

Matrix Multiplication using sampling.

Suppose A is $m \times n$ matrix

B is $n \times p$ matrix.

Goal: Find AB .

running time = $O(mnp)$.

Q: can we do faster?

\Rightarrow we will use sampling to get an approximate product faster than the traditional multiplication.

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

$\hookrightarrow m \times 1$ matrix.

$B(k, :)$ be k th row of B .

$\hookrightarrow 1 \times p$ matrix.

Then,

$$AB = \sum_{k=1}^n A(:, k) B(k, :) \quad (\text{why?})$$

Observe that for each k ,

$A(:, k) B(k, :)$ is an $m \times p$ matrix

each element of it is a single product of elements of A and B .

~~Define a random variable Z that takes values in $\{1, 2, \dots, n\}$. ~~Let~~ and $\mathbb{P}(Z = k) =: p_k$.~~

~~\leftarrow will choose p_k later.~~

$$\sum_{k=1}^n p_k = 1$$

$$\sum_{k=1}^n P_k = 1$$

Define a random (matrix) variable

$$X = \frac{1}{P_k} A(:, k) B(k, :)$$

with probability P_k .

will choose later.

Fact: $\mathbb{E}[X] = AB$.

proof: $\mathbb{E}[X] = \sum_{k=1}^n P_k \cdot \frac{1}{P_k} A(:, k) B(k, :) = AB$.

We are interested in bounding

$$\text{Var}[X] = \mathbb{E}(\|AB - X\|_F^2)$$

i,j th entry of matrix X .

Lemma: $\text{Var}[X] = \sum_{i=1}^m \sum_{j=1}^p \text{Var}[X_{i,j}] = \sum_k \frac{1}{P_k} \|A(:, k)\|_2^2 \|B(k, :)\|_2^2 - \|AB\|_F^2$.

proof: Exercise!

What is the best choice of P_k to minimize $\text{Var}[X]$, i.e. to minimize $\sum_k \frac{1}{P_k} \|A(:, k)\|_2^2 \|B(k, :)\|_2^2$ as $\|AB\|_F^2$ is fixed?

Exercise: Suppose c_1, c_2, \dots, c_n are nonnegative. Show that the minimum of $\sum_{k=1}^n \frac{c_k}{P_k}$ subject to the constraints $P_k \geq 0$ and $\sum_k P_k = 1$ is attained when P_k is proportional to $\sqrt{c_k}$.

Length squared sampling techniques.

pick $P_k = \frac{\|A(:, k)\|_2^2}{\|A\|_F^2}$

$$E[\|AB - X\|_F^2] = \text{Var}[X] \leq \|A\|_F^2 \sum_k \|B(k, :)\|_2^2 = \|A\|_F^2 \|B\|_F^2$$

$$\Rightarrow \text{Var}[X] = E[\|AB - X\|_F^2] \leq \|A\|_F^2 \|B\|_F^2$$

How to reduce the variance?

⇒ mean of estimators.

Consider Δ independent trials of X , says

$$X_1, X_2, \dots, X_\Delta$$

and take $\frac{1}{\Delta} \sum_{i=1}^{\Delta} X_i$ as our estimate of AB .

$$\Rightarrow \text{Var}\left[\frac{1}{\Delta} \sum_{i=1}^{\Delta} X_i\right] = \frac{1}{\Delta} \text{Var}[X] \leq \frac{1}{\Delta} \|A\|_F^2 \|B\|_F^2$$

Let $k_1, k_2, \dots, k_\Delta$ be the ~~column~~ index k 's chosen in each trial. Then.

$$\frac{1}{\Delta} \sum_{i=1}^{\Delta} X_i = \frac{1}{\Delta} \left(\frac{A(:, k_1) B(k_1, :)}{P_{k_1}} + \frac{A(:, k_2) B(k_2, :)}{P_{k_2}} + \dots + \frac{A(:, k_\Delta) B(k_\Delta, :)}{P_{k_\Delta}} \right)$$

$$\text{Let } C = \begin{bmatrix} \frac{A(:, k_1)}{\sqrt{\Delta P_{k_1}}} & \frac{A(:, k_2)}{\sqrt{\Delta P_{k_2}}} & \dots & \frac{A(:, k_\Delta)}{\sqrt{\Delta P_{k_\Delta}}} \end{bmatrix}$$

the $m \times \Delta$ matrix consisting of columns which are scaled versions of the chosen columns of A .

