

# Deck 3: Matrix multiplication: preasymptotic estimation

Math 7870: Topics in Randomized Numerical Linear Algebra

Spring 2026

Akil Narayan

## Recall: matmat

We proposed a randomized algorithm for approximating  $\mathbf{AB}$  using *uniform sampling*.

The basic idea was to write  $\mathbf{AB}$  as a sum of rank-1 outer products, and form an (unbiased) estimator by uniformly at random summing  $N$  of the rank-1 matrices.

We identified, in principle, the type of distribution that the estimator has: by the CLT, a normal centered random variable with a total variance scaling like  $1/N$ .

## Recall: matmat

We proposed a randomized algorithm for approximating  $\mathbf{AB}$  using *uniform sampling*.

The basic idea was to write  $\mathbf{AB}$  as a sum of rank-1 outer products, and form an (unbiased) estimator by uniformly at random summing  $N$  of the rank-1 matrices.

We identified, in principle, the type of distribution that the estimator has: by the CLT, a normal centered random variable with a total variance scaling like  $1/N$ .

What needs to be done: guarantees, and pre-asymptotic estimation.

## A simplification of matmat

It'll be convenient for us to get the crux of the ideas by simplifying the problem:

Given vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^k$ , let's approximate  $\langle \mathbf{b}, \mathbf{a} \rangle = \mathbf{a}^T \mathbf{b}$  using the same idea as before.

The goal is to *not* sample the entire set of entries of both vectors.

## A simplification of matmat

It'll be convenient for us to get the crux of the ideas by simplifying the problem:

Given vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^k$ , let's approximate  $\langle \mathbf{b}, \mathbf{a} \rangle = \mathbf{a}^T \mathbf{b}$  using the same idea as before.

The goal is to *not* sample the entire set of entries of both vectors.

The procedure now is a little more transparent:

$$\mathbf{a}^T \mathbf{b} = \sum_{j \in [k]} a_j b_j \implies p_X(ka_j b_j) = \frac{1}{k},$$

so that,

$$\mathbb{E}X = \sum_{j \in [k]} p_X(ka_j b_j)ka_j b_j = \sum_{j \in [k]} a_j b_j = \mathbf{a}^T \mathbf{b}$$

$$\text{E.g. } \text{Var } X = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \sum_j \frac{1}{k} k^2 |a_j b_j|^2 - |\mathbf{a}^T \mathbf{b}|^2 = k \sum_{j \in [k]} |a_j b_j|^2 - \left( \sum_{j \in [k]} a_j b_j \right)^2$$

## Concentration for inner products

We can explicitly compute,

$$\text{Var}X = \mathbb{E}X^2 - (\mathbb{E}X)^2 = k \sum_{j \in [k]} (a_j b_j)^2 - \left( \sum_{j \in [k]} a_j b_j \right)^2.$$

Therefore, if  $X_n \stackrel{\text{iid}}{\sim} X$ , then by the LLN + CLT,

$$\lim_{N \uparrow \infty} \frac{1}{N} \sum_{n \in [N]} X_n = \mathbf{a}^T \mathbf{b}, \quad \sqrt{N} \left( \frac{1}{N} \sum_{n \in [N]} X_n - \mathbf{a}^T \mathbf{b} \right) \stackrel{N \uparrow \infty}{\sim} \mathcal{N}(0, \text{Var}(X)).$$

This is, again, only asymptotic.

# Concentration for inner products

We can explicitly compute,

$$\text{Var}X = \mathbb{E}X^2 - (\mathbb{E}X)^2 = k \sum_{j \in [k]} (a_j b_j)^2 - \left( \sum_{j \in [k]} a_j b_j \right)^2.$$

Therefore, if  $X_n \stackrel{\text{iid}}{\sim} X$ , then by the LLN + CLT,

$$\lim_{N \uparrow \infty} \frac{1}{N} \sum_{n \in [N]} X_n = \mathbf{a}^T \mathbf{b}, \quad \sqrt{N} \left( \frac{1}{N} \sum_{n \in [N]} X_n - \mathbf{a}^T \mathbf{b} \right) \stackrel{N \uparrow \infty}{\sim} \mathcal{N}(0, \text{Var}(X)).$$

This is, again, only asymptotic.

