

Lecture 13

Sampling and Random Algorithms for Matrix Multiplication

Bao-Jun Cai, 5/27/2026

Introduction to Algorithms for Data Science and Physics IMP@Fudan, 2026

Topics of this lecture:

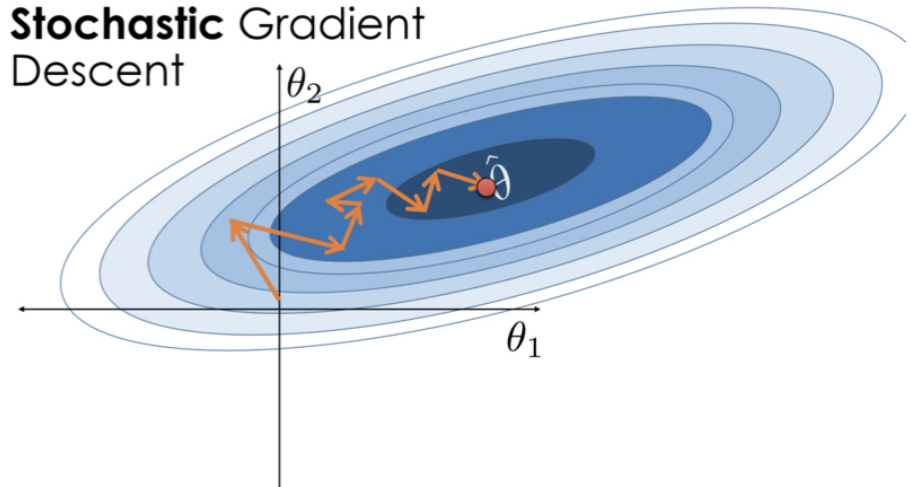
- reducing the accuracy to speed up the calculation(s)
- sampling a matrix with probability p_k
- length squared sampling scheme $p_k \sim \|\mathbf{A}(:, k)\| \cdot \|\mathbf{B}(k, :)\|$
- LR decomposition of \mathbf{AB} $\mathcal{O}(mn) \rightarrow \mathcal{O}(mdn)$
- CUR or LUR decomposition of \mathbf{A} error bound $\lesssim 2\|\mathbf{A}\|_F^2 \left(\frac{1}{\sqrt{r}} + \frac{r}{t} \right)$
- birthday paradox $K \sim \Theta(\sqrt{n})$

Motivation and examples

1. issue related to fast store in random access memory (RAM)
2. machine learning: we have a large population and want to take a small sample to perform some optimization and then argue that the optimal solution on the sample will be approximately optimal over the all population, example: stochastic gradient descent



Stochastic Gradient
Descent



2

Matrix multiplication revisited

$$\mathbf{A} \in \mathbb{R}^{m \times d}, \mathbf{B} \in \mathbb{R}^{d \times n} \rightarrow \mathbf{AB} \in \mathbb{R}^{m \times n} \sim \mathcal{O}(mdn)$$

divide-and-conquer

$$\text{Strassen's algorithm: } \sim \mathcal{O}(n^{\log_2 7})$$

Can we do better?

approximated algorithm by paying the price that some accuracy is lost

$$\text{Ex.: Show that } \mathbf{AB} = \sum_{k=1}^d \mathbf{A}(:, k)\mathbf{B}(k, :)$$

Idea of sampling

$$\mathbf{A} \in \mathbb{R}^{m \times d}, \mathbf{B} \in \mathbb{R}^{d \times n}$$

$$\mathbf{M} = \frac{1}{p_k} \overbrace{\mathbf{A}(:, k)}^{\text{kth column}} \underbrace{\mathbf{B}(k, :)}_{\text{kth row}} \leftrightarrow p_k: \text{some probability}$$

$$\mathbb{E}[\mathbf{M}] = \sum_{k=1}^d p_k \frac{1}{p_k} \mathbf{A}(:, k) \mathbf{B}(k, :) = \sum_{k=1}^d \mathbf{A}(:, k) \mathbf{B}(k, :) = \mathbf{A} \mathbf{B}$$

How can one select the probability p_k ?

