

Matrix Sketching

Mar 24, 2025

Data Mining

Jeff M. Phillips

Input

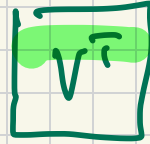
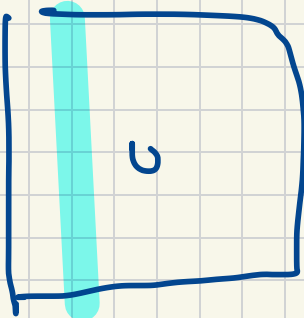
$$A \in \mathbb{R}^{n \times d}$$

rows $\{a_1, a_2, \dots, a_n\} \subset \mathbb{R}^d$



SVD

=



$\sigma_j = \sigma_j$ sing. val

← right sing vectors
 $v_j \in \mathbb{R}^d$

↑ left sing vectors
 $u_j \in \mathbb{R}^n$

$$A_k = \underset{\substack{B \\ \text{rank}(B)=k}}{\text{arg min}} \|A - B\|_2$$

or

$$\|A - B\|_F$$

$$A = \sum_{j=1}^d u_j \sigma_j v_j^T$$

$$A_k = \sum_{j=1}^k u_j \sigma_j v_j^T$$

$$\max_{k \leq d} \|(A - A_k)_{k+1}\|_2 = \|A - A_k\|_2 = \sigma_{k+1}$$

$$\|A - A_k\|_F^2 = \sum_{j=k+1}^d \sigma_j^2$$

How to compute A_k ?

↳ call "SVD" language

↳ LAPACK

runtime $O(nd \cdot \min\{n, d\}) \approx O(nd^2)$

2021
Jack Dongarra
Turing Award.

Streaming

see data in order

space $\leq n \cdot d$

$n \gg d$

$\underbrace{a_1, a_2, a_3, \dots}_{d \text{ space}}$

Setting

$n = 100$ million

$d = 100 - 1000$.

Summed Covariance

space
2

$$C \in \mathbb{R}^{d \times d}$$

$$C_{j,j} = 0$$

for $a_i \quad i = 1$ to n

$$C = C + a_i a_i^T$$

$$\text{eig}(C) = V, L$$

$$V_k \iff \text{top } k \text{ RSU}$$

$$\sqrt{L} \implies \text{sing values}$$

$$C = AA^T$$

$$\text{eig}(C) = \text{sod}(A)^2$$

Setting

$n = 100$ million

$d = 10k - 100k$

wants $k = 2, 10, 50$

d^2 too big

Need approximation parameter $l \approx \frac{1}{\epsilon}, k \approx \frac{1}{\epsilon}$

Frequent Directions

Initializing $B \in \mathbb{R}^{2^d \times d}$

for a_i : st. $i = 1$ to n

insert a_i into empty row of B

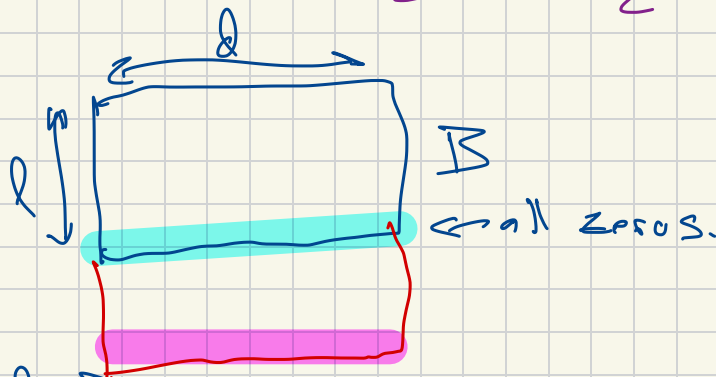
if (no all zero row)

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

$d = \sigma_1^2$

$S' = \text{diag}(\sqrt{\sigma_1^2 - \sigma}, \sqrt{\sigma_1^2 - \sigma}, \sqrt{\sigma_1^2 - \sigma}, \dots, 0, \sigma_1, 0, \dots, 0)$

$B = S' V^T$



runtime

$$O(n/k) = \frac{(2l)^2 \cdot d}{k}$$

$$= O(nd \cdot l)$$

return B .

$$FD(A) \rightarrow B \quad B \in \mathbb{R}^{l \times d}$$

$$\textcircled{1} \quad \forall x \in \mathbb{R}^d \quad \|x\|=1$$

$$0 \leq \|Ax\|^2 - \|Bx\|^2 \leq \frac{\|A - A_{\mathbb{R}}\|_F^2}{l-k} \quad l = k + \frac{k}{\epsilon}$$

$$\leq (1+\epsilon) \|A - A_{\mathbb{R}}\|_F^2$$

$$\textcircled{2} \quad \|A - \Pi_{B_k}(A)\|_F^2 \leq \frac{l}{l-k} \|A - A_{\mathbb{R}}\|_F^2 \quad l = k + \frac{k}{\epsilon}$$

$$\leq (1+\epsilon) \|A - A_{\mathbb{R}}\|_F^2$$

For Regressor

$$\text{want} \quad (1-\epsilon) \leq \frac{\|Ax\|}{\|Bx\|} < (1+\epsilon)$$

Does not hold for FD.

