



Численные методы

Курс «Численные методы»

ВОЛКОВ Василий Михайлович, Минск, БГУ
v.volkov@tut.by

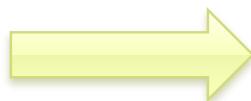
Минск, 4 сентября 2019

Зачем изучать численные методы?

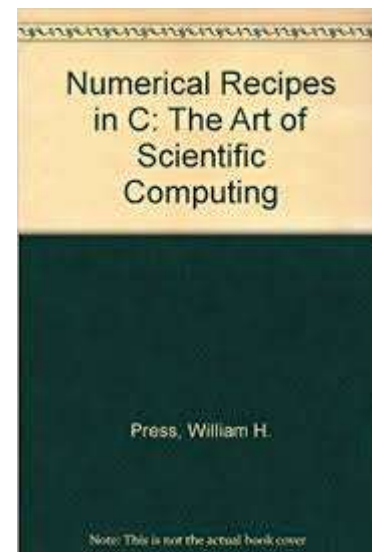
«Потому, что все оттенки смысла

Умное число передает...» Н.Гумилев 1921 г.

XXI век – век цифровых технологий



Численные методы – это цифровые технологии в математике, как языке современной науки.



Floating point operations per second (FLOPS). TOP 500



2016. Тяньхэ-2: 33.86(54.9) петафлопс (10^{15} flops)
2017 TaihuLight 93 petaflops

1.0 Компьютерная арифметика

Вычислительная погрешность

1.0.1 Формат чисел с плавающей запятой (floating point)

$$a = \pm d^s \left(\alpha_1 d^{-1} + \alpha_2 d^{-2} + \dots + \alpha_m d^{-m} \right),$$

МАНТИССА

$$\alpha_1 \neq 0, \quad 0 \leq \alpha_k < d$$

Например: $d=10$ (десятичная система), **0.31415e+001**

Длина мантиссы определяет относительную погрешность представления чисел: «машинное **ε**» минимальное число, которое, будучи прибавленным к единице, делает сумму отличной от единицы. Другими словами, это относительная погрешность представления чисел с плавающей запятой (точкой)

1.0.2 Форматы single и double. Стандарт IEEE 756

	длина	ϵ	∞
single	32 бит	$2^{-23} = 1.19e-7$	$>3.4028e+38$
double	64 бит	$2^{-52} = 2.22e-16$	$>1.7977e+308$

Особенности компьютерных вычислений (double)

Точная арифметика	Компьютерная арифметика
$2 + \epsilon - 2 = \epsilon$	$2 + \epsilon - 2 = 0$
$\sum_{n=1}^{\infty} \frac{1}{n} = \infty$	$\sum_{n=1}^{\infty} \frac{1}{n} < 49$
$10^{-10} \sum_{n=1}^{10^{10}} 0.2 - 0.2 = 0$	$10^{-10} \sum_{n=1}^{10^{10}} 0.2 - 0.2 = 3.2625e - 08$
$\sqrt{(10^{160})^2} = 10^{160}$	$\sqrt{(10^{160})^2} = Inf$

Вопросы

Что получится при вычислениях в классе double?:

$$\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n, \quad n = 10^{20}.$$



Лекция 1.

1.1 Нормы векторов и матриц. Оценка погрешности решения систем ЛАУ. Число обусловленности

Аксиомы нормы

Вектор $\mathbf{x} = (x_1, x_2, x_3, \dots, x_N)^T$, $\mathbf{x} \in R^N$

Норма вектора: $\|\cdot\| = R^N \rightarrow R$:

1. $\|\mathbf{x}\| \geq 0$, $\|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = 0$ – положительная определенность;
2. $\|\alpha\mathbf{x}\| = |\alpha| \|\mathbf{x}\|$, $\forall \alpha \in R$ – линейность при умножении на скаляр
3. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ – неравенство треугольника.

Примеры:

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{m=1}^N x_m^2} \quad \text{– Евклидова норма}$$

$$\|\mathbf{x}\|_\infty = \max_m |x_m| \quad \text{– максимальная норма}$$

$$1 / \sqrt{N} \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \quad \text{– Эквивалентность норм}$$

Матрицы и нормы матрицы

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1M} \\ a_{21} & a_{22} & \cdots & a_{2M} \\ \cdots & \cdots & a_{km} & \cdots \\ a_{K1} & a_{K2} & \cdots & a_{KM} \end{bmatrix}. \quad A \in R^{N \times N}.$$

Норма матрицы A , подчиненная векторной норме $\|\cdot\|$, определяется числом

$$\|A\| = \sup_{\mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \sup_{\|\mathbf{x}\|=1} \|A\mathbf{x}\|.$$