Error produced by the sampling

$$\begin{aligned} \text{var}[\mathbf{M}] &\equiv \mathbb{E} [\|\mathbf{AB} - \mathbf{M}\|_{\text{F}}^2] && \text{Frobenious norm: } \|\mathbf{A}\|_{\text{F}} \equiv \sqrt{\sum_{i,j} a_{ij}^2} \\ &= \sum_{i=1}^m \sum_{j=1}^n \text{var}[m_{ij}] = \sum_{i,j} (\mathbb{E}[m_{ij}^2] - \mathbb{E}^2[m_{ij}]) && \mathbf{AB} \in \mathbb{R}^{m \times n} \\ &= \left(\sum_{i,j} \sum_k p_k \frac{1}{p_k^2} a_{ik}^2 b_{kj}^2 \right) - \|\mathbf{AB}\|_{\text{F}}^2 && \text{Ex.: for } f = \sum_k c_k/p_k \text{ with } \sum_k p_k = 1, \\ &= \sum_k \frac{1}{p_k} \left(\sum_i a_{ik}^2 \right) \left(\sum_j b_{kj}^2 \right) = \sum_k \frac{1}{p_k} \|\mathbf{A}(:, k)\|^2 \cdot \|\mathbf{B}(k, :)\|^2 && \text{what is the optimal } p_k \text{ for minimizing } f? \end{aligned}$$

Length squared sampling

probability $p_k \sim \|\mathbf{A}(:, k)\| \cdot \|\mathbf{B}(k, :)\|$

$$\mathbf{B} = \mathbf{A}^\top \rightarrow p_k \sim \|\mathbf{A}(:, k)\|^2 \rightarrow p_k = \frac{\|\mathbf{A}(:, k)\|^2}{\|\mathbf{A}\|_F^2}$$

$$\text{var}[\mathbf{M}] = \sum_k \frac{1}{p_k} \|\mathbf{A}(:, k)\|^2 \cdot \|\mathbf{B}(k, :)\|^2$$

$$\text{var}[\mathbf{M}] \leq \|\mathbf{A}\|_F^2 \sum_k \|\mathbf{B}(k, :)\|^2 = \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2$$

More trials

$$\begin{pmatrix} \mathbf{A} \\ m \times d \end{pmatrix} \begin{pmatrix} \mathbf{B} \\ d \times n \end{pmatrix} \approx \begin{pmatrix} \text{sampled} \\ \text{scaled} \\ \text{columns} \\ \text{of } \mathbf{A} \\ m \times t \end{pmatrix} \begin{pmatrix} \text{sampled scaled} \\ \text{rows of } \mathbf{B} \\ t \times n \end{pmatrix}$$

$$\begin{aligned} \langle \mathbf{M} \rangle_t &= t^{-1} \sum_{i=1}^t \mathbf{M}_i \\ &= \frac{1}{t} \left(\frac{\mathbf{A}(;, k_1) \mathbf{B}(k_1, ;)}{p_{k_1}} + \frac{\mathbf{A}(;, k_2) \mathbf{B}(k_2, ;)}{p_{k_2}} + \dots + \frac{\mathbf{A}(;, k_t) \mathbf{B}(k_t, ;)}{p_{k_t}} \right) \\ &= \underbrace{\left(\frac{\mathbf{A}(;, k_1)}{\sqrt{t p_{k_1}}}, \frac{\mathbf{A}(;, k_2)}{\sqrt{t p_{k_2}}}, \dots, \frac{\mathbf{A}(;, k_t)}{\sqrt{t p_{k_t}}} \right)}_{\mathbf{L} \in \mathbb{R}^{m \times t}, \text{ columns}} \cdot \underbrace{\left(\frac{\mathbf{B}(k_1, ;)}{\sqrt{t p_{k_1}}}, \frac{\mathbf{B}(k_2, ;)}{\sqrt{t p_{k_2}}}, \dots, \frac{\mathbf{B}(k_t, ;)}{\sqrt{t p_{k_t}}} \right)}_{\mathbf{R} \in \mathbb{R}^{t \times n}, \text{ rows}} \\ &= \mathbf{LR} \end{aligned}$$

$$\mathbf{AB} \approx \mathbf{LR}$$

Ex.: Show that $\mathbb{E}[\mathbf{L}\mathbf{L}^\top] = \mathbf{A}\mathbf{A}^\top$, $\mathbb{E}[\mathbf{R}^\top\mathbf{R}] = \mathbf{B}^\top\mathbf{B}$.