Setting

want: ① Have sparsities, want to maintain

$$A \rightarrow B$$

size $\neq n \cdot d \Rightarrow \text{nnz}(A)$
number of non-zeros.

② want to maintain exemplar rows.

sample l rows from $A \rightarrow$ as rows of B .

(weighted) Reservoir sampling?

row weight

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

$$\begin{aligned} & \|A - \Pi_B(A)\|_F \\ & \leq \|A - A_l\|_F + \epsilon \|A\|_F \\ & l = O\left(\frac{r}{\epsilon^2}\right) \end{aligned}$$

Better Row Sampling

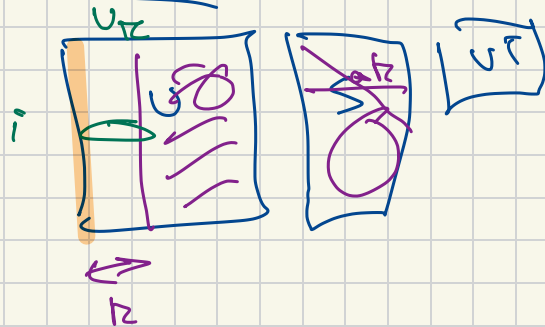
w/ Leverage Scores

$$A \xrightarrow{\text{svd}} U S V^T$$

$$\text{lev}(a_i) = \|U_{i:}^T\|^2$$

if sample $l = \mathcal{O}\left(\frac{k}{\epsilon}\right) \sim \text{lev}(a_i)$

$$\text{the } \|A - \tilde{A}_l\|_F^2 \leq (1 + \epsilon) \|A - A_k\|_F^2$$



Random Projection

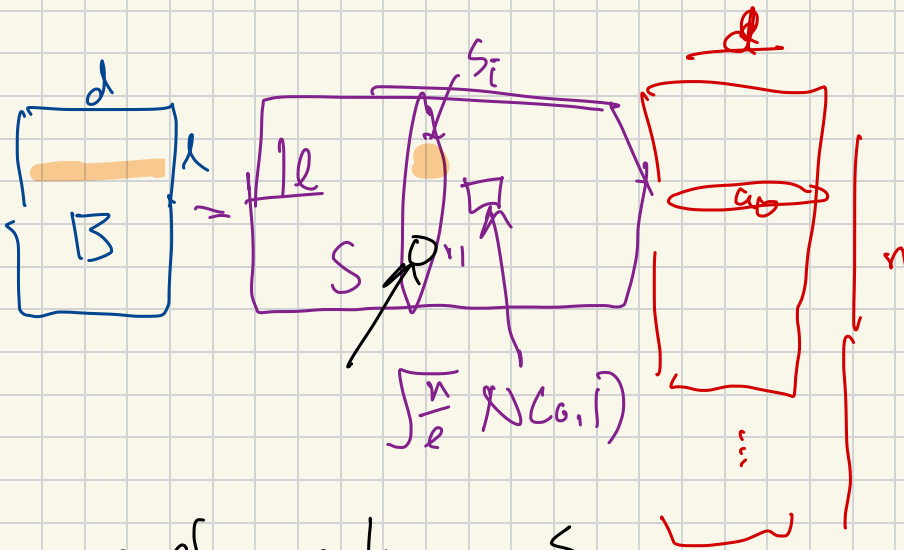
$$B = SA$$

predefine $S \in \mathbb{R}^{l \times n}$

$$S_{ij} \sim \text{iid } \mathcal{N}(0, 1) \sqrt{\frac{n}{l}}$$

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

$$l \ll n$$



$$(1-\epsilon) \leq \frac{\|A\|}{\|B\|} \leq (1+\epsilon)$$

each column s

have 1 non-zero

$$S_{ij} = \{-1, +1\}$$

Count Sketch

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

sparsify
 $O(n \cdot z(A) \cdot \frac{1}{\epsilon^2})$

Count Sketch Alg.

init $B = \mathbb{R}^{l \times d}$

$$l = O\left(\frac{d^2}{\epsilon^2}\right)$$

$$\|Bx\|_1 = \sqrt{\sum_{j=1}^d \langle b_j, x \rangle^2}$$

for $a_i \quad i=1$ to n

choose $j \sim \text{Unif}(1, d)$

choose $s_i \sim \text{Unif}(\{-1, +1\})$

$b_j = \underline{b_j + s_i \cdot a_i}$

return B .

Gaussian free $\forall x$

$$(1 - \epsilon) \leq \frac{\|Ax\|_1}{\|Bx\|_1} \leq (1 + \epsilon)$$