

Основы матричных вычислений

Весенний семестр 2021

Лекция 10: Умножение матриц и вычислительная  
устойчивость

Максим Рахуба

Высшая Школа Экономики

# План

## Матричное умножение

Метод Штрассена

BLAS

## Устойчивость и обусловленность

Машинные числа

Вычислительная устойчивость

Обусловленность

## Матричное умножение: сложность

$$C = A B$$

$\begin{matrix} \in \mathbb{R}^{n \times n} & & \\ & \in \mathbb{R}^{n \times n} & \end{matrix}$

— сложность:  $2n^3 + O(n^2)$

память:  $O(n^2)$

можно ли быстрее, чем за  $O(n^3)$ ?

# План

## Матричное умножение

Метод Штрассена

BLAS

## Устойчивость и обусловленность

Машинные числа

Вычислительная устойчивость

Обусловленность

# Матричное умножение: метод Штрассена<sup>1</sup>

$$\begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}, \quad A_{ij}, B_{ij} - \frac{n}{2} \times \frac{n}{2} \text{ матрицы}$$

“Строка на столбец”:

$$C_{11} = A_{11}B_{11} + A_{12}B_{21}$$

$$C_{12} = A_{11}B_{12} + A_{12}B_{22}$$

$$C_{21} = A_{21}B_{11} + A_{22}B_{21}$$

$$C_{22} = A_{21}B_{12} + A_{22}B_{22}$$

8 умножений и 4 сложения

$$2 \left( \frac{n}{2} \right)^3 \cdot 8 -$$

Нет выигрыша

Штрассен:

$$M_1 = (A_{11} + A_{22})(B_{11} + B_{22})$$

$$M_2 = (A_{21} + A_{22})B_{11}$$

$$M_3 = A_{11}(B_{12} - B_{22})$$

$$M_4 = A_{22}(B_{21} - B_{11})$$

$$M_5 = (A_{11} + A_{12})B_{22}$$

$$M_6 = (A_{21} - A_{11})(B_{11} + B_{12})$$

$$M_7 = (A_{12} - A_{22})(B_{21} + B_{22})$$

$$C_{11} = M_1 + M_4 - M_5 + M_7$$

$$C_{12} = M_3 + M_5$$

$$C_{21} = M_2 + M_4$$

$$C_{22} = M_1 + M_3 - M_2 + M_6$$

7 умножений и 18 сложений

<sup>1</sup>Strassen, V. (1969). Gaussian elimination is not optimal. Numerische mathematik, 13(4), 354-356. (<https://link.springer.com/content/pdf/10.1007/BF02165411.pdf>)

## Матричное умножение: метод Штрассена

$$M(n) = 7 M\left(\frac{n}{2}\right) = 7 \cdot 7 M\left(\frac{n}{4}\right) = \dots = 7^{\log_2 n} = n^{\log_2 7 \approx 2.81}$$

$$A(n) = 7 A\left(\frac{n}{2}\right) + 18 \cdot \left(\frac{n}{2}\right)^2 \xrightarrow{\text{мастер троп.}} A(n) = O(n^{\log_2 7})$$

$$7 n^{\log_2 7} < 2 n^3 \quad \text{при } n \geq 700$$

$\Rightarrow$  не надо "раскруз." рекурсивно до конца

---

**Мастер троп**  $T(n) = a \cdot T\left(\frac{n}{b}\right) + O(n^c)$  и  $c < \log_b a \Rightarrow$

$$T(n) = O(n^{\log_b a})$$

# Матричное умножение: метод Штрассена

$$\begin{bmatrix} C_1 & C_2 \\ C_3 & C_4 \end{bmatrix} = \begin{bmatrix} A_1 & A_2 \\ A_3 & A_4 \end{bmatrix} \begin{bmatrix} B_1 & B_2 \\ B_3 & B_4 \end{bmatrix}$$

$$C_1 = A_1 B_1 + A_2 B_3$$

$$C_2 = A_1 B_2 + A_2 B_4 \Rightarrow$$

$$C_3 = A_3 B_1 + A_4 B_3$$

$$C_4 = A_3 B_2 + A_4 B_4$$

$$C_k = \sum_{i,j=1}^4 x_{kij} A_i B_j$$

$$x_{kij} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$x_{kij} = \sum_{\alpha=1}^{R=7} u_{i\alpha} v_{j\alpha} w_{\alpha k}$$

$$C_k = \sum_{\alpha=1}^7 w_{\alpha k} \left( \sum_{i=1}^4 u_{i\alpha} A_i \right) \cdot \left( \sum_{j=1}^4 v_{j\alpha} B_j \right)$$

