$A = (a_1, a_2, \dots, a_n) \quad a_i \in \mathbb{R}^d$$

n very very big, $d = \text{manageable} \approx 1000$

$C \in \text{zeros}(d, d)$

for $(a_i \in A)$

$$C = C + a_i a_i^T \in \mathbb{R}^{d \times d}$$

return C

↑ outer product

space $O(d^2)$

time $O(nd^2)$

$$C = A^T A \in \mathbb{R}^d$$

eigenvectors $(C) =$ right sing vect (A)

eigenvalue $(C) =$ squared sing, val (A)

Frequent Directions

$\Rightarrow d^2$ very very big

n very, very big
 $\Rightarrow 100$ mill

d very big
 $\Rightarrow 100$ thousand

Rank- k apx

k small
 $k \Rightarrow 10$

$\Rightarrow \tilde{O}(dk)$ space
 $\tilde{O}(ndk)$ time

Freq. Dir

$B = \text{zeros}(d \times l)$

for $(a_i \in A)$

insert a_i into B

if (B no all zero rows)

$$[U, S, V^T] = \text{svd}(B)$$

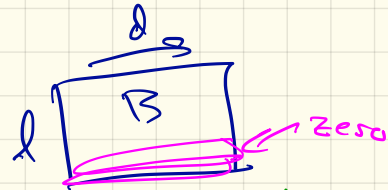
set $\delta = \text{sg}^2 \leftarrow \text{all sing val}(B)^2$

$$S' = \text{diag}(\sqrt{s_{11}^2 - \delta}, \sqrt{s_{22}^2 - \delta}, \dots, \sqrt{s_{k+1, k+1}^2 - \delta}, 0, \dots, 0)$$

Return $B = S' V^T$

$$l = k + 1/2$$

$$\text{or } k + k/\epsilon$$



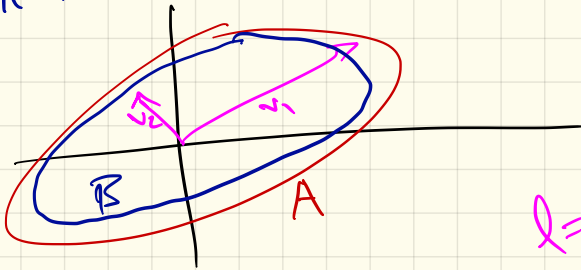
$\rightarrow O(ndl^2)$
 $\rightarrow O(ndl)$

$$B \leftarrow \text{FD}_\ell(A)$$

$$\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n a_{ij}^2 = \sum_{i=1}^n \sum_{j=1}^n c_{ij}^2$$

(1) $\forall x$
 $\|x\|=1$

$$0 \leq \|A x\|^2 - \|B x\|^2 \leq \frac{\|A - A_{\ell}\|_F^2}{\ell - \ell_2}$$



$$\ell = \frac{1}{2} + \ell_2$$

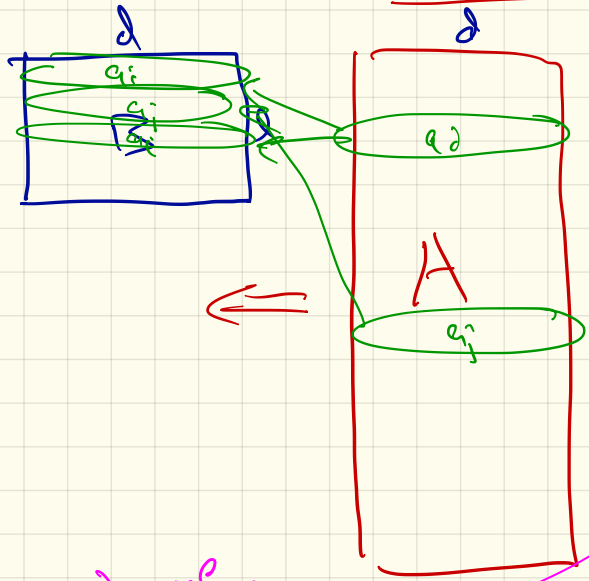
$$\ell = \ell_1/2 + \ell_2$$

$$\leq \frac{\|A\|_F^2}{\ell} \quad \ell = \frac{1}{2}$$

$$(2) \|A - A \Pi_B\|_F^2 \leq \frac{\ell}{\ell - \ell_2} \|A - A_{\ell}\|_F^2$$

projection onto ℓ_2 sing. vectors of B

Row Sampling



rows of B are
"interpretable"

