# L17: Matrix Sketching

March 25, 2020

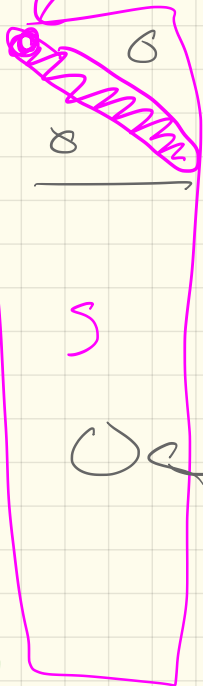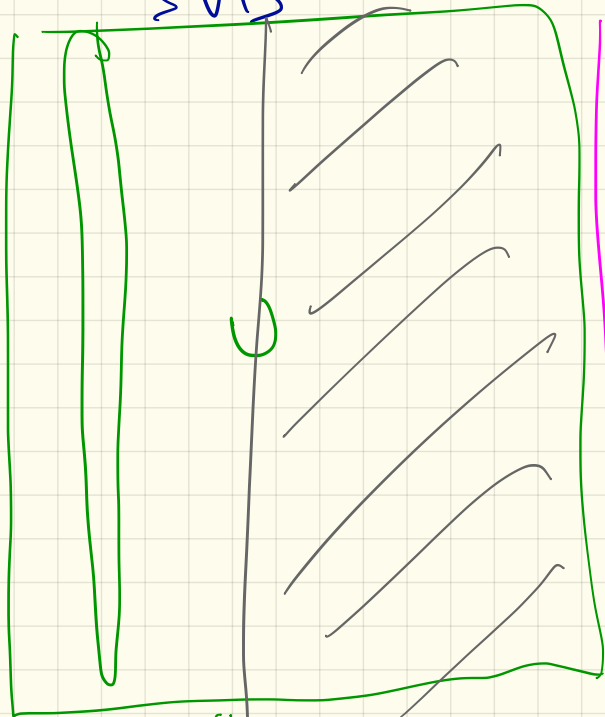Input matrix $A \in \mathbb{R}^{n \times d}$

SVD

$\sigma_1 > \sigma_2 > \dots$

$V^T$

orthogonal

$d$

$n$ | $A$ | $=$

$U$

orthogonal

$S$

$n = 10$ million
$d = 100,000$

# Eigen Value Decomposition

Input: square matrix $M \in \mathbb{R}^{d \times d}$

$$M\underline{v} = \underline{v}\lambda \quad \leftarrow \text{eigenvalue}$$

eigenvector $v \in \mathbb{R}^d$

$$\|v\| = 1$$

$$M = VLV^{-1} \qquad V \text{ orthogonal} \quad V^{-1} = V^T$$

$$L = \begin{bmatrix} \lambda_1 \lambda_2 & 0 \\ 0 & \ddots \lambda_d \end{bmatrix} \quad \begin{matrix} \lambda \geq 0 \\ \text{real} \end{matrix} \text{ if } M \begin{matrix} \text{positive} \\ \text{semidefinite} \end{matrix}$$

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

$$M_L = A A^T \in \mathbb{R}^{n \times n} \qquad \longrightarrow U, S^2$$

$\longrightarrow$ positive semidefinite

$$M_R = A^T A \, \boxed{V}$$
$$A = U S \boxed{V^T}$$
$$= (V S U^T)(U S V^T) V$$
$$\underset{I}{} \quad \underset{I}{}$$
$$= \boxed{V} S^2$$

$$S^2 = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_d^2 \end{bmatrix}$$

right sing. vectors $v_j$ of $A$

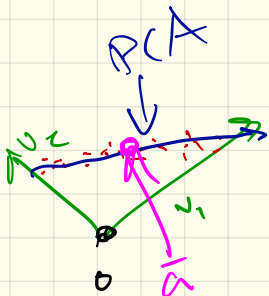are eigenvectors of $M_R$

sing values squared $\sigma_i^2 = \lambda_i$

eigenvalues of $M_R$

Find subspace $B$ ($k$-dim) $V_B = \{v_1, \ldots v_k\}$

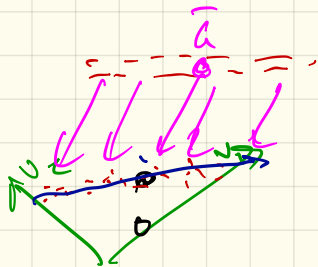PCA minimize $SSE(A, B) = \sum_{i=1}^{n} \| a_i - \pi_B(a_i) \|^2$

(if say $B$ contains $0$, the SUD opt)
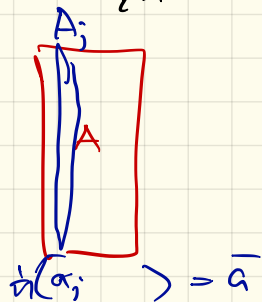


· Sol·· first center the data

For each dimension $j \in [1 \ldots d]$
find average value $\bar{a}_j = \frac{1}{n} \sum_{i=1}^{n} A_{ij}$

$$\bar{a} = (\bar{a}_1, \bar{a}_2, \ldots, \bar{a}_d)$$

$$\hat{A} = [A_{ij} - \bar{a}_j]_{ij}$$



$\pi(a_j) = \bar{a}$

# Centering Matrix

$$C_n = I_n - \frac{1}{n} \mathbb{1} \mathbb{1}^T \qquad \mathbb{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$n$ identity