7 матриц

## Матричное умножение: метод Штрассена

- ▶ Мировой рекорд [J. Alman, V.V. Williams, 2020]<sup>2</sup>:  $< \mathcal{O}(n^{2.37286})$ . Но большая константа в  $\mathcal{O}(\cdot)$ .
- ▶ Неизвестен минимальный показатель  $\alpha$  в числе операций  $\mathcal{O}(n^\alpha)$  (очевидно,  $\alpha \geq 2$ ).
- ▶ Алгоритм Штрассена не часто используется на практике. В недавней статье<sup>3</sup> (2016) утверждается, что алгоритм Штрассена может быть эффективен и для небольших матриц.

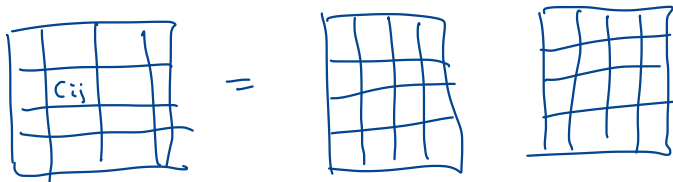
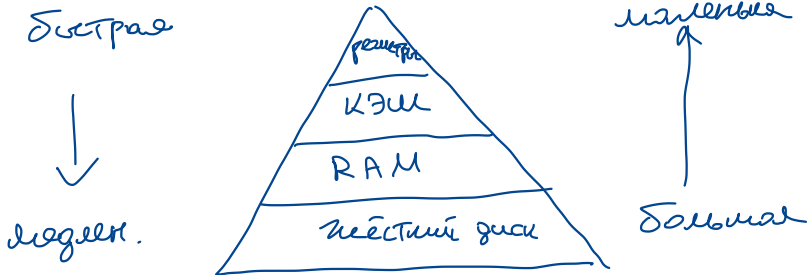
---

<sup>2</sup><https://arxiv.org/pdf/2010.05846.pdf>

<sup>3</sup><http://jianyuhuang.com/papers/sc16.pdf>



# Матричное умножение: иерархия памяти



$C_{ij} = A_{ik} B_{kj}$ , надо, чтобы  $C_{ij}, A_{ik}, B_{kj}$  помещались в КЭШ или регистры

# План

## Матричное умножение

Метод Штрассена

BLAS

## Устойчивость и обусловленность

Машинные числа

Вычислительная устойчивость

Обусловленность

# BLAS (Basic Linear Algebra Subprograms)

Оригинальная версия BLAS: 1979 год на fortran. С того времени переписан множество раз, но интерфейс функций стандартизован.

Разделяют 3 уровня операций в BLAS (далее  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$ ):

## Уровень 1

$\mathcal{O}(n)$  flops,  $\mathcal{O}(n)$  memops:

$$\frac{\text{flops}}{\text{memops}} = \text{const}$$

$$\text{(AXPY)} \quad y \leftarrow \alpha x + y, \quad x, y \in \mathbb{F}^n, \quad \alpha \in \mathbb{F},$$

## Уровень 2

$\mathcal{O}(n^2)$  flops,  $\mathcal{O}(n^2)$  memops:

$$\text{(MV)} \quad y \leftarrow \alpha Ax + \beta y, \quad x, y \in \mathbb{F}^n, \quad A \in \mathbb{F}^{n \times n}, \quad \alpha, \beta \in \mathbb{F},$$

обращение треугольных матриц, ранг-1 апдейт матрицы, и т.д.

## Уровень 3

$\mathcal{O}(n^3)$  flops,  $\mathcal{O}(n^2)$  memops:

$$\text{(MM)} \quad C \leftarrow \alpha AB + \beta C, \quad A, B, C \in \mathbb{F}^{n \times n}, \quad \alpha, \beta \in \mathbb{F}$$

Надо стараться записывать алгоритмы через матричное произведение.

# BLAS (Basic Linear Algebra Subprograms)

Оригинальная версия BLAS: 1979 год на fortran. С того времени переписан множество раз, но интерфейс функций стандартизован.