∴ we can show that $E[C] = 0$

$$E[CC^T] = AA^T$$

$$\text{Let } R = \begin{bmatrix} - & \frac{B(k_{1,:})}{\sqrt{\Delta P_{k_1}}} & - \\ - & \frac{B(k_{2,:})}{\sqrt{\Delta P_{k_2}}} & - \\ & \vdots & \\ - & \frac{B(k_{s,:})}{\sqrt{\Delta P_{k_s}}} & - \end{bmatrix}$$

We can also show that $E[R^T R] = B^T B$.

$$\Rightarrow \frac{1}{\Delta} \sum_{i=1}^{\Delta} X_i = CR$$

$$\begin{bmatrix} \frac{A(:,k_1)}{\sqrt{\Delta P_{k_1}}} & \frac{A(:,k_2)}{\sqrt{\Delta P_{k_2}}} & \dots & \frac{A(:,k_s)}{\sqrt{\Delta P_{k_s}}} \end{bmatrix} \begin{bmatrix} - & \frac{B(k_{1,:})}{\sqrt{\Delta P_{k_1}}} & - \\ - & \frac{B(k_{2,:})}{\sqrt{\Delta P_{k_2}}} & - \\ & \vdots & \\ - & \frac{B(k_{s,:})}{\sqrt{\Delta P_{k_s}}} & - \end{bmatrix} \approx AB$$

C R

Thm: Suppose A is an $m \times n$ matrix and B $n \times p$ matrix. Then AB can be estimated by CR . The error is bounded by

$$(*) \quad E(\|AB - CR\|_F^2) \leq \frac{\|A\|_F^2 \|B\|_F^2}{\Delta}$$

Thus, to ensure $E(\|AB - CR\|_F^2) \leq \epsilon^2 \|A\|_F^2 \|B\|_F^2$, take $\Delta \geq \frac{1}{\epsilon^2}$.

Note: The multiplication CR can be carried out in time $O(mp)$ if ϵ is large enough.

Q: When is the error bound (*) good and when is it not? (still an open problem in general.)

Let's consider $B = A^T$.

1) $A = I$, then the (*) is not very good.

In this case,

$$\|AB\|_F^2 = \|AA^T\|_F^2 = \|II^T\|_F^2 = \|I\|_F^2 = n$$

$$\text{But RHS of (*)} = \frac{\|A\|_F^2 \|B\|_F^2}{\lambda} = \frac{\|I\|_F^2 \|I^T\|_F^2}{\lambda} = \frac{n^2}{\lambda}$$

\Rightarrow We would need $\lambda > n$ for the bound to be better than approximating AB by the zero matrix.

2) For general A .

Suppose that $\sigma_1, \sigma_2, \dots$ are singular values of A .

$\Rightarrow \sigma_1^2, \sigma_2^2, \dots$ are singular values of $\underbrace{AA^T}_{AB}$.

$$\text{and } \|A\|_F^2 = \sum_i \sigma_i^2 \quad \|AA^T\|_F^2 = \sum_i \sigma_i^4$$

$$\|A^T\|_F^2 = \sum_i \sigma_i^2$$

$$(*) \Rightarrow \mathbb{E}(\|AA^T - CR\|_F^2) \leq \frac{\|A\|_F^2 \|A^T\|_F^2}{\lambda} = \frac{(\sum_i \sigma_i^2)^2}{\lambda}$$

$$\text{if } \lambda \geq \frac{(\sigma_1^2 + \sigma_2^2 + \dots)^2}{\sigma_1^4 + \sigma_2^4 + \dots} = \frac{\|A\|_F^2 \|A^T\|_F^2}{\|AA^T\|_F^2} \quad \text{then}$$

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

If $\text{rank}(A) = r$, then by Cauchy-Schwarz inequality,

$$\begin{aligned} (\sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2)^2 &\leq (\sigma_1^4 + \sigma_2^4 + \dots + \sigma_r^4) \cdot r \\ \Rightarrow \frac{(\sigma_1^2 + \dots + \sigma_r^2)^2}{\sigma_1^4 + \dots + \sigma_r^4} &\leq r. \end{aligned}$$

Hence, in general s need to be at least r .

If A is full rank, this means sampling will not gain us anything over taking the whole matrix!

