

Прямые методы решения систем
линейных алгебраических ур-ий.

$$A = \{a_{ij}\}_{i,j=1}^n; \quad \vec{f} = (f_1, f_2, \dots, f_n);$$

$$\vec{u} = (u_1, u_2, \dots, u_n);$$

$$(1) \quad A \vec{u} = \vec{f};$$

$$(1') \quad \sum_{j=1}^n a_{ij} u_j = f_i; \quad i = \overline{1, n};$$

Метод Гаусса. Прямой ход:

$$a_{i,1} \left\{ \begin{array}{l} u_1 + \sum_{j=2}^n a_{1,j}^1 u_j = f_1^1 \quad \leftarrow \\ \sum_{j=2}^n a_{i,j}^1 u_j = f_i^1 \quad (i=\overline{2,n}). \end{array} \right.$$

$$\Downarrow \left\{ \begin{array}{l} u_2 + \sum_{j=3}^n a_{2,j}^2 u_j = f_2^2 \quad \leftarrow \\ \sum_{j=3}^n a_{i,j}^2 u_j = f_i^2 \quad (i=\overline{3,n}) \end{array} \right.$$

$$(2) \left\{ \begin{array}{l} u_i + \sum_{j=i+1}^n a_{i,j}^i u_j = f_i^i \quad (i=\overline{1,n-1}). \\ u_n = f_n^n \end{array} \right.$$

Обратный ход.

$$(3) \begin{cases} u_n = f_n^n; \\ u_i = f_i^i - \sum_{j=i+1}^n a_{ij}^i u_j; \quad i = n-1, n-2, \dots, 1. \end{cases}$$

Теорема 1. Метод Гаусса (2), (3) может быть реализован

\Leftrightarrow все главные миноры матрицы A отличны от 0 \oplus

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}$$

Метод прогонки.

$$(4) \quad \begin{cases} B_1 u_1 - C_1 u_2 = f_1 \\ -A_i u_{i-1} + B_i u_i - C_i u_{i+1} = f_i \quad (i = \overline{2, n-1}) \\ -A_n u_{n-1} + B_n u_n = f_n. \end{cases}$$

$$\begin{pmatrix} B_1 & -C_1 & & & & 0 \\ -A_2 & B_2 & -C_2 & & & \\ 0 & -A_3 & & & & \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & & & 0 & -A_n & B_n \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ \vdots \\ u_n \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ \vdots \\ f_n \end{pmatrix}$$

Метод Гаусса для системы (4) = метод прогонки.

$$u_1 = \frac{c_1}{B_1} u_2 + \frac{f_1}{B_1} \equiv d_1 u_2 + \beta_1;$$

$$-A_2(\underline{d_1 u_2 + \beta_1}) + B_2 u_2 - C_2 u_3 = f_2$$

$$u_2 = \frac{C_2}{B_2 - A_2 d_1} \cdot u_3 + \frac{f_2 + \beta_1 A_2}{B_2 - A_2 d_1} \equiv d_2 u_3 + \beta_2;$$

Пусть $u_{i-1} = d_{i-1} u_i + \beta_{i-1}$; Находим " u_i ".

$$-A_i(\underline{d_{i-1} u_i + \beta_{i-1}}) + B_i u_i - C_i u_{i+1} = f_i;$$

$$u_i = \frac{C_i}{B_i - A_i d_{i-1}} \cdot u_{i+1} + \frac{f_i + \beta_{i-1} A_i}{B_i - A_i d_{i-1}} \equiv d_i u_{i+1} + \beta_i;$$

u_n найдем из системы ур-ий:

$$\begin{cases} u_{n-1} - \alpha_{n-1} u_n = \beta_{n-1}; \\ -A_n u_{n-1} + B_n u_n = f_n. \end{cases}$$

\Downarrow

$$u_n = \frac{f_n + \beta_{n-1} A_n}{B_n - A_n \alpha_{n-1}};$$

Стандартная (монотонная) прогонка:

I. $\alpha_1 = \frac{C_1}{B_1}; \quad \beta_1 = \frac{f_1}{B_1};$

(5) $\alpha_i = \frac{C_i}{B_i - A_i \alpha_{i-1}}; \quad \beta_i = \frac{f_i + \beta_{i-1} A_i}{B_i - A_i \alpha_{i-1}}; \quad i = \overline{2, n-1};$

α_i, β_i — "прогоночные коэффициенты".

$$\text{II. } u_n = \frac{f_n + \beta_{n-1} A_n}{B_n - A_n d_{n-1}};$$

$$(6) \quad u_i = d_i u_{i+1} + \beta_i, \quad i = n-1, n-2, \dots, 1.$$

Определение 1. ① Прогонка (5), (6) наз. "корректной",
если $B_1 \neq 0$; $B_i - A_i d_{i-1} \neq 0$ ($i = \overline{2, n-1}$).

② Прогонка (5), (6) наз. "устойчивой", если;

$$|d_i| \leq 1 \quad (i = \overline{1, n-1}) \quad \oplus$$

Пусть d_i, β_i ($i = \overline{1, n-1}$) - вычисл. точно; $\bar{u}_n = u_n + \varepsilon_n$;

$$(6) \Rightarrow \bar{u}_i = u_i + \varepsilon_i; \quad (6) \Rightarrow \varepsilon_i = d_i \varepsilon_{i+1}; \quad i = n-1, n-2, \dots$$

$$|\varepsilon_i| \leq |d_i| \cdot |\varepsilon_{i+1}| \leq |\varepsilon_{i+1}|.$$

Теорема 2. Рассмотрим систему ур-ний (4). Пусть:

① $B_1 \neq 0, B_n \neq 0; A_i \neq 0, C_i \neq 0 \quad (i = \overline{2, n-1})$.

② $|B_1| \geq |C_1|, |B_i| \geq |A_i| + |C_i| \quad (i = \overline{2, n-1});$

$$|B_n| \geq |A_n|.$$

\Downarrow и хотя бы одно из этих нер-в строгое.

стандартная прогонка (5), (6) — корректна и устойчива.

Доказательство.

II. Предположим, что:

$$B_j - A_j d_{j-1} \neq 0 \text{ и } |d_j| \leq 1 \text{ для } j = 1, 2, \dots, i-1$$

при некотором фиксированном $i \leq n-1$.

Докажем, что: $B_i - A_i d_{i-1} \neq 0$ и $|d_i| \leq 1$.

$$|B_i - A_i d_{i-1}| \geq |B_i| - |A_i| \cdot |d_{i-1}| \stackrel{\textcircled{2}}{\geq} |A_i| + |C_i| - |A_i| \cdot |d_{i-1}|$$

$$= (1 - |d_{i-1}|) \cdot |A_i| + |C_i| \geq |C_i| \stackrel{\textcircled{1}}{>} 0.$$

\Downarrow $B_i - A_i d_{i-1} \neq 0$ \forall

$$|d_i| = \frac{|C_i|}{|B_i - A_i d_{i-1}|} \leq \frac{|C_i|}{|C_i|} = 1 \quad \textcircled{+}$$

Осталось показать: $B_n - A_n d_{n-1} \neq 0$:

$$|B_n - A_n d_{n-1}| \geq |B_n| - |A_n| \cdot \underbrace{|d_{n-1}|}_{\leq 1} \geq |B_n| - |A_n| > 0. \quad (+)$$

Процесс Дюрра.

$$\begin{cases} \alpha_i = \frac{C_i}{B_i - A_i d_{i-1}}; & \beta_i = \frac{f_i + A_i \beta_{i-1}}{B_i - A_i d_{i-1}}; \end{cases}$$

$$u_i = d_i \cdot u_{i+1} + \beta_i;$$

$$d_i \rightarrow 1 - d_i.$$

$$\begin{aligned} 1 - d_i &= \frac{C_i}{B_i - A_i(1 - d_{i-1})} \Rightarrow d_i = 1 - \frac{C_i}{B_i - A_i + A_i d_{i-1}} = \\ &= \frac{B_i - A_i - C_i + A_i d_{i-1}}{B_i - A_i + A_i d_{i-1}} = \frac{S_i + A_i d_{i-1}}{C_i + S_i + A_i d_{i-1}}; \end{aligned}$$

$$\text{где } S_1 \equiv B_1 - C_1; \quad S_i \equiv B_i - A_i - C_i \quad (i = \overline{2, n-1});$$

$$S_n \equiv B_n - A_n.$$

Обычно выполняются: $S_i \geq 0 \quad (i = \overline{1, n}); \quad C_i > 0; \quad A_i > 0$
 $(i = \overline{2, n-1})$

В этом случае знаменатели в ф-лах прогонки Форра содержат суммы положительных величин.