Разделяют 3 уровня операций в BLAS (далее  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$ ):

## Уровень 1

$\mathcal{O}(n)$  flops,  $\mathcal{O}(n)$  memops:

$$(AXPY) \quad y \leftarrow \alpha x + y, \quad x, y \in \mathbb{F}^n, \quad \alpha \in \mathbb{F},$$

## Уровень 2

$\mathcal{O}(n^2)$  flops,  $\mathcal{O}(n^2)$  memops:

$$\frac{\text{flops}}{\text{memops}} = \text{const}$$

$$(MV) \quad y \leftarrow \alpha Ax + \beta y, \quad x, y \in \mathbb{F}^n, \quad A \in \mathbb{F}^{n \times n}, \quad \alpha, \beta \in \mathbb{F},$$

обращение треугольных матриц, ранг-1 апдейт матрицы, и т.д.

## Уровень 3

$\mathcal{O}(n^3)$  flops,  $\mathcal{O}(n^2)$  memops:

$$(MM) \quad C \leftarrow \alpha AB + \beta C, \quad A, B, C \in \mathbb{F}^{n \times n}, \quad \alpha, \beta \in \mathbb{F}$$

Надо стараться записывать алгоритмы через матричное произведение.

# BLAS (Basic Linear Algebra Subprograms)

Оригинальная версия BLAS: 1979 год на fortran. С того времени переписан множество раз, но интерфейс функций стандартизован.

Разделяют 3 уровня операций в BLAS (далее  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$ ):

## Уровень 1

$\mathcal{O}(n)$  flops,  $\mathcal{O}(n)$  memops:

$$(AXPY) \quad y \leftarrow \alpha x + y, \quad x, y \in \mathbb{F}^n, \quad \alpha \in \mathbb{F},$$

## Уровень 2

$\mathcal{O}(n^2)$  flops,  $\mathcal{O}(n^2)$  memops:

$$(MV) \quad y \leftarrow \alpha Ax + \beta y, \quad x, y \in \mathbb{F}^n, \quad A \in \mathbb{F}^{n \times n}, \quad \alpha, \beta \in \mathbb{F},$$

обращение треугольных матриц, ранг-1 апдейт матрицы, и т.д.

## Уровень 3

$\mathcal{O}(n^3)$  flops,  $\mathcal{O}(n^2)$  memops:

*flops =  $\mathcal{O}(n^3)$  - хорошо, т.к.  
memops  $t_{mem} \gg t_{flop}$*

$$(MM) \quad C \leftarrow \alpha AB + \beta C, \quad A, B, C \in \mathbb{F}^{n \times n}, \quad \alpha, \beta \in \mathbb{F}$$

Надо стараться записывать алгоритмы через матричное произведение.

# BLAS (Basic Linear Algebra Subprograms)



## Названия операций

- ▶ DOT:  $x^T y$
- ▶ AXPY:  $y \leftarrow \alpha x + y$
- ▶ MV:  $y \leftarrow \alpha Ax + \beta y$
- ▶ MM:  $C \leftarrow \alpha AB + \beta C$
- ▶ R:  $A \leftarrow \alpha xy^T + A$   
(добавить ранг-1)
- ▶ ...

## Типы матриц

- ▶ GE – general
- ▶ GB – general band
- ▶ SY – symmetric
- ▶ SB – symm. band
- ▶ TR – triangular
- ▶ ...

## Precision:

- ▶ S – single
- ▶ D – double
- ▶ C – single complex
- ▶ Z – double complex

**Пример:** ZGEMM (матрично-матричное умножение с произвольными плотными матрицами из комплексных чисел в двойной точности)

# BLAS (Basic Linear Algebra Subprograms)

## Релевантные пакеты программ

- ▶ LAPACK (Linear Algebra PACKage): матричные факторизации, решение линейных систем, SVD, ... Использует BLAS.
- ▶ Intel MKL (Math Kernel Library): оптимизованные под Intel процессоры BLAS и LAPACK.
- ▶ OpenBLAS: оптимизированный BLAS. Базируется на GotoBLAS, который долгое время был рекордсменом (К. Гото написал GotoBLAS во время саббатикала в 2002).
- ▶ ATLAS (Automatically Tuned Linear Algebra Software): автоматически оптимизирует BLAS под конкретную систему.
- ▶ cuBLAS (CUDA BLAS): имплементация для GPU от NVIDIA.