- do not sample uniform
- norm-squared sampling

sample u_i w.p. $p_i \sim \|a_i\|^2$

$$\|A - AB\|_F^2 \leq \|A - A_{\text{rand}}\|_F^2 + \epsilon \|A\|_F^2 \frac{\sum p_i}{\|A\|_F^2}$$

$$l = \left(\frac{1}{\epsilon^2}\right) \log \frac{1}{\delta} \quad \text{w/ } \delta \text{ prob of fail}$$

- leverage score sampling
 $s_i(a_i) = \|U_{\text{left}}^T a_i\|^2$

$$\|A - AB\|_F^2 \leq (1 + \epsilon) \|A - A_{\text{rand}}\|_F^2$$

time
(and loss)

Reservoir Sampling

Sampling 1 item from stream

proportional to weight w_i

$$W_t = w_i$$

$$b \leftarrow a_i$$

for ($a_i \in A$)

$$w_i = w_{i-1} + w_i$$

$$v \sim \text{Unif}(0, 1)$$

$$\text{if } \left(v \geq \frac{w_i}{W_t} \right)$$

do nothing

else

$$b \leftarrow a_i$$

$$A = \langle a_1, a_2, \dots, a_t, \dots, a_n \rangle$$

\downarrow \downarrow \dots \downarrow \dots \downarrow

$$w_1 \quad w_2 \quad \dots \quad w_t \quad \dots \quad w_n$$

$$\rightarrow W_t = \sum_{i=1}^t w_i$$

Row Sampling

$$w_i = \|a_i\|^2$$

$$v \sim \text{Unif}(0, 1)$$

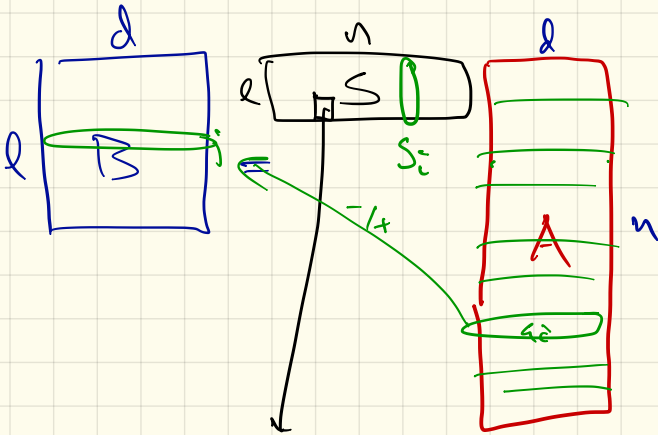
random float in $[0, 1]$

pdf



Count Sketch

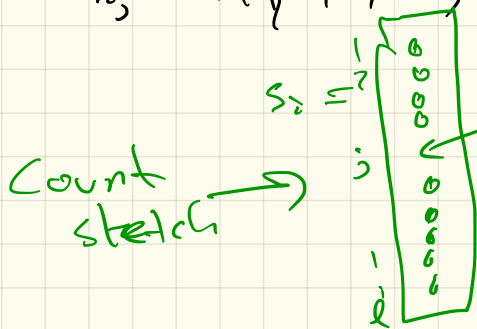
→ time $O(nd)$
 if $n \neq 0$
 number of
 non-zeros
 time $O(\text{nnz}(A))$



JL

$$s_{ij} \sim N(0, 1)$$

$$s_{ij} = \text{Unif}(\beta^{-1}, 0, \beta)$$



$$l = (d/\epsilon^2)$$

$$\forall x \in \mathbb{R}^d$$

$$(1-\epsilon) \leq \frac{\|Ax\|}{\|x\|} \leq (1+\epsilon)$$

$$l = \frac{d^2}{\epsilon^2}$$

Count sketch

$$\approx 10,000$$