Формулы прогонки Форра:

$$\text{I.} \quad \alpha_1 = \frac{C_1}{B_1} = \frac{C_1}{S_1 + C_1}; \quad \beta_1 = \frac{f_1}{C_1 + S_1};$$

$$\alpha_i = \frac{S_i + A_i \alpha_{i-1}}{C_i + S_i + A_i \alpha_{i-1}}; \quad \beta_i = \frac{f_i + A_i \beta_{i-1}}{C_i + S_i + A_i \alpha_{i-1}}; \quad i = \overline{2, n-1}$$

$$\text{II. } u_n = \frac{f_n + A_n d_{n-1}}{S_n + A_n d_{n-1}} ;$$

$$u_i = (1 - d_i) \cdot u_{i+1} + \beta_i ; i = n-1, n-2, \dots, 1. \oplus$$

Задача. Доказать, что число операций умножения и деления в прогонках "стандартная" и "Дорра" имеет порядок $O(n)$ при $n \rightarrow \infty$.

Пример
$$\begin{cases} u_1 = \varphi, \\ u_{i-1} + u_i - u_{i+1} = 0, i = \overline{2, n-1}. \\ u_n = \psi. \end{cases}$$

Решение:
$$u_i = \frac{\varphi \cdot \sin \frac{\pi(n-i)}{3} + \psi \cdot \sin \frac{\pi(i-1)}{3}}{\sin \frac{\pi(n-1)}{3}} ; n \neq 3k+1.$$

Вычисляем $d_i = \frac{C_i}{B_i - A_i d_{i-1}}$;

$$d_1 = \frac{C_1}{B_1} = 1; \quad d_2 = \frac{C_2}{B_2 - A_2 d_1} = \frac{1}{1 - 1 \cdot 1} = \frac{1}{0} (?)$$

$$B_i - A_i - C_i = 1 - 1 - 1 = -1 < 0; \quad i = \overline{2, n-1};$$

↑ условия (2) диагонального преобладания
не выполнены.

Погрешности округления.

$$(1) \quad q = (-1)^S \cdot M \cdot z^P;$$

$z \in \mathbb{Z} \cap (1, +\infty)$; основание системы счисления;

$P \in \mathbb{Z} \cap (-P_*, P^*)$ - порядок числа "q"; P_*, P^* - целые
и положительные (границы порядка)

$M = \sum_{i=1}^t q_i z^{-i}$ - мантисса числа "q";

$q_i \in \mathbb{Z} \cap [0, z-1]$ - разряды мантиссы; $i = \overline{1, t}$;

t - целое и положительное (разрядность)

$q_1 \neq 0$ - для нормализованного представления.

$S = 0 \vee 1$ - знак числа q

Q^+ — положительная компонента мн-ва (1).

q_{\min} и q_{\max} — во мн-ве Q^+ :