Линейно-алгебраические операции в `scipy` и `numpy` – обертки для функций из BLAS и LAPACK. В Anaconda Python Distribution (версия 2.5 и старше) и MATLAB по умолчанию используется MKL.

# План

## Матричное умножение

Метод Штрассена

BLAS

## Устойчивость и обусловленность

Машинные числа

Вычислительная устойчивость

Обусловленность



# Машинные числа

мантисса

$$\mathbb{FP} = \left\{ \pm \left( \frac{d_1}{b} + \frac{d_2}{b^2} + \dots + \frac{d_m}{b^m} \right) b^e, \quad d_i = 0, 1, \dots, b-1, \quad e_{\min} \leq e \leq e_{\max} \right\}$$

- ▶  $\mathbb{FP} \subset \mathbb{R}$  – конечное множество машинных чисел
- ▶  $b \in \mathbb{N}$  – основание (base) арифметики
- ▶  $m \in \mathbb{N}$  – длина мантиссы
- ▶  $e \in \mathbb{Z}$  – порядок (экспонент) конкретного  $x \in \mathbb{FP}$
- ▶  $d_i \in \{0, \dots, b-1\}$  – разряды числа  $x \in \mathbb{FP}$

$b^{1-m}$  – иногда называют машинным эпитом или разряд округ. *машинным эпитом*



IEEE<sup>4</sup> стандарт 754

Был принят в 1985 для унификации представления чисел и операций с ними.

Точность	$b$	$m$	$e_{\max}$	$e_{\min}$
single	2	23	127	-126
double	2	52	1023	-1022

- ▶ Задаёт правило округления  $\text{fl}: \mathbb{R} \rightarrow \mathbb{FP}$ .
- ▶ Задаёт правило для разрешения неопределённостей. Например, для  $\frac{0}{0}$ .
- ▶ Арифметические и другие операции.
- ▶ ...

---

<sup>4</sup>Institute of Electrical and Electronics Engineers

# Машинные числа: округление

Для операции округления

$$\text{fl}: \mathbb{R} \rightarrow \mathbb{FP}$$

можно записать

$$\text{fl}(x) = x(1 + \epsilon), \quad |\epsilon| \leq \epsilon_{\text{machine}},$$

где  $\epsilon_{\text{machine}}$  (машинная эпсилон) — точная верхняя грань для  $|\epsilon|$ . То есть мы “привязали” определение  $\epsilon_{\text{machine}}$  к  $\text{fl}^5$ .

## Фундаментальная аксиома машинной арифметики

Для любых  $x, y \in \mathbb{FP}$  существует  $\epsilon : |\epsilon| \leq \epsilon_{\text{machine}}$ , такое что для любой операции  $\text{op} \in \{+, -, \times, /\}$  выполняется:

$$\text{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \epsilon), \quad |\epsilon| \leq \epsilon_{\text{machine}}$$

Для вычислений на компьютере, построенном по принципу этой аксиомы, будет удобно строить теоретический анализ ошибок округления.

---

<sup>5</sup>Для “школьного”  $\text{fl}$  — отбрасывания лишних цифр (truncation),  $\epsilon_{\text{machine}} = \frac{1}{2}b^{1-m}$ . В литературе также встречается  $\epsilon_{\text{machine}} = b^{1-m}$ .

# Устойчивость и обусловленность

Перейдем к обсуждению двух ключевых понятий численного анализа: обусловленность и устойчивость.

## Важно помнить

1. Устойчивость определяется для алгоритма.
2. Обусловленность определяется для задачи.

# План

## Матричное умножение

Метод Штрассена

BLAS

## Устойчивость и обусловленность

Машинные числа

Вычислительная устойчивость

Обусловленность

# Вычислительная устойчивость

конкретномер., нормир.

Пусть задана задача  $f: X \rightarrow Y$  и  $\tilde{f}: X \rightarrow Y$  – некоторый алгоритм ее решения.

## Прямая устойчивость

Алгоритм обладает свойством прямой устойчивости (forward stability), если

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \leq C \epsilon_{\text{machine}} = \mathcal{O}(\epsilon_{\text{machine}}).$$

Сложно анализировать (нужно следить за каждой операцией) и нужно учитывать, что при большой ошибке “плохой” может оказаться задача, а не алгоритм.