Подчиненные нормы обладают свойством мультипликативности:

$$\|A\mathbf{x}\| \leq \|A\| \cdot \|\mathbf{x}\|$$

Примеры матричных норм

$$\|A\|_{\infty} = \max_k \left\{ \sum_{m=1}^N |a_{km}| \right\}$$

$$\|A\|_2 = \max \{ \lambda(AA^T)^{1/2} \}$$

$\max \{ \lambda(AA^T) \}$ Максимальное сингулярное число

Оценка погрешности решения систем ЛАУ

Система линейных алгебраических уравнений в матричном виде

$$\mathbf{A}\mathbf{x} = \mathbf{f}. \quad \sum_{m=1}^N a_{km}x_m = f_k, \quad k = \overline{1, N}.$$

$$\mathbf{A} \in \mathbf{R}^{N \times N}; \quad \mathbf{x}, \mathbf{f} \in \mathbf{R}^N,$$

Пусть найдено приближенное решение системы

$$\tilde{\mathbf{x}} = \mathbf{x} + \delta\mathbf{x}$$

Приближенное решение исходной системы является точным решением некоторой **возмущенной** системы

$$\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{f} + \delta\mathbf{f}$$

$$\delta\mathbf{f} = \mathbf{A}(\mathbf{x} + \delta\mathbf{x}) - \mathbf{f} \quad - \text{невязка.}$$

Оценка погрешности возмущенной задачи

$$A\delta\mathbf{x} = \delta\mathbf{f} \Rightarrow \|\delta\mathbf{x}\| = \|A^{-1}\delta\mathbf{f}\| \Rightarrow \|\delta\mathbf{x}\| \leq \|A^{-1}\| \cdot \|\delta\mathbf{f}\|$$

Проделявая то же самое с исходной системой будем иметь

$$A\mathbf{x} = \mathbf{f} \Rightarrow \|A\| \cdot \|\mathbf{x}\| \geq \|\mathbf{f}\| \Rightarrow \|\mathbf{x}\| \geq \|A\|^{-1} \cdot \|\mathbf{f}\|$$

Деление первого неравенства на второе дает оценку

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta\mathbf{f}\|}{\|\mathbf{f}\|}$$

Число обусловленности матрицы

$$M_A = \text{cond}(A) = \|A\| \cdot \|A^{-1}\|$$

Геометрический смысл плохой обусловленности

$$\begin{cases} a_{11}x_1 + a_{12}x_2 = f_1, \\ a_{21}x_1 + a_{22}x_2 = f_2. \end{cases}$$

$$\tan(\alpha_1) = -\frac{a_{11}}{a_{12}}, \quad \tan(\alpha_2) = -\frac{a_{21}}{a_{22}}$$

Определитель матрицы

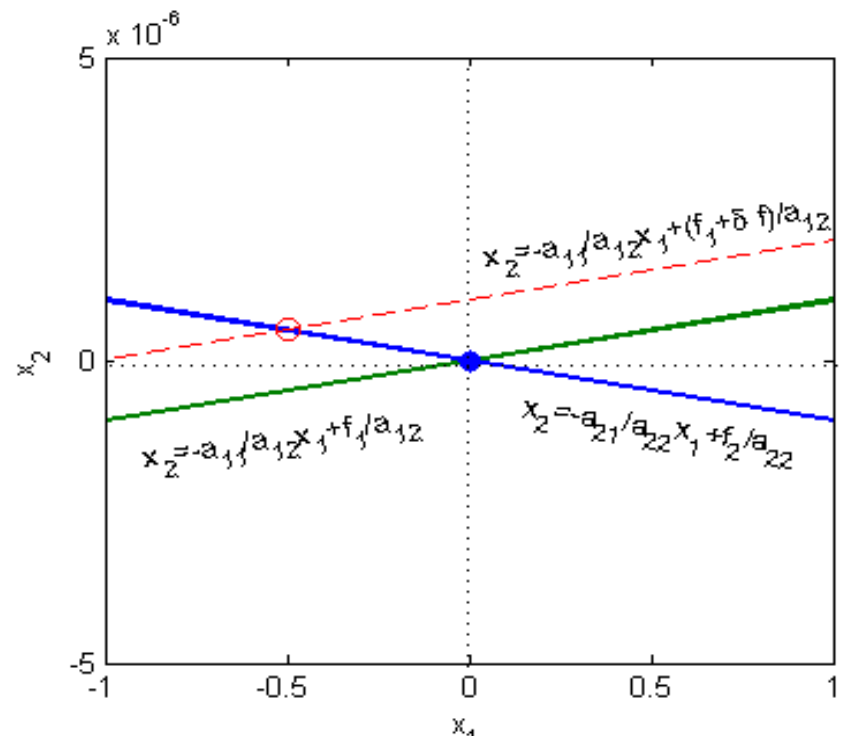
$$\det(A) = a_{11}a_{22} - a_{12}a_{21}.$$