$$q = z^p \cdot \left(\sum_{i=1}^t q_i z^{-i} \right)$$

$$-p_* < p < p^* ; \quad 0 \leq q_i \leq z-1 ;$$

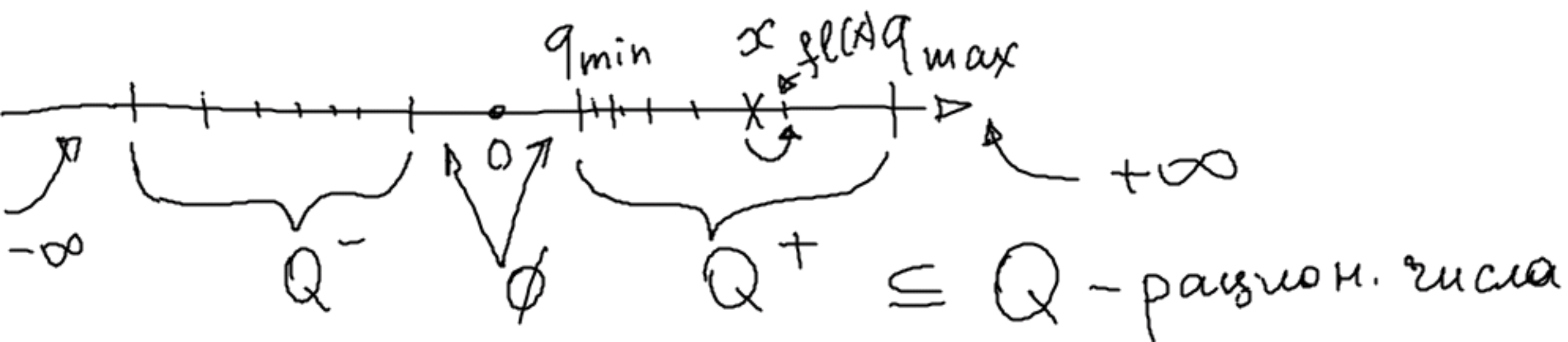
$$q_{\min} : \quad p = -p_* ; \quad q_1 = 1 ; \quad q_i = 0 \quad (i = \overline{2, t})$$

$$q_{\min} = z^{-p_*} \cdot z^{-1} = \underline{\underline{z^{-(p_*+1)}}}$$

$$q_{\max} : \quad p = p^* ; \quad q_i = z-1 ; \quad i = \overline{1, t} ;$$

$$q_{\max} = z^{p^*}$$

$$q_{\max} = z^{p^*} \cdot \left(\sum_{i=1}^t z^{-i} \right) (z-1) = \frac{z^{p^*} \cdot (1 - z^{-t})}{z-1}$$
$$\sum_{i=1}^t \left(\frac{1}{z} \right)^i = \frac{1 - \left(\frac{1}{z} \right)^{t+1}}{1 - 1/z} - 1 = \frac{1/z - (1/z)^{t+1}}{1 - 1/z} = \frac{1 - z^{-t}}{z-1} ;$$



$$x \in \mathbb{R} \Rightarrow \textcircled{1} \quad |x| \geq q_{\max} \Rightarrow x = +\infty.$$

$$\textcircled{2} \quad |x| \leq q_{\min} \Rightarrow x = \emptyset;$$

$$\textcircled{3} \quad q_{\min} < |x| < q_{\max} \Rightarrow$$

$$\Rightarrow x \rightarrow fl(x) \in Q^- \cup Q^+$$

$$\frac{fl(x) - x}{x}$$

- относительная погрешность округления.

Теорема 1. $q_{\min} < |x| < q_{\max} \Rightarrow \left| \frac{fl(x) - x}{x} \right| \leq \frac{\tau^{1-t}}{2};$

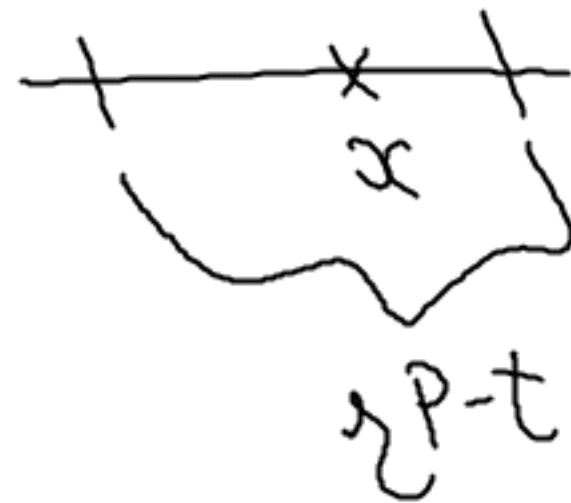
Док-во. $q_{\min} < x < q_{\max}$ - считаем, без огр. общности.