However, we do have a preasymptotic quantitative understanding:  $\text{Var} \frac{1}{N} \sum_{n \in [N]} X_n = \frac{1}{N} \text{Var}X$ .

We can compute this variance. Define a vector  $\mathbf{c}$  as,

$$\mathbf{c} = \mathbf{a} \odot \mathbf{b} \in \mathbb{R}^k, \quad c_j = a_j b_j, \quad \mathbb{E}X = \mathbf{1}^T \mathbf{c}.$$

Then we have,

$$\text{Var}X = k \|\mathbf{c}\|_2^2 - |\mathbf{1}^T \mathbf{c}|^2$$

## Best- and worst-case variance

$$\mathbf{c} = \mathbf{a} \odot \mathbf{b}, \quad \text{Var}X = k\|\mathbf{c}\|_2^2 - \left|\mathbf{1}^T \mathbf{c}\right|^2$$

Good algorithmic performance:  $\text{Var}X$  is small, relative to the (squared) oracle value.  
What kinds of vectors  $\mathbf{c}$  maximize/minimize the variance?

## Best- and worst-case variance

$$\mathbf{c} = \mathbf{a} \odot \mathbf{b}, \quad \text{Var}X = k\|\mathbf{c}\|_2^2 - \left|\mathbf{1}^T \mathbf{c}\right|^2$$

Good algorithmic performance:  $\text{Var}X$  is small, relative to the (squared) oracle value.  
What kinds of vectors  $\mathbf{c}$  maximize/minimize the variance?

$$\frac{1}{\left|\mathbf{1}^T \mathbf{c}\right|^2} \text{Var}X = \frac{k}{\left|\mathbf{1}^T \hat{\mathbf{c}}\right|^2}, \quad \hat{\mathbf{c}} := \frac{\mathbf{c}}{\|\mathbf{c}\|_2}.$$

## Best- and worst-case variance

$$\mathbf{c} = \mathbf{a} \odot \mathbf{b}, \quad \text{Var}X = k\|\mathbf{c}\|_2^2 - \left|\mathbf{1}^T \mathbf{c}\right|^2$$

Good algorithmic performance:  $\text{Var}X$  is small, relative to the (squared) oracle value.  
What kinds of vectors  $\mathbf{c}$  maximize/minimize the variance?

$$\frac{1}{\left|\mathbf{1}^T \mathbf{c}\right|^2} \text{Var}X = \frac{k}{\left|\mathbf{1}^T \hat{\mathbf{c}}\right|^2}, \quad \hat{\mathbf{c}} := \frac{\mathbf{c}}{\|\mathbf{c}\|_2}.$$

**The best case:**  $\hat{\mathbf{c}} = \frac{1}{\sqrt{k}} \mathbf{1}$ . Then  $\text{Var}X = 0$ .  
(I.e., each  $X_j$  takes a single value, equal to  $\mathbf{a}^T \mathbf{b}$ , with probability 1.)

## Best- and worst-case variance

$$\mathbf{c} = \mathbf{a} \odot \mathbf{b}, \quad \text{Var}X = k\|\mathbf{c}\|_2^2 - \left| \mathbf{1}^T \mathbf{c} \right|^2$$

Good algorithmic performance:  $\text{Var}X$  is small, relative to the (squared) oracle value.  
What kinds of vectors  $\mathbf{c}$  maximize/minimize the variance?

$$\frac{1}{\left| \mathbf{1}^T \mathbf{c} \right|^2} \text{Var}X = \frac{k}{\left| \mathbf{1}^T \hat{\mathbf{c}} \right|^2}, \quad \hat{\mathbf{c}} := \frac{\mathbf{c}}{\|\mathbf{c}\|_2}.$$

**The best case:**  $\hat{\mathbf{c}} = \frac{1}{\sqrt{k}} \mathbf{1}$ . Then  $\text{Var}X = 0$ .  
(I.e., each  $X_j$  takes a single value, equal to  $\mathbf{a}^T \mathbf{b}$ , with probability 1.)