(Matrix)
* Elementwise sampling.

Goal: Given a matrix A , find another matrix B such that $\|A - B\|$ is small and that B is much sparser than A .

"sparse matrix" = matrix with a lot of zero entries.

Consider any $m \times n$ matrix A .

Let $A_{i,j}$ be the $m \times n$ matrix whose entries are all zeros except entry (i,j) which is set to $a_{i,j}$.

E.g. $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$.

$$A_{11} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$A_{12} = \begin{bmatrix} 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$A = A_{11} + A_{12} + A_{13} + A_{21} + A_{22} + A_{23}$$

→ For any $m \times n$ matrix

$$\Rightarrow A = \sum_{i,j} A_{i,j}$$

valued random variables

Let's ~~define~~ define a matrix B as follows.

$$B = \frac{1}{P_{i,j}} A_{i,j} \text{ with probability } P_{i,j}.$$

then

$$E[B] = \sum_{i,j} \underbrace{P(B = \frac{1}{P_{i,j}} A_{i,j})}_{P_{i,j}} \cdot \frac{1}{P_{i,j}} A_{i,j} = \sum_{i,j} A_{i,j} = A.$$

Since we cannot hope to approximate A with a matrix with only 1 nonzero,

$$A \quad \rightarrow \quad \text{Find} \quad B$$

$$\begin{bmatrix} x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \end{bmatrix} \quad \begin{bmatrix} & & & x & \\ & x & & x & x \\ & & x & x & \\ x & & & x & \end{bmatrix}$$

Consider B_1, \dots, B_n independent copies of X .

$$\text{let } B = \frac{1}{n} \sum_{k=1}^n B_k$$

$$\text{then } \mathbb{E}[B] = \frac{1}{n} \sum_{k=1}^n \mathbb{E}[B_k] = \frac{1}{n} \sum_{k=1}^n \mathbb{E}[X] = \mathbb{E}[X] = A$$

$$\Rightarrow \mathbb{E}[B] = A$$

We would like to show that B is close to its mean A with high probability.

Generalized
version of
Chernoff's

Lemma: (Matrix Bernstein inequality).

Let X_1, \dots, X_n be independent $m \times n$ matrix valued random variables such that

$$\mathbb{E}[X_k] = 0 \quad \text{and} \quad \|X_k\|_2 \leq R \quad \text{for } k=1, \dots, n.$$

Set $\sigma^2 = \max \left\{ \left\| \sum_{k=1}^n \mathbb{E}[X_k X_k^T] \right\|_2, \left\| \sum_{k=1}^n \mathbb{E}[X_k^T X_k] \right\|_2 \right\}$. Then

$$\mathbb{P} \left(\left\| \sum_{k=1}^n X_k \right\|_2 > t \right) \leq (m+n) e^{-\frac{t^2}{\sigma^2 + Rt/3}}$$

(Candès - Recht '12)

Tropp '15

To use the lemma, we need to write $B-A$ in term of a sum of mean zero matrices.

$$B - A = \frac{1}{\Delta} \sum_{k=1}^{\Delta} B_k - A = \sum_{k=1}^{\Delta} \frac{(B_k - A)}{\Delta}$$

$$\text{Let } X_k = \frac{B_k - A}{\Delta}$$

$$\text{Then } \mathbb{E}[X_k] = \mathbb{E}\left[\frac{B_k - A}{\Delta}\right] = \frac{\mathbb{E}[B_k] - A}{\Delta} = 0.$$

Let Denote $|A|_1 = \sum_{i,j} |a_{i,j}|$ sum of absolute value of all entries of A .

$$\text{Set } p_{i,j} = \frac{|a_{i,j}|}{|A|_1} \Rightarrow \sum_{i,j} p_{i,j} = \sum_{i,j} \frac{|a_{i,j}|}{|A|_1} = 1.$$

$$\text{Let's bound } R = \max_k \|X_k\|_2$$

$$\max_k \|X_k\|_2 = \max_k \left\| \left(\frac{A_{i,j} - A}{p_{i,j}} \right) / \Delta \right\|_2$$

for any $n \times n$ matrix

M and N ,

$$\|M + N\|_2 \leq \|M\|_2 + \|N\|_2$$

$$\leq \max_k \left\| \frac{A_{i,j}}{p_{i,j}} \right\|_2 \cdot \frac{1}{\Delta} + \frac{\|A\|_2}{\Delta}$$

$$\frac{A_{i,j}}{|a_{i,j}|} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & \frac{a_{i,j}}{|a_{i,j}|} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \leq \frac{|A|_1}{\Delta} + \frac{\|A\|_2}{\Delta}$$

$$\Rightarrow R \leq \max_k \|X_k\|_2 \leq \frac{|A|_1}{\Delta} + \frac{\|A\|_2}{\Delta}$$