$$\tilde{A} = C_n A = A - \frac{1}{n} \mathbb{1} \mathbb{1}^T A$$

$$\text{svd}(\tilde{A}) = \hat{U} \hat{S} \tilde{V}^T$$

store
$\bar{a} \in \mathbb{R}^d$

PCA

$\hat{V}^T =$ principal components

$\hat{S} =$ principal value

# Very large scale

$n \supset d$

SVD take $O(nd^2)$ time



small space $\rightarrow \hat{S} \hat{V}^T$

if $d^2$ fits in memory

$B = zeros(d \times d)$

for $i = 1$ to $n$

$\quad B \mathrel{+}= a_i a_i^T \in \mathbb{R}^{d \times d}$

Return $B = M_R$

$d$

$a_i$

$A$

$n$

If $d^2$ too big but

$l = k/\varepsilon$     $(10, dk)$ ok to fit in memory

## Frequent Direction     (Misra-Gries) but for Matrix

0. $B$ zeros $(2l \times d)$

1. for $a_i \in A$

2. Insert $a_i$ into all zero row $l$ $B$

   $\cdot$ runtime $O(nd\,l^2)$

   $\searrow O(ndl)$

3. if (no more zero rows) $\leftarrow$

4. $[U; S, V^T] = svd(B)$

5. set $\delta_i = \sigma_\ell^2$ ← still $\ell$

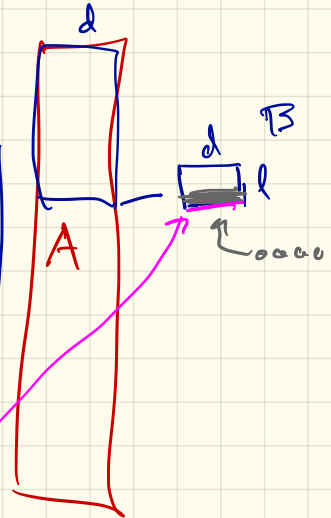6. set $S' = diag\left(\sqrt{\sigma_1^2 - \delta_i}, \sqrt{\sigma_2^2 - \delta_i}, \cdots\right)$

7. $B = S'V^T$

8. Return $B$

zero!

$l$ $\ell$

000

0

## Freq. Dim     $B = (2\ell \times d)$

for all    unit vectors $x \in \mathbb{R}^d$

$$0 \leq \|Ax\|^2 - \|Bx\|^2 \leq \frac{\|A - A_k\|_F^2}{(\ell - k)}$$

$$\ell = k + \frac{1}{\epsilon} \qquad \epsilon \cdot \|A - A_k\|_F^2$$

$$\|A - A\Pi_{B_k}\|_F^2 \leq \frac{\ell}{\ell - k} \|A - A_k\|_F^2$$

$$\underbrace{\quad}_{A_k^2}$$

$$\ell = k + k/a$$

$$\leq a \|A - A_k\|_F^2$$

# Row Sampling

$d$   $B$

Interprobable

$A$

Leverage
Scores $\bar{w}_i$

$\|A - A\pi_B\|_F$

$\leq \|A - A_k\|_F + \varepsilon \|A\|_F^2$

$\ell = \frac{k}{\varepsilon^2} \log \frac{1}{\delta}$

Sample $a_i$ proportional
to $\|a_i\|^2 = w_i$

(Reservior Sampling — sampling in stream)

$\bar{w}_i = \sum_{\xi=1}^{i} w_i$

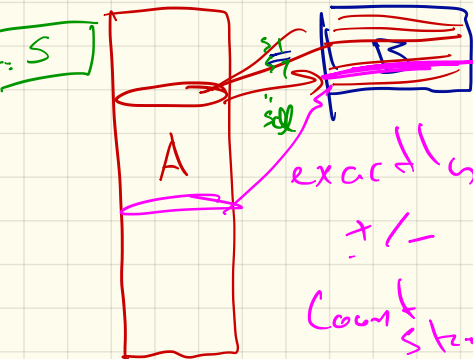replace w/ prob $\frac{w_i}{\bar{w}_i}$

Do $\ell$ times independently

Fix duplicate problem
w/ Priority Sampling

# Random Projection / Count Stretch

$$S \in \mathbb{R}^{\ell \times n} \qquad S_{ij} \sim N(0,1)\sqrt{\tfrac{n}{\ell}}$$

$$\text{sketch} \quad B = SA \in \mathbb{R}^{\ell \times d}$$



Let $\left[\hat{A} \hat{U}\right]_{12} = \underset{\text{rank}}{\text{best}} \;\; U \;\; RSU \;\; B$

$$\left\| A - \left[\hat{A} \hat{U}\right]_{12} V^{T}\right\|_{F} \leq (1+a) \left\| A - A_{12}\right\|_{F}^{2}$$

$$\ell \approx \tfrac{k}{\varepsilon}$$

$$\underbrace{\left[\hat{A} \hat{U}\right]_{12}}_{\tilde{A}_{12}}$$

exactly +/-
Count Stretch

$$(1-\varepsilon) \leq \frac{\|Ax\|}{\|Bx\|} \leq (1+a) \qquad \ell = \frac{d}{\varepsilon^2}$$

$$\forall x \in \mathbb{R}^{d} \qquad\qquad\qquad\qquad \ell = \frac{d^2}{\varepsilon^2}$$