**The worst case:**  $\mathbf{c} \perp \mathbf{1}$ , i.e.,  $\mathbf{c}$  has positive and negative components of approximately the same mass. Then  $\text{Var}X = k\|\mathbf{c}\|_2^2$   
(I.e.,  $\sum_n X_n$  sums positive and negative components with similar “mass”.)

## Best- and worst-case variance

$$\mathbf{c} = \mathbf{a} \odot \mathbf{b}, \quad \text{Var}X = k\|\mathbf{c}\|_2^2 - \left| \mathbf{1}^T \mathbf{c} \right|^2$$

Good algorithmic performance:  $\text{Var}X$  is small, relative to the (squared) oracle value.  
What kinds of vectors  $\mathbf{c}$  maximize/minimize the variance?

$$\frac{1}{\left| \mathbf{1}^T \mathbf{c} \right|^2} \text{Var}X = \frac{k}{\left| \mathbf{1}^T \hat{\mathbf{c}} \right|^2}, \quad \hat{\mathbf{c}} := \frac{\mathbf{c}}{\|\mathbf{c}\|_2}.$$

**The best case:**  $\hat{\mathbf{c}} = \frac{1}{\sqrt{k}} \mathbf{1}$ . Then  $\text{Var}X = 0$ .  
(I.e., each  $X_j$  takes a single value, equal to  $\mathbf{a}^T \mathbf{b}$ , with probability 1.)

**The worst case:**  $\mathbf{c} \perp \mathbf{1}$ , i.e.,  $\mathbf{c}$  has positive and negative components of approximately the same mass. Then  $\text{Var}X = k\|\mathbf{c}\|_2^2$   
(I.e.,  $\sum_n X_n$  sums positive and negative components with similar “mass”.)

**A near-worst case:**  $\hat{\mathbf{c}} = \mathbf{e}_j$ , so that  $\text{Var}X = (k-1)\|\mathbf{c}\|_2^2$ .  
(I.e.,  $\mathbf{a}^T \mathbf{b}$  has a bunch of zero summands, which we randomly sample with nonzero probability....)

## Importance sampling, I

The near-worst case reveals a qualitative issue: sampling entries uniformly can provide suboptimal results.

An alternative: sampling based on knowledge of entries of  $\mathbf{a}, \mathbf{b}$ .

In particular, we can generalize our random variable to have a different mass function:

$$p_X \left( \frac{1}{p_j} a_j b_j \right) = p_j \text{ with } \sum_{j \in [k]} p_j = 1 \implies \mathbb{E}X = \mathbf{a}^T \mathbf{b}. \quad (\text{Cf. } p_j = \frac{1}{k})$$

We can craft the  $p_j$  values to improve performance. E.g., by minimizing variance.

$$\mathbb{E}X = \sum_j p_X \left( \frac{1}{p_j} a_j b_j \right) \cdot \frac{1}{p_j} a_j b_j = \sum_j p_j \frac{1}{p_j} a_j b_j = \mathbf{a}^T \mathbf{b}$$

## Importance sampling, I

The near-worst case reveals a qualitative issue: sampling entries uniformly can provide suboptimal results.

An alternative: sampling based on knowledge of entries of  $\mathbf{a}, \mathbf{b}$ .

In particular, we can generalize our random variable to have a different mass function:

$$p_X \left( \frac{1}{p_j} a_j b_j \right) = p_j \text{ with } \sum_{j \in [k]} p_j = 1 \implies \mathbb{E}X = \mathbf{a}^T \mathbf{b}.$$

We can craft the  $p_j$  values to improve performance. E.g., by minimizing variance.