Now we compute σ^2 .

$$\begin{aligned} \left\| \sum_k \mathbb{E}[X_k X_k^T] \right\|_2 &= \left\| \sum_k \mathbb{E} \left[\frac{(B_k - A)(B_k - A)^T}{\Delta^2} \right] \right\|_2 \\ &= \left\| \sum_k \mathbb{E} \frac{B_k B_k^T - B_k A^T - A B_k^T + A A^T}{\Delta^2} \right\|_2 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[B_k A^T] &= \mathbb{E}[B_k] A^T = A A^T \\ \mathbb{E}[A B_k^T] &= A A^T \end{aligned} \quad \Rightarrow \quad \left\| \sum_k \frac{\mathbb{E}[B_k B_k^T] - A A^T}{\Delta^2} \right\|_2$$

$$= \left\| \frac{\mathbb{E}[B_k B_k^T] - A A^T}{\Delta} \right\|_2$$

$$\leq \frac{\|\mathbb{E}[B_k B_k^T]\|_2}{\Delta} + \frac{\|A A^T\|_2}{\Delta}$$

$$= \frac{\|\mathbb{E}[B_k B_k^T]\|_2}{\Delta} + \frac{\|A\|_2^2}{\Delta}$$

To compute $\mathbb{E}[B_k B_k^T]$, we observe that

Recall that $B_k = \frac{1}{P_{ij}} A_{ij}$ with probability P_{ij} .

the (i,j) entry is a_{ij} , all other entries are zero.

$$\Rightarrow B_k B_k^T = \frac{1}{P_{ij}^2} \overset{\leftarrow}{A}_{ij} \vec{A}_{ij} \quad \text{with probability } P_{ij}$$

$$= \frac{1}{P_{ij}^2} \begin{bmatrix} 0 & 0 & 0 \\ 0 & a_{ij} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 \\ a_{ij} & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$= |A|_{ij}^2 \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$\Rightarrow B_k B_k^T = |A|_{\perp}^2 E_{i,i}$ with prob. $P_{i,j}$.
 where $E_{i,i}$ is a matrix such that (i,i) entry is 1 and other entries are 0.

$$\therefore E[B_k B_k^T] = \sum_{i,j} P(B_k B_k^T = |A|_{\perp}^2 E_{i,i}) \cdot |A|_{\perp}^2 E_{i,i}$$

$$= \sum_{i,j} P_{i,j} |A|_{\perp}^2 E_{i,i}$$

$$= \sum_{i,j} \frac{a_{i,j}}{|A|_{\perp}} |A|_{\perp}^2 E_{i,i}$$

$$= |A|_{\perp} \sum_{i,j} a_{i,j} E_{i,i}$$

$$\Rightarrow \|E[B_k B_k^T]\|_2 = |A|_{\perp} \left\| \sum_{i,j} a_{i,j} E_{i,i} \right\|_2$$

$$= |A|_{\perp} \left\| \sum_i \left(\sum_j a_{i,j} \right) E_{i,i} \right\|_2$$

$$\leq |A|_{\perp} \max_i \sum_j |a_{i,j}|$$

$$= |A|_{\perp} \|A\|_1$$

$$\therefore P(\|A - B\| > t) \leq (m+n) e^{-\frac{\lambda t^2}{|A|_{\perp} \|A\|_1 + \|A\|_2^2 + |A|_{\perp} t/3 + \|A\|_2 t/3}}$$

Set $t = \epsilon \|A\|$, and demand a failure probability of at most δ we get

$$\lambda \geq \frac{\log((m+n)/\delta)}{\epsilon^2} \left(\frac{|A|_{\perp} \|A\|_1}{\|A\|_2^2} + 1 + \frac{\epsilon \|A\|_1}{3 \|A\|_2} + \frac{\epsilon}{3} \right)$$