

# 13 matrix sketching

Monday, October 19, 2020 12:44 PM

Last time we finished off by showing that the SVD provides a natural way to approximate matrices by low-rank matrices. Why is this important?

Recall that the SVD  $A = VDU^T = \begin{bmatrix} V \end{bmatrix} \begin{bmatrix} D \end{bmatrix} \begin{bmatrix} U^T \end{bmatrix} = \sum_{i=1}^r \sigma_i v_i u_i^T$

gives an optimal rank-k approximation

$$A_k = \sum_{i=1}^k \sigma_i v_i u_i^T = \begin{bmatrix} V_k \end{bmatrix} \begin{bmatrix} D_k \end{bmatrix} \begin{bmatrix} U_k^T \end{bmatrix}$$

which is formed by taking only the top-k singular values and singular vectors.

But computing an SVD requires computing singular vectors and values, which can be slow and memory intensive. Can we do something faster?

## Thm 6.9 [Foundations of Data Science, 2020, Blum, Hopcroft, Kannan]

Let  $A \in \mathbb{R}^{m \times n}$ , and  $r, s \in \mathbb{Z}^+$ .

$$A = [A^1 \dots A^n] = \begin{bmatrix} A_1 \\ \vdots \\ A_m \end{bmatrix}, \text{ where } A^i \text{ are cols of } A, A_i \text{ are rows of } A.$$

Let  $C \in \mathbb{R}^{m \times s} = [C^1 \dots C^s]$  chosen by randomly sampling  $A^i$  as follows:

$$C^i = A^j \text{ with probability } \frac{\|A^j\|_2^2}{\|A\|_F^2}.$$

Let  $R \in \mathbb{R}^{r \times n} = \begin{bmatrix} R^1 \\ \vdots \\ R^r \end{bmatrix}$  chose by randomly sample rows  $A_i$  as follows:

$$R^i = A_i \text{ with probability } \frac{\|A_i\|_2^2}{\|A\|_F^2}.$$

$$R^i = A_j \text{ with probability } \frac{\|A_j\|_2^2}{\|A\|_F^2}.$$

Then there exists  $U \in \mathbb{R}^{s \times r}$  s.t.

$$\mathbb{E} \left( \|A - CUR\|_2^2 \right) \leq \|A\|_F^2 \left( \frac{2}{\sqrt{s}} + \frac{2r}{s} \right).$$

If we fix  $s$ , we minimize error with  $s^{2/3}$ .

$$\text{Choose } s = \frac{1}{\epsilon^3} \text{ and } r = \frac{1}{\epsilon^2}. \text{ Then } \mathbb{E} \left( \|A - CUR\|_2^2 \right) = O(\epsilon) \|A\|_F^2.$$

$$\text{i.e. } A = \begin{bmatrix} A \\ n \times m \end{bmatrix} \approx \begin{bmatrix} \text{sample} \\ \text{columns} \\ n \times s \\ C \end{bmatrix} \begin{bmatrix} \text{multiplier} \\ s \times r \\ U \end{bmatrix} \begin{bmatrix} \text{sample rows} \\ r \times m \\ R \end{bmatrix}$$

## Matrix multiplication through sampling

Let  $A \in \mathbb{R}^{m \times n}$

$B \in \mathbb{R}^{n \times p}$ . We want to approximate  $AB$  is less than  $O(mnp)$  time.

$$AB = [A^1 \dots A^n] \begin{bmatrix} B_1 \\ \vdots \\ B_n \end{bmatrix} = \sum_{k=1}^n A^k B_k \quad \left( \text{sum of outer products} \right).$$

Let's try to sample  $AB$  by taking components with prob  $p_k$ .

i.e. let  $z = k$  w.p.  $p_k$  for  $k \in \{1, \dots, n\}$  a random variable.

Define  $X = \frac{1}{P_z} A^z B_z$ , a matrix r.v.

Then the entry-wise expectation

$$\mathbb{E}X = \sum_{k=1}^n P(z=k) \cdot \frac{1}{P_k} A^k B_k = \sum_{k=1}^n A^k B_k = AB.$$

But when using an estimator, we care about both mean and variance.

Def.  $\text{Var}(X) = \mathbb{E}(\|AB - X\|_F^2)$ , the entry-wise variance

$$\text{Then } \text{Var}(X) = \sum_{i=1}^m \sum_{j=1}^P \text{Var}(x_{ij}) = \sum_{ij} \mathbb{E}(x_{ij}^2) - \mathbb{E}(x_{ij})^2 = \left( \sum_{i,j} \sum_{k=1}^n P_k \cdot \frac{1}{P_k^2} \cdot a_{ik}^2 b_{kj}^2 \right) - \|AB\|_F^2$$

doesn't matter for optimizing  $P_k$ .

