

23. Matrix sketching

Tuesday, November 2, 2021 9:18 PM

Recall: SVDs are one way to approximate a matrix (MAT1850)

Today: Sampling for matrix products.

Let $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$. We want to approximate $AB \in \mathbb{R}^{m \times p}$

Naive approach: Matrix multiplication $O(mnp)$ time. (maybe slightly faster with Strassen, etc.)

Let $A(:, k)$ be the k th column of A
 $B(k, :)$ be the k th row of B .

Then $AB = \sum_{k=1}^n A(:, k) B(k, :)$ (outer prod.) $m \begin{bmatrix} \end{bmatrix} \begin{bmatrix} \end{bmatrix}^p$

Let's sample AB by taking components with prob. $p_k > 0$. ← arbitrary for now.

i.e. Let $z = k$ w.p. p_k for $k \in [n]$, a r.v.

Define: $X = \frac{1}{p_z} A(:, z) B(z, :)$, a matrix r.v.

Then the entry-wise expectation

$$\mathbb{E} X = \sum_{k=1}^n \mathbb{P}(z=k) \frac{1}{p_k} A(:, k) B(k, :) = \sum_{k=1}^n p_k \frac{1}{p_k} A(:, k) B(k, :) = AB$$

↙ cancellation is reason we scaled by $\frac{1}{p_k}$.

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

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

↘ doesn't matter for minimizing variance

Want to choose p_k to minimize variance.

$$\sum_{ij} \sum_k p_k \cdot \frac{1}{p_k^2} a_{ik}^2 b_{kj}^2 = \sum_k \frac{1}{p_k} \left(\sum_i a_{ik}^2 \right) \left(\sum_j b_{kj}^2 \right) = \sum_k \frac{1}{p_k} \|A(:, k)\|^2 \|B(k, :)\|^2$$

↙ Euclidean norm

Lemma: Let $f(p_1, \dots, p_n) = \sum_{k=1}^n \frac{c_k}{p_k}$, where $c_k \geq 0$ and $p_k > 0$.

Then subject to the condition $p_1 + \dots + p_n = 1$, the minimum of f is achieved by $p_k \propto \sqrt{c_k}$.

↙ $1 = \frac{c_1}{\sqrt{c_1}} + \dots + \frac{c_n}{\sqrt{c_n}}$ (an unconstrained optimization)

is achieved by $p_k \sim \sqrt{c_k}$.

proof. WLOG say $c_1 > 0$ and remove p_1 by $f(p_2, \dots, p_n) = \frac{c_1}{1 - (p_2 + \dots + p_n)} + \sum_{k=2}^n \frac{c_k}{p_k}$. (an unconstrained optimization)

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

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

Thus, we want to pick $p_k \sim |A(:, k)| |B(k, :)|$.

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

Even if $B \neq A^T$, can still use as easy to analyze upper bound.

Use
$$p_k = \frac{|A(:, k)|^2}{\|A\|_F^2}$$

$$\Rightarrow \mathbb{E}(\|AB - X\|_F^2) = \text{Var}(X) \leq \|A\|_F^2 \sum_k |B(k, :)|^2 = \|A\|_F^2 \|B\|_F^2.$$

Repeat with s ind. trials, getting X_1, \dots, X_s .

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

$$\begin{bmatrix} A \\ m \times n \end{bmatrix} \begin{bmatrix} B \\ n \times p \end{bmatrix} \approx \begin{bmatrix} \text{Scaled sampled cols of } A \\ m \times s \end{bmatrix} \begin{bmatrix} \text{Corresponding scaled rows of } B \\ s \times p \end{bmatrix}$$

$$\frac{1}{s} \sum_{i=1}^s X_i = \frac{1}{s} \left(\frac{A(:, k_1) B(k_1, :)}{p_{k_1}} + \dots + \frac{A(:, k_s) B(k_s, :)}{p_{k_s}} \right) = CR$$

C has cols $\frac{A(:, k_i)}{\sqrt{s p_{k_i}}}$

Note: $\mathbb{E}(CC^T) = AA^T$

R has rows $\frac{B(k_i, :)}{\sqrt{s p_{k_i}}}$

$\mathbb{E}(R^T R) = B^T B.$

Thm 6.5

Suppose $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$. The product AB can be estimated by CR given above, and the error is bounded by

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

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

To ensure $\mathbb{E}(\|AB - CR\|_F^2) \leq \epsilon^2 \|A\|_F^2 \|B\|_F^2$, suffices to choose $s \geq \frac{1}{\epsilon^2}$.

\Rightarrow CR can be computed in $O(\frac{1}{\epsilon^2} mp + mn + np)$ time.

compute CR
sample cols of A
sample rows of B

When is this a good estimate?

Consider case $B = A^T$ for simplicity.

Then if $A = I$, $\|II^T\|_F^2 = n$, but $\frac{\|I\|_F^2 \|I\|_F^2}{s} = \frac{n^2}{s}$,

\rightarrow so need $s > n$ for bound to be useful.

Trivial estimate of 0-matrix gives error $\|AA^T\|_F^2$, so need to be at least as good

Analysis via SVD

When is SVD approximation good?

$$(A = U\Sigma V^T)$$

When top p singular values take up a large constant fraction of Frobenius mass.

Recall: $\|A\|_F^2 = \sigma_1^2 + \dots + \sigma_r^2$, where $r = \text{rank}(A)$ and $\Sigma = \begin{bmatrix} \sigma_1 & \dots & \sigma_r \end{bmatrix}$.

Suppose $\exists 0 < c < 1$ and a small integer p s.t. for a matrix A ,

$$\sigma_1^2 + \dots + \sigma_p^2 \geq c(\sigma_1^2 + \dots + \sigma_r^2).$$

Note $\|AA^T\|_F^2 = \sum_{t=1}^r \sigma_t^4$ and $\|A\|_F^2 = \sum_{t=1}^r \sigma_t^2$.