$$q = \tau^p \left(\sum_{i=1}^t q_i \tau^{-i} \right) \in [\tau^{p-1}, \tau^p];$$

$x \in [\tau^{p-1}, \tau^p]$ с некоторым " p "

$$\text{Т.к. } [q_{\min}, q_{\max}] = \bigcup_{p=-p_*}^{p_*} [\tau^{p-1}, \tau^p];$$

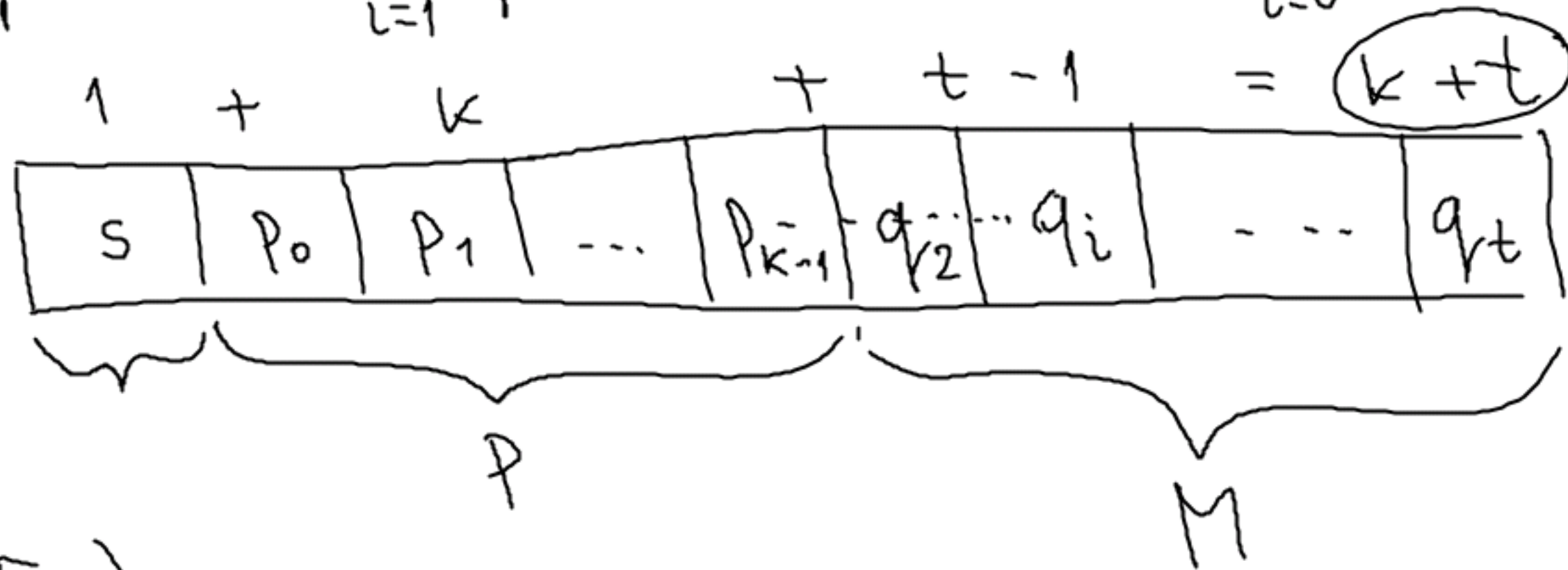
$$\left\{ \begin{array}{l} |fl(x) - x| \leq \frac{\tau^{p-t}}{2}; \\ |x| \geq \tau^{p-1}. \end{array} \right.$$



$$\frac{|fl(x) - x|}{|x|} \leq \frac{\tau^{p-t}}{2 \tau^{p-1}} = \frac{\tau^{1-t}}{2} \quad (+)$$

Форматы представления вещ. чисел в памяти комп.

$$q = (-1)^s \left(\sum_{i=1}^t q_i 2^{-i} \right) 2^p; \quad p = \sum_{i=0}^{k-1} p_i 2^i - p_*$$



(E1): $k=8, t=24, p_*=126$; (одинарная точность)

(E2) $k=11, t=53, p_*=1022$; (двойная точность)

$$p^* = 2^k - p_* - 1;$$

Для (E1): $q_{\min} = 2^{-(p_*+1)} = 2^{-127} \approx 10^{-38}$; — мин. кол.