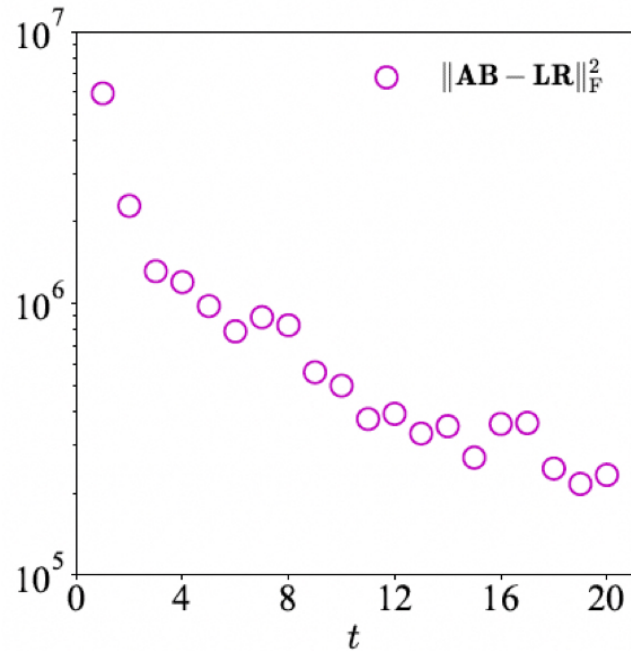
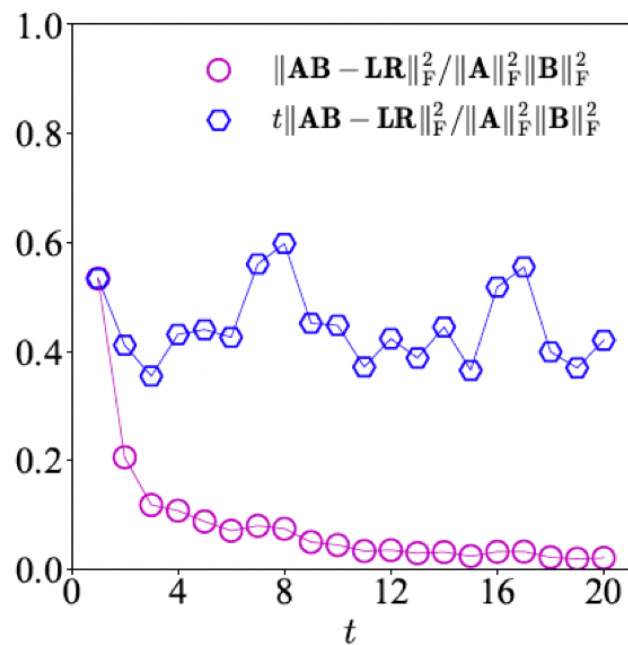
$$\mathbb{E} [\| \langle \mathbf{M} \rangle_t - \mathbf{AB} \|_F^2] \leq \frac{\| \mathbf{A} \|_F^2 \| \mathbf{B} \|_F^2}{t}$$

$$\mathbb{E} [\| \langle \mathbf{M} \rangle_t - \mathbf{AB} \|_F^2] \leq \epsilon^2 \| \mathbf{A} \|_F^2 \| \mathbf{B} \|_F^2$$

→ if $\epsilon \sim \mathcal{O}(1)$, then $t \sim \mathcal{O}(1) \rightarrow \mathcal{O}(mn)$

better than $\mathcal{O}(mnd)$ ↑

Simulation example



Ex.: How smooth the pattern of the curve in the right panel?

$$m = d = n = 10^2, a_{ij}, b_{ij} \sim \text{Unif}[0, 1]$$

$$\text{roughly, } \frac{t \|\mathbf{LR} - \mathbf{AB}\|_F^2}{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2} \gtrsim 0.5$$

Is the algorithm always useful?