Then $\mathbb{E}(\|AA^T - CR\|_F^2) \leq \frac{\|A\|_F^2 \|A^T\|_F^2}{s}$

For approx. to be good, we want $\frac{\|A\|_F^2 \|A^T\|_F^2}{s} \leq \|AA^T\|_F^2$

$$\Leftrightarrow (\sigma_1^2 + \dots + \sigma_r^2)^2 \leq s(\sigma_1^4 + \dots + \sigma_r^4)$$

$$\Leftrightarrow s \geq \frac{(\sigma_1^2 + \dots + \sigma_r^2)^2}{(\sigma_1^4 + \dots + \sigma_r^4)}$$

because max is when $\sigma_1 = \sigma_2 = \dots = \sigma_r$

$$\frac{r^2}{r} = r$$

$\max_{\sigma_1, \dots, \sigma_r} \frac{(\sigma_1^2 + \dots + \sigma_r^2)^2}{(\sigma_1^4 + \dots + \sigma_r^4)} \geq r$ So not good bound

$$\text{But } \frac{(\sigma_1^2 + \dots + \sigma_r^2)^2}{(\sigma_1^4 + \dots + \sigma_r^4)} \leq \frac{(\sigma_1^2 + \dots + \sigma_p^2)^2}{c^2(\sigma_1^4 + \dots + \sigma_r^4)} \leq \frac{(\sigma_1^2 + \dots + \sigma_p^2)^2}{c^2(\sigma_1^4 + \dots + \sigma_p^4)} \leq \frac{p}{c^2}$$

Thus, if $s \geq \frac{p}{c^2}$, then approx. is better than 0-matrix.

Thus, we don't need to sample that many columns if the mass is concentrated

Thus, we don't need to sample that many columns if the mass is concentrated in a few singular vectors.

Intuition is that we are sampling according to squared col length, so we will probably pick out cols with large singular components.

Matrix sketch

If $A \in \mathbb{R}^{m \times n}$ is a data matrix $\left[\begin{array}{c} \text{sample} \\ \text{feature} \end{array} \right]$ (e.g. expression of mRNA in a cell),

we may want a low-dimensional representation.

SVD $A = U \Sigma V^T$ is a natural solution, but takes $O(mn \cdot \min(n, m)) \approx O(n^3)$
 ↳ also, destroys sparsity and interpretability because singular vectors are mixes of samples

What about sketching by sampling cols?

Thm 6.9 Let $A \in \mathbb{R}^{m \times n}$ and $r, s \in \mathbb{Z}^+$.

Let $C \in \mathbb{R}^{m \times s}$ of s cols of A picked via length square sampling

Let $R \in \mathbb{R}^{r \times n}$ of r rows " " " " " " "

Then we can find from C & R a matrix $U \in \mathbb{R}^{s \times r}$ s.t.

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

If we fix s , we can minimize error with $r = s^{2/3}$

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

Recall: $\|A - A_k\|_2^2 = \sigma_{k+1}^2 \leq \frac{\sigma_1^2 + \dots + \sigma_k^2}{k} \leq \frac{\|A\|_F^2}{k}$ (similar form)

best rank- k approx via SVD

$$\begin{bmatrix} A \\ n \times m \end{bmatrix} = \begin{bmatrix} \text{Sample cols} \\ n \times s \end{bmatrix} \begin{bmatrix} \text{multiplier} \\ s \times r \end{bmatrix} \begin{bmatrix} \text{sample rows} \\ r \times m \end{bmatrix}$$

$C \qquad U \qquad R$

Lemma 6.6 If RR^T is invertible, the $P = R^T (RR^T)^{-1} R$ has properties (proof omitted)

Moore-Penrose pseudo-inverse R^+

Orthogonal projection operator $\left\{ \begin{array}{l} (i) P\vec{x} = \vec{x} \text{ for every vector } \vec{x} = R^T \vec{y} \text{ (if } \vec{x} \text{ is row space of } R) \\ (ii) \text{ If } \vec{v} \perp R^T \vec{v} \text{ then } P\vec{v} = 0 \end{array} \right.$

Orthogonal
projection
operator

$$\left\{ \begin{array}{l} \text{(i)} \quad P\vec{x} = \vec{x} \text{ for every vector } \vec{x} = R^T \vec{y} \quad (\text{if } \vec{x} \text{ is row space of } R) \\ \text{(ii)} \quad \text{If } \vec{x} \perp R^T \vec{y} \quad \forall \vec{y}, \text{ then } P\vec{x} = 0 \end{array} \right.$$

If RR^T not invertible, let $\text{rank}(RR^T) = r$ and $RR^T = \sum_{k=1}^r \sigma_k \vec{u}_k \vec{v}_k^T$ the SVD.

Then $P = R^T \left(\sum_{k=1}^r \frac{1}{\sigma_k} \vec{u}_k \vec{v}_k^T \right) R$ satisfies these properties.

$$\underbrace{\sum_{k=1}^r \frac{1}{\sigma_k} \vec{u}_k \vec{v}_k^T}_{R^+ \in \mathbb{R}^{r \times r}} \in \mathbb{R}^{r \times r}$$

Prop. 6.7

$$A \approx AP \quad \text{and} \quad \mathbb{E} \left(\|A - AP\|_2^2 \right) \leq \frac{1}{\sqrt{r}} \|A\|_F^2$$

(proof omitted)

Thus, we can sample s cols of A to form C , and choose correspondingly sampled rows of P to form a $s \times m$ matrix, which we can decompose into s rows of R^+ , multiplied by R .

(Use
Thm
6.5)

Matrix sketch follows from sampled matrix multiplication on AP .