Through a similar computation as before, we have,

$$\text{Var}X = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \sum_{j \in [k]} \frac{1}{p_j} (a_j b_j)^2 - \left( \sum_{j \in [k]} a_j b_j \right)^2 = \sum_{j \in [k]} \frac{1}{p_j} c_j^2 - |\mathbf{c}|^2$$

So we can attempt to solve the problem:

$$\min_{p_j} \text{Var}X \text{ subject to } \sum_{j \in [k]} p_j = 1, \quad p_j > 0$$

(Lagrange multipliers)

$$\min_{p_j} \sum_j \frac{1}{p_j} c_j^2 \quad \text{s.t.} \quad \sum_{j \in [k]} p_j = 1, \quad p_j > 0 \quad \forall j.$$

$$p_j = s_j^2 \quad \forall j.$$

$s_j$ : new, "slack" variables,  $s_j \neq 0$

$$\min_{p, s} f(p, s) \quad \text{s.t.} \quad g(p, s) = 0 \quad \text{and} \quad h_j(p_j, s_j) = 0 \quad \forall j \in [k]$$

Augmented Lagrangian:  $\mathcal{L}(p, s) = f(p, s) + \lambda g(p, s) + \mu^T h(p, s)$

Variables:  $p, s, \lambda, \mu$ , compute stationary pts.

$$\frac{\partial \mathcal{L}}{\partial p_j} = -\frac{1}{p_j^2} c_j^2 + \lambda \cdot 1 + \mu_j \cdot 1 = 0$$

$$\frac{\partial \mathcal{L}}{\partial s_j} = -2\mu_j s_j = 0 \quad (h_j = p_j - s_j^2)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_j p_j - 1 = 0$$

$$\frac{\partial \mathcal{L}}{\partial \mu_j} = p_j - \underbrace{s_j^2}_{\text{brace}} = 0$$

$$s_j = \sqrt{p_j} \neq 0$$

$$\Rightarrow \mu_j = 0$$

$$\Rightarrow -\frac{1}{p_j^2} c_j^2 + \lambda = 0 \quad (1)$$

$$\sum_j p_j = 1 \quad (2)$$

$$(1): \frac{p_j^2}{c_j^2} = \frac{1}{\lambda} \Rightarrow p_j = |c_j| \sqrt{\frac{1}{\lambda}}$$

$$(2): \sum_j p_j = 1 \Rightarrow \sqrt{\lambda} \left( \sum_j |c_j| \right) = 1$$

$$\sqrt{\lambda} = \|\underline{c}\|_1$$

$$\Rightarrow p_j = \frac{|c_j|}{\|\underline{c}\|_1}$$

$$\begin{aligned} \text{Var } X &= \sum_j p_j |c_j|^2 - \|\underline{c}\|_1^2 \\ &= \|\underline{c}\|_1^2 - \|\underline{c}\|_2^2 \end{aligned}$$

## Importance sampling, II

We have:

$$p_j = \frac{|c_j|}{\|\mathbf{c}\|_1} \implies \text{Var}X = \|\mathbf{c}\|_1^2 - (\mathbf{1}^T \mathbf{c})^2.$$

In this case, if  $\mathbf{c} = \mathbf{e}_j$ , then  $\text{Var}X = 0$ . (This was the “near” worst-case before.)

## Importance sampling, II

We have:

$$p_j = \frac{|c_j|}{\|\mathbf{c}\|_1} \implies \text{Var}X = \|\mathbf{c}\|_1^2 - (\mathbf{1}^T \mathbf{c})^2.$$

In this case, if  $\mathbf{c} = \mathbf{e}_j$ , then  $\text{Var}X = 0$ . (This was the “near” worst-case before.)