We want to choose  $P_k$ 's to minimize variance.

$$\text{Note: } \sum_{i,j,k} \frac{1}{P_k} a_{ik}^2 b_{kj}^2 = \sum_k \frac{1}{P_k} \|A^k\|_2^2 \|B_k\|_2^2$$

Lemma:  $\forall c_k \geq 0$ ,  $f(p_1, \dots, p_n) = \sum_{k=1}^n \frac{c_k}{p_k}$ , subject to the constraint  $p_1 + \dots + p_n = 1$ , is minimized by  $p_k \sim \sqrt{c_k}$ .

proof.

$$\text{So } f(p_2, \dots, p_n) = \frac{c_1}{1 - (p_2 + \dots + p_n)} + \sum_{k=2}^n \frac{c_k}{p_k}$$

$$\frac{\partial f}{\partial p_k} = \frac{c_1}{(1 - (p_2 + \dots + p_n))^2} - \frac{c_k}{p_k^2} = 0 \quad \text{at optimum}$$

$$\Rightarrow \frac{p_k}{1 - (p_2 + \dots + p_n)} = \sqrt{\frac{c_k}{c_1}}$$

$$\Rightarrow p_k = \sqrt{c_k} \cdot \frac{1 - (p_2 + \dots + p_n)}{\sqrt{c_1}} \quad \forall k \neq 1.$$



Thus, we want to pick  $p_k \sim \|A^k\|_2 \|B_k\|_2$ .

Note, when  $B = A^T$ ,  $p_k \sim \|A^k\|_2^2$  (squared length of cols)

Even if  $B \neq A^T$ , this is still an upper bound on  $\text{Var}(X)$ .

So use 
$$p_k = \frac{\|A^k\|_2^2}{\|A\|_F^2}$$

$$\Rightarrow \mathbb{E}(\|AB - X\|_F^2) = \text{Var}(X) \leq \|A\|_F^2 \sum_{k=1}^n \|B_k\|_2^2 = \|A\|_F^2 \|B\|_F^2$$

Repeat with  $s$  independent trials to get  $X_1, \dots, X_s$ . Let  $\bar{X} = \frac{1}{s} \sum_{i=1}^s X_i$ .

Then 
$$\text{Var}(\bar{X}) = \frac{1}{s} \text{Var}(X) \leq \frac{1}{s} \|A\|_F^2 \|B\|_F^2$$

Note:  $\frac{1}{s} \sum_{i=1}^s X_i = \frac{1}{s} \left( \frac{A^{k_1} B_{k_1}}{p_{k_1}} + \dots + \frac{A^{k_s} B_{k_s}}{p_{k_s}} \right)$ , where each  $k_i$  is an ind. choice of col.

= CR, where

$$C = \left[ \frac{A^{k_1}}{\sqrt{s p_{k_1}}} \quad \dots \quad \frac{A^{k_s}}{\sqrt{s p_{k_s}}} \right], \quad R = \begin{bmatrix} B_{k_1} \\ \sqrt{s p_{k_1}} \\ \vdots \\ B_{k_s} \\ \sqrt{s p_{k_s}} \end{bmatrix}$$

Thm 6.5 Suppose  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times p}$ . CR as given above is an estimator for  $AB$ , and the error is bounded

$$\mathbb{E}(\|AB - CR\|_F^2) \leq \frac{\|A\|_F^2 \|B\|_F^2}{s}$$

To ensure  $\mathbb{E}(\|AB - CR\|_F^2) \leq \epsilon^2 \|A\|_F^2 \|B\|_F^2$  for some  $\epsilon > 0$ , it suffices

to make  $s \geq \frac{1}{\epsilon^2}$ . Thus, CR can be computed in  $O(msp)$  time.  
(often  $< O(msp)$  time)

Lemma 6.6: Given  $R = VDU^T$  an SVD, let  $P = R^T R = UD^+V^T R$ .  $P$  is



Then by the matrix multiplication sampling,  $\mathbb{E} \|A - AP\|_2^2 \leq \frac{\|A\|_F^2}{r}$  □

Prop. 6.8  $\|P\|_F^2 \leq r$ , if we choose  $R$  by sampling  $r$  rows of  $A$  and  $P = R^+ R$ .

pf.  $P = U D^+ V^T V D U^T = U U^T = U I_q U^T$ , where  $q = \text{rank}(R) \leq r$ .

Then  $\|P\|_F^2 = q \leq r$  because all the singular values are 1. □

Proof of matrix sketch  $\mathbb{E} (\|A - CUR\|_2^2) \leq \|A\|_F^2 \left( \frac{2}{\sqrt{r}} + \frac{2r}{s} \right)$

Let's approximate  $AP$  by matrix multiplication sampling.

Let  $C = s$  sampled cols of  $A$ .

$P' = s$  sampled rows of  $P = R^+ R$

Then  $CP' = C \underbrace{\begin{bmatrix} 1 & & \\ & 1 & \\ & & \ddots \end{bmatrix}}_{\text{sampling matrix}} R^+ R \equiv CUR$ .

Then  $\mathbb{E} \|AP - CUR\|_2^2 \leq \mathbb{E} \|AP - CUR\|_F^2 \leq \frac{\|A\|_F^2 \|P\|_F^2}{s} \leq \frac{r}{s} \|A\|_F^2$ .

But  $\|A - CUR\|_2 \leq \|A - AP\|_2 + \|AP - CUR\|_2$

$$\begin{cases} c = a + b \\ c^2 \leq 2a^2 + 2b^2 \end{cases}$$

$\Rightarrow \|A - CUR\|_2^2 \leq 2\|A - AP\|_2^2 + 2\|AP - CUR\|_2^2$

$\Rightarrow \mathbb{E} \|A - CUR\|_2^2 \leq \frac{2}{\sqrt{r}} \|A\|_F^2 + \frac{2r}{s} \|A\|_F^2$  □