$$\mathbf{B} = \mathbf{A}^\top, \mathbf{A} = \vec{\mathbf{I}}_m$$

$$\mathbb{E} [\|\mathbf{AB} - \mathbf{LR}\|_{\mathbb{F}}^2] \leq \frac{\|\mathbf{A}\|_{\mathbb{F}}^2 \|\mathbf{B}\|_{\mathbb{F}}^2}{t}$$

$$(\|\mathbf{AA}^\top\|_{\mathbb{F}}^2 = m)$$

$$\Leftrightarrow m \leq \frac{m^2}{t} \Leftrightarrow t > m \text{ (no advantage)}$$

Spectrum of the matrices involved

if SV $\sigma_k \in \mathbf{A}$, then SV $\sigma_k^2 \in \mathbf{A}\mathbf{A}^\top$

therefore $\|\mathbf{A}\mathbf{A}^\top\|_F^2 = \sum_k \sigma_k^4$, $\|\mathbf{A}\|_F^2 = \sum_k \sigma_k^2$

so $E[\|\mathbf{A}\mathbf{A}^\top - \mathbf{L}\mathbf{R}\|_F^2] \leq \|\mathbf{A}\mathbf{A}^\top\|_F^2 \leftrightarrow t \geq \frac{(\sum_k \sigma_k^2)^2}{\sum_k \sigma_k^4}$

$\text{rank}\mathbf{A} = r \rightarrow$ best upper bound of $\frac{(\sum_k \sigma_k^2)^2}{\sum_k \sigma_k^4} \approx r \rightarrow t \gtrsim r \rightarrow$ if \mathbf{A} is full rank, no advantage

related case:

if $\sigma_1^2 + \dots + \sigma_z^2 \geq f(\sigma_1^2 + \dots + \sigma_r^2)$

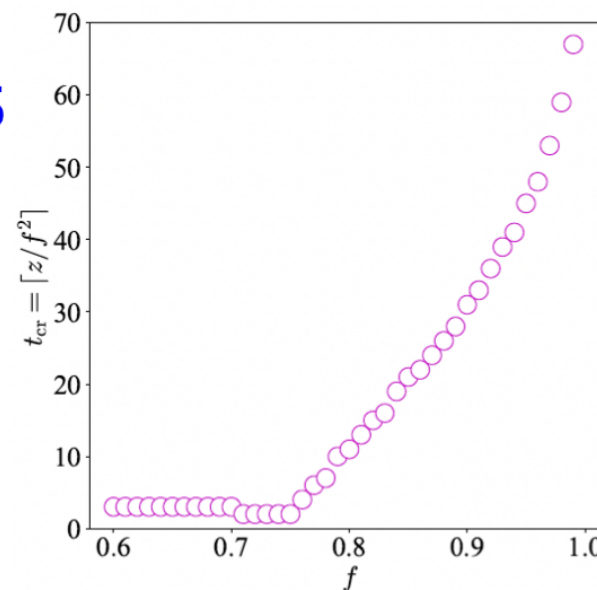
$t_{cr} \approx 3 \sim 5$

then $\frac{(\sigma_1^2 + \dots + \sigma_r^2)^2}{\sigma_1^4 + \dots + \sigma_r^4} \leq \frac{1}{f^2} \frac{(\sigma_1^2 + \dots + \sigma_z^2)^2}{\sigma_1^4 + \dots + \sigma_z^4} \leq \frac{z}{f^2}$

Ex.: Try to approximate $\mathbf{A}\vec{1}$.

$$t \gtrsim t_{cr} \approx \left\lceil \frac{z}{f^2} \right\rceil$$

10



Better choice for the identity matrix

define $\mathbf{P} = \mathbf{R}^\top (\mathbf{R}\mathbf{R}^\top)^{-1} \mathbf{R}$ and $\mathbf{x} = \mathbf{R}^\top \mathbf{y}$ (row space of \mathbf{R})

then $\mathbf{P}\mathbf{x} = \mathbf{R}^\top (\mathbf{R}\mathbf{R}^\top)^{-1} \mathbf{R}\mathbf{R}^\top \mathbf{y} = \mathbf{x} \rightarrow \mathbf{P}$ plays the role of $\vec{\mathbf{I}}$

consider decompose $\mathbf{A}\mathbf{P}$

* sample t columns of \mathbf{A} to form $\mathbf{L} \in \mathbb{R}^{m \times t}$