$$\det(A) = a_{11}a_{22} - a_{12}a_{21} = 0 \iff \tan(\alpha_1) = \tan(\alpha_2)$$

$$a_{11} = -a_{21} = -10^{-6},$$

$$a_{22} = a_{12} = 0,$$

$$f_1 = f_2 = 1$$

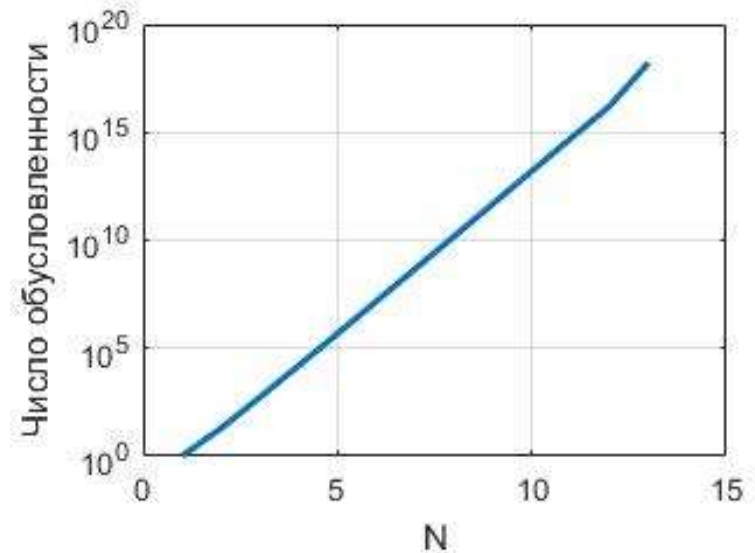


Пример плохо обусловленной матрицы

Матрица Гильберта

$$H_n = \begin{pmatrix} 1 & 1/2 & 1/3 & \dots & 1/n \\ 1/2 & 1/3 & 1/4 & \dots & 1/(n+1) \\ 1/3 & 1/4 & 1/5 & \dots & 1/(n+2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1/n & 1/(n+1) & 1/(n+2) & \dots & 1/(2n-1) \end{pmatrix}$$

Зависимость числа обусловленности от размерности матрицы Гильберта



Решение системы ЛАУ с матрицей Гильберта 14x14 имеем предупреждение :

Warning: Matrix is close to singular or badly scaled. Results may be inaccurate. RCOND = 1.602620e-18.

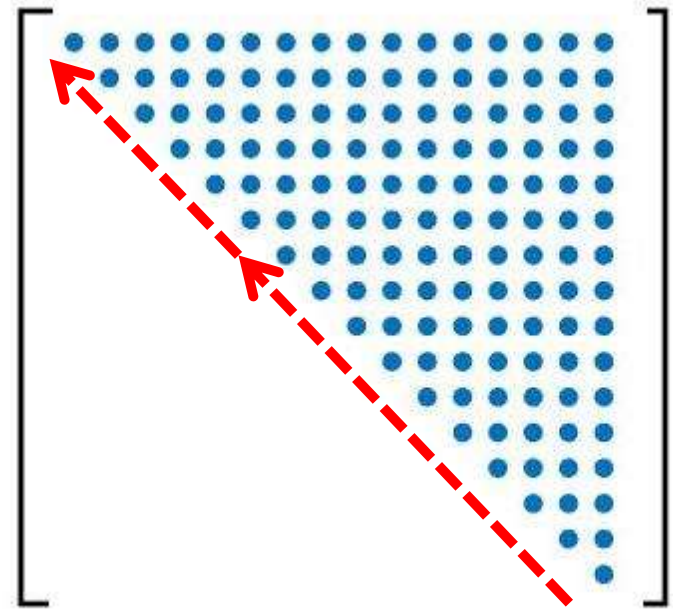
Метод Гаусса

Предпосылкой для метода Гаусса является тот факт, что для системы с треугольной матрицей решение находится относительно просто.

$$a_{k \leq m} \neq 0, \quad a_{k > m} = 0$$

$$\sum_{m=k}^N a_{km} x_m = f_k$$

$$x_N = f_N / a_{NN}$$



Обратный ход метода Гаусса

$$x_k = \left(f_k - \sum_{m=k+1}^N a_{km} x_m \right) a_{kk}^{-1}, \quad k = N-1, N-2, \dots, 2, 1$$

Приведение матрицы к треугольному виду

Приведение матрицы системы к треугольному виду основано на элементарных преобразованиях:

- *Умножение строки на ненулевое число (масштабирование строки);*
- *Перестановка строк;*
- *Замена строки на ее линейную комбинацию с другими строками.*

Каждое из элементарных преобразований можно выразить в виде умножения левой и правой части системы на соответствующую (квадратную не вырожденную) матрицу преобразований. Матрицы элементарных преобразований имеют ту же размерность , что и матрица системы.