This analysis can be lifted to the case when the inner product is  $(m \times n)$ -valued (i.e., a matrix). Like before, with  $\mathbf{A} \in \mathbb{R}^{m \times k}$  and  $\mathbf{B} \in \mathbb{R}^{k \times n}$ ,

$$\mathbf{AB} = \sum_{j \in [k]} \mathbf{a}_j \mathbf{b}_j^T.$$

With  $\mathbf{X} \in \mathbb{R}^{m \times n}$  the random matrix,

$$p_{\mathbf{X}} \left( \frac{1}{p_j} \mathbf{a}_j \mathbf{b}_j^T \right) = p_j \implies \mathbb{E}X = \mathbf{AB}.$$

## Importance sampling, III

A direct computation yields that the expected Frobenius norm error is,

$$\begin{aligned}\mathbb{E} \left\| \mathbf{AB} - \frac{1}{N} \sum_{q \in [N]} \mathbf{X}_q \right\|_F^2 &= \frac{1}{N} \text{trace}(\text{Var}(\text{vec}(\mathbf{X}))) = \frac{1}{N} \sum_{(i,j) \in [m] \times [n]} \text{Var}((AB)_{i,j}) \\ &= \frac{1}{N} \left( \sum_{j \in [k]} \frac{1}{p_j} \|\mathbf{a}_j\|_2^2 \|\mathbf{b}_j\|_2^2 - \|\mathbf{AB}\|_F^2 \right).\end{aligned}$$

This quadratic norm is minimized by choosing,

$$p_j = \frac{\|\mathbf{a}_j\|_2 \|\mathbf{b}_j\|_2}{\sum_{q \in [k]} \|\mathbf{a}_q\|_2 \|\mathbf{b}_q\|_2} \implies \text{trace}(\text{Var}(\text{vec}(\mathbf{X}))) = \left( \sum_{q \in [k]} \|\mathbf{a}_q\|_2 \|\mathbf{b}_q\|_2 \right)^2 - \|\mathbf{AB}\|_F^2.$$

This does give preasymptotic quantitative understanding of first- and second-moments of the estimator.

## A note on practicality

We sort-of have a chicken vs egg problem: To compute  $\|\mathbf{a}_q\|_2$  and  $\|\mathbf{b}_q\|_2$  for all  $q$  naively, we require  $k$ -dependent complexity, which we're trying to avoid.

Sometimes there is exploitable structure in matrices that allow us to compute these values.

Alternatively, if we can *approximate* these values, then we can still achieve similar results.

## A note on practicality

We sort-of have a chicken vs egg problem: To compute  $\|\mathbf{a}_q\|_2$  and  $\|\mathbf{b}_q\|_2$  for all  $q$  naively, we require  $k$ -dependent complexity, which we're trying to avoid.

Sometimes there is exploitable structure in matrices that allow us to compute these values.

Alternatively, if we can *approximate* these values, then we can still achieve similar results.

Namely, if we can choose the probabilities  $p_j$  so that for some  $\tau \leq 1$ ,

$$p_j \geq \tau \frac{\|\mathbf{a}_j\|_2 \|\mathbf{b}_j\|_2}{\sum_{q \in [k]} \|\mathbf{a}_q\|_2 \|\mathbf{b}_q\|_2}$$

then the resulting quadratic expected error suffers a multiplicative  $1/\tau$  penalty.

## A note on practicality

We sort-of have a chicken vs egg problem: To compute  $\|\mathbf{a}_q\|_2$  and  $\|\mathbf{b}_q\|_2$  for all  $q$  naively, we require  $k$ -dependent complexity, which we're trying to avoid.

Sometimes there is exploitable structure in matrices that allow us to compute these values.

Alternatively, if we can *approximate* these values, then we can still achieve similar results.

Namely, if we can choose the probabilities  $p_j$  so that for some  $\tau \leq 1$ ,

$$p_j \geq \tau \frac{\|\mathbf{a}_j\|_2 \|\mathbf{b}_j\|_2}{\sum_{q \in [k]} \|\mathbf{a}_q\|_2 \|\mathbf{b}_q\|_2}$$

$\|\mathbf{a}_q\|_2$

then the resulting quadratic expected error suffers a multiplicative  $1/\tau$  penalty.

The point: we can sample *near*-optimally and get near-optimal results.

## Moments to probabilities, I

We've computed the *expectation* of the error.