* sample r rows of \mathbf{A} to form \mathbf{R}

* construct $\mathbf{P} = \mathbf{R}^\top (\mathbf{R}\mathbf{R}^\top)^{-1} \mathbf{R}$

(orthogonal projection onto the row space of \mathbf{R})

($\mathbf{A}\mathbf{P}$: projection of \mathbf{A} onto the sampled row subspace)

* so $\mathbf{A}\mathbf{P} = \mathbf{B}\mathbf{C}$ with $\mathbf{B} = \mathbf{A}\mathbf{R}^\top$, $\mathbf{C} = (\mathbf{R}\mathbf{R}^\top)^{-1} \mathbf{R}$

(each column of \mathbf{B} is from sampled row of \mathbf{A})

* construct \mathbf{C} , these rows form \mathbf{U}

$$\begin{pmatrix} \mathbf{A} \\ m \times d \end{pmatrix} \approx \begin{pmatrix} \text{sampled} \\ \text{scaled} \\ \text{columns} \\ \text{of } \mathbf{A} \\ m \times t \end{pmatrix} \begin{pmatrix} \quad \quad \\ t \times r \end{pmatrix} \begin{pmatrix} \text{sampled scaled} \\ \text{rows of } \mathbf{A} \\ r \times d \end{pmatrix}$$

LUR or CUR decomposition

*Algorithm analysis

$$\|\mathbf{P}\|_F^2 \leq r$$

$$\mathbb{E} [\|\mathbf{AP} - \mathbf{LUR}\|_2^2] \leq \mathbb{E} [\|\mathbf{AP} - \mathbf{LUR}\|_F^2] \leq \frac{\|\mathbf{A}\|_F^2 \|\mathbf{P}\|_F^2}{t} \leq \frac{r}{t} \|\mathbf{A}\|_F^2$$

*bound of $\|\mathbf{A} - \mathbf{AP}\|_2^2$

$$\|\mathbf{A} - \mathbf{AP}\|_2^2 = \max_{\|\mathbf{x}\|=1} \|(\mathbf{A} - \mathbf{AP})\mathbf{x}\|^2$$

\mathbf{x} is in row space of $\mathbf{R} \rightarrow \mathbf{Px} = \vec{0} \rightarrow (\mathbf{A} - \mathbf{AP})\mathbf{x} = \mathbf{Ax}$

$$\|(\mathbf{A} - \mathbf{AP})\mathbf{x}\|^2 = \|\mathbf{Ax}\|^2 = \mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} = \mathbf{x}^\top (\mathbf{A}^\top \mathbf{A} - \mathbf{R}^\top \mathbf{R}) \mathbf{x} \leq \|\mathbf{A}^\top \mathbf{A} - \mathbf{R}^\top \mathbf{R}\|_2 \underbrace{\|\mathbf{x}\|^2}_1$$

$$\rightarrow \|\mathbf{A} - \mathbf{AP}\|_2^2 \leq \|\mathbf{A}^\top \mathbf{A} - \mathbf{R}^\top \mathbf{R}\|_2 \rightarrow \|\mathbf{A} - \mathbf{AP}\|_2^2 \leq \frac{\|\mathbf{A}\|_F^2}{\sqrt{r}}$$

$\frac{1}{\sqrt{r}}$: statistical error

*Algorithm analysis: continues

$$\begin{aligned} \|\mathbf{A} - \mathbf{LUR}\|_2 &\leq \|\mathbf{A} - \mathbf{AP}\|_2 + \|\mathbf{AP} - \mathbf{LUR}\|_2 & \|\mathbf{A} - \mathbf{AP}\|_2^2 &\leq \frac{\|\mathbf{A}\|_F^2}{\sqrt{r}} \\ \rightarrow \|\mathbf{A} - \mathbf{LUR}\|_2^2 &\leq 2\|\mathbf{A} - \mathbf{AP}\|_2^2 + 2\|\mathbf{AP} - \mathbf{LUR}\|_2^2 & \mathbb{E} [\|\mathbf{AP} - \mathbf{LUR}\|_2^2] &\leq \frac{r}{t} \|\mathbf{A}\|_F^2 \end{aligned}$$

$$\mathbb{E} [\|\mathbf{A} - \mathbf{LUR}\|_2^2] \leq 2\|\mathbf{A}\|_F^2 \left(\frac{1}{\sqrt{r}} + \frac{r}{t} \right)$$