## Обратная устойчивость

Алгоритм обладает свойством обратной устойчивости, если

$$\tilde{f}(x) = f(\tilde{x}) \text{ для некоторого } \tilde{x}: \frac{\|\tilde{x} - x\|}{\|x\|} = \mathcal{O}(\epsilon_{\text{machine}}).$$

То есть, мы хотим заменить вычисленную величину как точное вычисление с возмущенными данными. Подход удобен для анализа (обратный анализ ошибок). Будем использовать в след. лекциях.

# Вычислительная устойчивость

## (Смешанная) устойчивость

Алгоритм является устойчивым, если

$$\frac{\|\tilde{f}(x) - f(\tilde{x})\|}{\|f(\tilde{x})\|} = \mathcal{O}(\epsilon_{\text{machine}}) \text{ для некоторого } \tilde{x}: \frac{\|\tilde{x} - x\|}{\|x\|} = \mathcal{O}(\epsilon_{\text{machine}}).$$

Пример

backward  
stable

⇐

$$\frac{|x+y - \underbrace{(x+y)}_{\tilde{x} + \tilde{y}}(1+\epsilon)|}{|x+y|} = |\epsilon| \leq \epsilon_{\text{machine}} \Rightarrow \text{forward stable}$$
$$\begin{aligned} \tilde{x} &= x(1+\epsilon) \\ \tilde{y} &= y(1+\epsilon) \end{aligned} \Rightarrow \frac{|\hat{x} - x|}{|x|} = |\epsilon| \leq \epsilon_{\text{machine}}$$



# План

## Матричное умножение

Метод Штрассена

BLAS

## Устойчивость и обусловленность

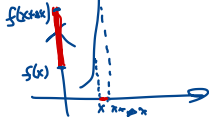
Машинные числа

Вычислительная устойчивость

Обусловленность

# Обусловленность

$$f: X \rightarrow Y$$



$$f(x + \Delta x) = f(x) + f'(x) \Delta x + O(\|\Delta x\|^2)$$

$$\frac{\|f(x + \Delta x) - f(x)\|}{\|f(x)\|} \approx \frac{\|f'(x)\|}{\|f(x)\|} \|\Delta x\|$$

$\| \text{def} \text{ } \text{много } \text{объяс. } f - \text{ } \text{не } \text{ра-} \text{убежит. } \text{за } \text{границ}$

$\text{Cond}(f, x)$

**Пример**

$$f(x) = Ax, \quad X = Y = \mathbb{C}^n$$

$$f'(x) = A$$

$$\text{Cond}_2(A) = \|A\| \|A^{-1}\|_2 = \| |I| \| = \|A A^{-1}\| \leq \|A\| \|A^{-1}\|$$

$$\text{Cond}(f, x) = \frac{\|A\|}{\|Ax\|} \|x\| = \frac{\|A\|}{\|Ax\|} \|A^{-1}Ax\| \leq \frac{\|A\|}{\|Ax\|} \|A^{-1}\| \|Ax\| = \|A\| \|A^{-1}\| = \text{Cond}(A)$$

*много объяс. матрицы A*

# Обусловленность

$$\begin{aligned} \text{forward-err} &= \frac{\| \tilde{f}(x) - f(x) \|}{\| f(x) \|} = \frac{\| f(\tilde{x}) - f(x) \|}{\| f(x) \|} \leq \\ &\leq \left( \text{cond}(f, x) + \underbrace{\mathcal{O}(\|\tilde{x} - x\|)}_{\mathcal{O}(\epsilon_{\text{machine}})} \right) \frac{\|\tilde{x} - x\|}{\|x\|} \approx \\ &\approx \text{cond}(f, x) \cdot \text{backward-err} \end{aligned}$$

## Литература

- ▶ N. Higham “Accuracy and Stability of Numerical Algorithms”, SIAM, 2002.
- ▶ Тыртышников, Е.Е. Матричный анализ и линейная алгебра, Москва, Физматлит, 2007. – 477 с
- ▶ Тыртышников Е. Е. Методы численного анализа. – Издательский центр Академия Москва, 2007. – 320 с.
- ▶ Trefethen, L. N., & Bau III, D. (1997). Numerical linear algebra. (Vol. 50). Siam. Philadelphia.