More practical information, such as the probability of failure, require more analysis. A simple, suboptimal strategy is to use e.g., *Markov's inequality*,

$$\Pr(X \geq t) \leq \frac{\mathbb{E}X}{t} \quad (X \geq 0 \text{ wp1})$$

$$\begin{aligned} \text{("Pf": } \mathbb{E}X &= P(X \geq t) \mathbb{E}[X | X \geq t] + P(X < t) \mathbb{E}[X | X < t] \\ &\geq P(X \geq t) \mathbb{E}[X | X \geq t] \geq t P(X \geq t) \end{aligned}$$

# Moments to probabilities, I

We've computed the *expectation* of the error.

More practical information, such as the probability of failure, require more analysis. A simple, suboptimal strategy is to use e.g., *Markov's inequality*,

$$\Pr(X \geq t) \leq \frac{\mathbb{E}X}{t} \quad (X \geq 0 \text{ wp1})$$

To use this in our matrix multiplication setting, let,

$$Z = \left\| \mathbf{AB} - \frac{1}{N} \sum_{q \in [N]} \mathbf{x}_q \right\|_F^2, \quad \mathbb{E}Z = \frac{1}{N} \left[ \left( \sum_{q \in [k]} \|\mathbf{a}_q\|_2 \|\mathbf{b}_q\|_2 \right)^2 - \|\mathbf{AB}\|_F^2 \right] =: \frac{\beta}{N}.$$

Our goals are:

- Given  $\epsilon > 0$ , ensure that  $Z \leq \epsilon\beta$ .
- Given  $\delta > 0$ , ensure failure of the above with probability at most  $\delta$ .

i.e., given  $\epsilon, \delta$ , when is it true that  $\Pr(Z \geq \epsilon\beta) \leq \delta$ ?

$$\Pr(Z \geq \epsilon\beta) \leq \frac{\mathbb{E}Z}{\epsilon\beta} = \frac{\beta}{\epsilon\beta N} \stackrel{?}{\leq} \delta$$

## Moments to probabilities, II

We require  $N \geq 1/(\delta\epsilon)$  for this to occur.

This is a precise sample complexity to achieve prescribed accuracy with prescribed error.

This means: if we choose  $N \geq 1/(\delta\epsilon)$ , then  $\mathbb{E}Z \leq \epsilon\beta$  with probability at least  $1 - \delta$ .

(To achieve simplicity, we're kind of cheating here: this is a bound for *quadratic* error. Really we should worry about  $\sqrt{Z}$ . By using Jensen's inequality,  $N \geq 1/(\delta\epsilon)^2$  is the sample requirement for  $\epsilon$ -relative accuracy on  $\sqrt{Z}$ .)

# Matrix multiplication summary

Using a simple concentration strategy, we have a random sampling algorithm (with probabilistic weights depending on the column/row norms of  $\mathbf{A}, \mathbf{B}$ ) that achieves a prescribed error with a prescribed probability.

- We can explicitly compute moments.
- A variance-like quadratic deviation can be minimized by choosing appropriate probabilities (that require knowledge of  $\mathbf{A}, \mathbf{B}$ ).
- These moments can be transformed into failure probabilities through inequalities. (We used Markov's inequality.)
- This results in precise sample requirements to achieve (error, success).

# Matrix multiplication summary

Using a simple concentration strategy, we have a random sampling algorithm (with probabilistic weights depending on the column/row norms of  $\mathbf{A}, \mathbf{B}$ ) that achieves a prescribed error with a prescribed probability.

- We can explicitly compute moments.
- A variance-like quadratic deviation can be minimized by choosing appropriate probabilities (that require knowledge of  $\mathbf{A}, \mathbf{B}$ ).
- These moments can be transformed into failure probabilities through inequalities. (We used Markov's inequality.)
- This results in precise sample requirements to achieve (error, success).
- The resulting sampling complexity is not that great: ensuring an at-most 10% failure rate with 10% relative error requires 100 samples. (And this is to guarantee achieving the *quadratic* variance-type error.)
- We can do better...with some more work. The way we've transformed moments into probabilities is a very loose translation. Stronger, sharper results require more precise estimates of *concentration*.