*if t is fixed $\rightarrow r = t^{2/3}$, choosing $t \sim \mathcal{O}(1/\epsilon^3)$ and $r \sim \mathcal{O}(1/\epsilon^2) \rightarrow \text{bound} \sim \mathcal{O}(\epsilon) \|\mathbf{A}\|_F^2$

$$\sigma_1^2 = 1 \text{ (scaling)} \rightarrow \mathbb{E} [\|\mathbf{A} - \mathbf{LUR}\|_2^2] \leq \epsilon \sum_i \sigma_i^2(\mathbf{A})$$

if top k SVs of $\mathbf{A} \sim \mathcal{O}(1)$ for $k \gg m^{1/3} \rightarrow \sum_i \sigma_i^2(\mathbf{A}) \gg m^{1/3}$, then $\epsilon \sim \mathcal{O}(m^{-1/3}) \rightarrow t > m$

*if just the few SVs of \mathbf{A} are large: quite good

Quiz 4: 5/27/2026

Quiz 4.1:

Explain the role(s) of ℓ -1 loss, ℓ -2 loss, ℓ -1 penalty and ℓ -2 penalty. What is LASSO? What is Huber Loss? Explain/describe them.

Quiz 4.2:

What is difference between the k -median and k -center clusterings? Write down the definitions.

Quiz 4.3:

What is the motivation for doing spectral clustering? Give an example.

Quiz 4.4:

Show that two random vectors in d D are probably orthogonal.

Quiz 4.5:

Under which condition that a best-fit subspace for $A \in \mathbb{R}^{d \times d}$ exists? Use SVD to write down it.

Birthday problem: randomized design/analysis

One asks how many people must there be in a room before there is a 50% chance that two of them were born on the same day of the year?

$B_K = \bigcap_{i=1}^K A_i$, $B_K = B_{K-1} \cap A_K$ $A_i =$ person i 's birthday is different from person j 's for all $j < i$

$$\begin{aligned}
 \text{Prob}(B_K) &= \text{Prob}(B_{K-1})\text{Prob}(A_K|B_{K-1}) \\
 &= \text{Prob}(B_{K-2})\text{Prob}(A_{K-1}|B_{K-2})\text{Prob}(A_K|B_{K-1}) = \dots \\
 &= \text{Prob}(B_1)\text{Prob}(A_2|B_1)\text{Prob}(A_3|B_2) \dots \text{Prob}(A_K|B_{K-1}) \\
 &= 1 \cdot \frac{n-1}{n} \cdot \frac{n-2}{n} \dots \frac{n-K+1}{n} \\
 &= 1 \cdot \left(1 - \frac{1}{n}\right) \cdot \left(1 - \frac{2}{n}\right) \dots \left(1 - \frac{K-1}{n}\right) \\
 &\leq \exp\left(-\frac{1}{n}\right) \exp\left(-\frac{2}{n}\right) \exp\left(-\frac{K-1}{n}\right) \\
 &= \exp\left[-\left(\frac{1}{n} + \frac{2}{n} + \dots + \frac{K-1}{n}\right)\right] = \exp\left(-\frac{K(K-1)}{2n}\right)
 \end{aligned}$$

$$\text{Prob}(A_1) = \text{Prob}(B_1) = 1$$

$$\begin{aligned}
 \text{Prob}(b_i = d) &= \frac{1}{n}, \text{Prob}(b_i = d, b_j = d) = \frac{1}{n^2} \\
 \text{Prob}(b_i = b_j) &= \sum_{d=1}^n \text{Prob}(b_i = d, b_j = d) = \frac{1}{n}
 \end{aligned}$$

$$\text{Prob}(B_K) \leq 50\% \rightarrow K \geq \frac{1}{2}(1 + \sqrt{1 + 8n \log 2}) \approx 23$$

$$K \sim \Theta(\sqrt{n})$$