Масштабирование и перестановка строк.

Элементарная матрица масштабирования k -й строки получается из единичной матрицы, в которой единичному элементу $a_{kk} = 1$ присваивается значение масштабирующего множителя $a_{kk} = s$.
Произведение произвольного числа элементарных матриц масштабирования тоже является матрицей масштабирования.

Элементарная матрица перестановок k -й и m -й строк получается из единичной матрицы путем перестановки в ней соответствующих (k -й и m -й) строк.

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 5 & 9 & 13 \\ 2 & 6 & 10 & 14 \\ 3 & 7 & 11 & 15 \\ 4 & 8 & 12 & 16 \end{bmatrix} = \begin{bmatrix} 1 & 5 & 9 & 13 \\ 4 & 12 & 20 & 2 \\ 3 & 7 & 11 & 15 \\ 4 & 8 & 12 & 16 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 1 & 5 & 9 & 13 \\ 2 & 6 & 10 & 14 \\ 3 & 7 & 11 & 15 \\ 4 & 8 & 12 & 16 \end{bmatrix} = \begin{bmatrix} 1 & 5 & 9 & 13 \\ 4 & 8 & 12 & 16 \\ 3 & 7 & 11 & 15 \\ 2 & 6 & 10 & 14 \end{bmatrix}$$

Элементарная треугольная матрица

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2N} \\ \cdots & \cdots & a_{km} & \cdots \\ a_{N1} & a_{N2} & \cdots & a_{NN} \end{bmatrix}, \quad L_{(k)} = \begin{bmatrix} 1 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & \ddots & & 0 & \cdots & 0 \\ 0 & 0 & 1/a_{kk} & 0 & \cdots & 0 \\ 0 & 0 & -a_{k+1,k}/a_{kk} & 1 & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \ddots & 0 \\ 0 & \cdots & -a_{N,k}/a_{kk} & \cdots & 0 & 1 \end{bmatrix}.$$

При умножении на матрицу $L_{(k)}$ элементы k -го столбца ниже диагонали обращаются в нуль, а элементы левее этого столбца остаются неизменными.

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ -2 & 1 & 0 & 0 \\ -3 & 0 & 1 & 0 \\ -4 & 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 5 & 9 & 13 \\ 2 & 6 & 10 & 14 \\ 3 & 7 & 11 & 15 \\ 4 & 8 & 12 & 16 \end{bmatrix} = \begin{bmatrix} 1 & 5 & 9 & 13 \\ 0 & -4 & -8 & -12 \\ 0 & -8 & -16 & -24 \\ 0 & -12 & -24 & -36 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -1/4 & 0 & 0 \\ 0 & -2 & 1 & 0 \\ 0 & -3 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 5 & 9 & 13 \\ 0 & -4 & -8 & -12 \\ 0 & -8 & -16 & -24 \\ 0 & -12 & -24 & -36 \end{bmatrix} = \begin{bmatrix} 1 & 5 & 9 & 13 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Приведение матрицы к треугольному виду (Прямой ход метода Гаусса)

$$A_{(1)} = L_{(1)}A,$$

$$A_{(2)} = L_{(2)}A_{(1)},$$

.....

$$A_{(N)} = L_{(N)}A_{(N-1)}$$

Приведение матрицы к треугольному виду сводится к N матричным умножениям. Каждое такое умножение требует порядка $O(N^2)$ арифметических операций. Общая вычислительная сложность алгоритма $O(N^3)$. Для обратного хода метода Гаусса несложно показать, что его вычислительная сложность $O(N^2)$, что на порядок ниже, чем у прямого хода.

Выбор главного элемента

Элемент a_{kk} , определяющий элемент $l_{kk} = 1 / a_{kk}$ элементарной треугольной матрицы $L_{(k)}$, принято называть **главным**. Очевидно, что если $a_{kk} = 0$, то метод Гаусса становится некорректным (деление на ноль). Избежать данной ситуации можно с помощью перестановки строк таким образом, чтобы $|a_{m>k,k}| \leq |a_{k,k}|$. Выбор главного элемента не только позволяет избежать деления на ноль, но и предотвращает деление на малое число, близкое к нулю, что препятствует катастрофическому возрастанию абсолютной погрешности. Если оказывается, что главный элемент и все элементы в столбце ниже главного равны нулю,

$$\max_m \{a_{m \leq k, k}\} = 0.$$

то это говорит о том, что преобразуемая матрица **вырожденная**.



СПАСИБО ЗА ВНИМАНИЕ!