машин. бескон.: $q_{\max} = 2^{p^*} = 2^{129} \approx 10^{39}$;

машин. "ε" : $\varepsilon = 2^{-t} = 2^{-24} \approx 10^{-7}$;

Для (E2): $q_{\min} = 2^{-1023} \approx 10^{-308}$; (\emptyset)

$q_{\max} = 2^{1025} \approx 10^{309}$; (∞)

$\varepsilon = 2^{-53} \approx 10^{-16}$;

$x, y \in \mathbb{R} \Rightarrow fl(x), fl(y) \in \mathbb{Q}^- \cup \mathbb{Q}^+ \Rightarrow$

$\Rightarrow fl(x) \oplus fl(y) \Rightarrow fl[fl(x) \oplus fl(y)]$;

$p_1 + p_2 \geq p^*$ — "переполнение".

$-p_1 - p_2 \leq -(p_* + 1)$ — "потеря порядка".

\oplus - операция "+" и "-"

$$x, y \in \mathbb{R}; \quad \frac{f_l(x) - x}{x} = \varepsilon_x \Rightarrow \underline{f_l(x) = (\varepsilon_x + 1) \cdot x};$$

$$\frac{f_l(y) - y}{y} = \varepsilon_y \Rightarrow \underline{f_l(y) = (\varepsilon_y + 1) \cdot y};$$

$$\frac{f_l[f_l(x) \oplus f_l(y)] - f_l(x) \oplus f_l(y)}{f_l(x) \oplus f_l(y)} = \varepsilon$$

$$\underline{f_l[f_l(x) \oplus f_l(y)] = (\varepsilon + 1) \cdot [f_l(x) \oplus f_l(y)]};$$

$$\frac{f_l[f_l(x) \oplus f_l(y)] - (x \oplus y)}{x \oplus y} = \frac{(\varepsilon + 1) [(\varepsilon_x + 1)x \oplus (\varepsilon_y + 1)y] - x \oplus y}{x \oplus y} =$$

$$= \varepsilon + (\varepsilon + 1) \frac{\varepsilon_x \cdot x \oplus \varepsilon_y \cdot y}{x \oplus y} \equiv E:$$

$$\underline{\underline{\oplus = +}}; \quad \varepsilon, \varepsilon_x, \varepsilon_y \approx \frac{z^{1-t}}{2};$$

$$|E| \leq \frac{z^{1-t}}{2} + \left(1 + \frac{z^{1-t}}{2}\right) \cdot \frac{z^{1-t}}{2} \approx \frac{z^{1-t}}{2} \quad (t \rightarrow \infty)$$

$$\underline{\underline{\oplus = -}}; \quad \frac{\varepsilon_x \cdot x - \varepsilon_y \cdot y}{x - y} \approx E$$

проблема: $x - y$ - малое.

Пример. $z^2 - 200z + 1 = 0$ - найти наименьший корень.

$$\Rightarrow z = 100 - \sqrt{9999};$$

$$\boxed{z=10 \text{ и } t=4;}$$

$$\begin{cases} x = 100; \\ y = \sqrt{9999} = \underline{99,99499987} \dots; \\ x - y = 0,00500012 \dots \end{cases}$$

$$\begin{cases} fl(x) = 100; \\ fl(y) = 99,99; \\ fl(x) - fl(y) = 0,01; \\ fl[fl(x) - fl(y)] = 0,01; \end{cases}$$

$$E = \frac{0,01 - 0,00500012 \dots}{0,00500012 \dots} = 0,99992499 \dots \approx 1;$$

$$\varepsilon_x = 0; \varepsilon_y = 5 \cdot 10^{-5}$$

$$x - y = 5 \cdot 10^{-3}; y \approx 100$$

$$z = \frac{1}{100 + \sqrt{9999}};$$

$$f_l(x) + f_l(y) = 199,99;$$

$$f_l[f_l(x) + f_l(y)] = 200;$$

$$\boxed{z = 0,005;}$$