

Санкт-Петербургский государственный университет

ТАМАСЯН Григорий Шаликович
ПРОСОЛУПОВ Евгений Викторович

АНАЛИЗ АЛГОРИТМОВ ПРОЕКТИРОВАНИЯ ТОЧКИ
НА СТАНДАРТНЫЙ СИМПЛЕКС

III международная конференция
«УСТОЙЧИВОСТЬ И ПРОЦЕССЫ УПРАВЛЕНИЯ»
посвященная 85-летию со дня рождения профессора,
чл.-корр. РАН В. И. Зубова

г. Санкт-Петербург
5 – 9 октября 2015 г.

СОДЕРЖАНИЕ ДОКЛАДА

- Постановка задачи
- Краткая история вопроса
- Алгоритм Малозёмова – Певного
- Алгоритм Maculan – de Paula
- Численные эксперименты

ПОСТАНОВКА ЗАДАЧИ

Задача ортогонального проектирования точки $c = (c_1, \dots, c_n)$ на стандартный симплекс $\Lambda \subset \mathbb{R}^n$, определяемый условиями

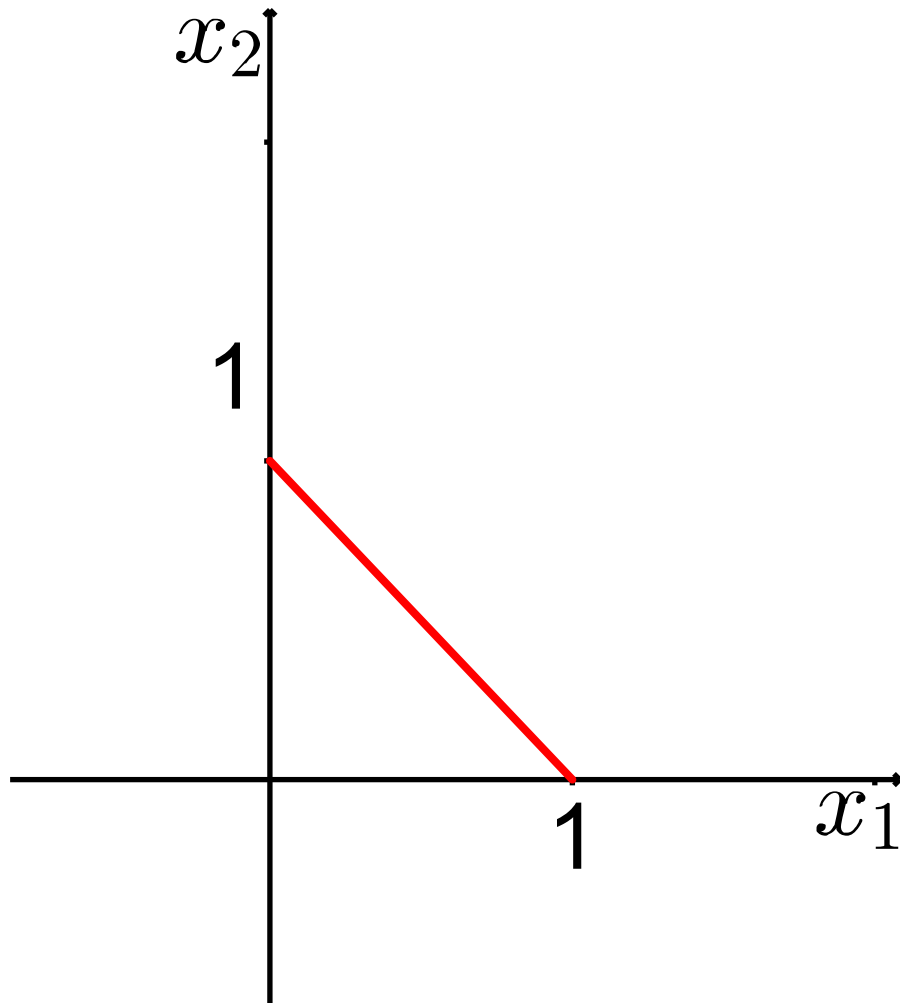
$$\sum_{i=1}^n x_i = 1; \quad x_i \geq 0, \quad i \in 1 : n,$$

ставится следующим образом:

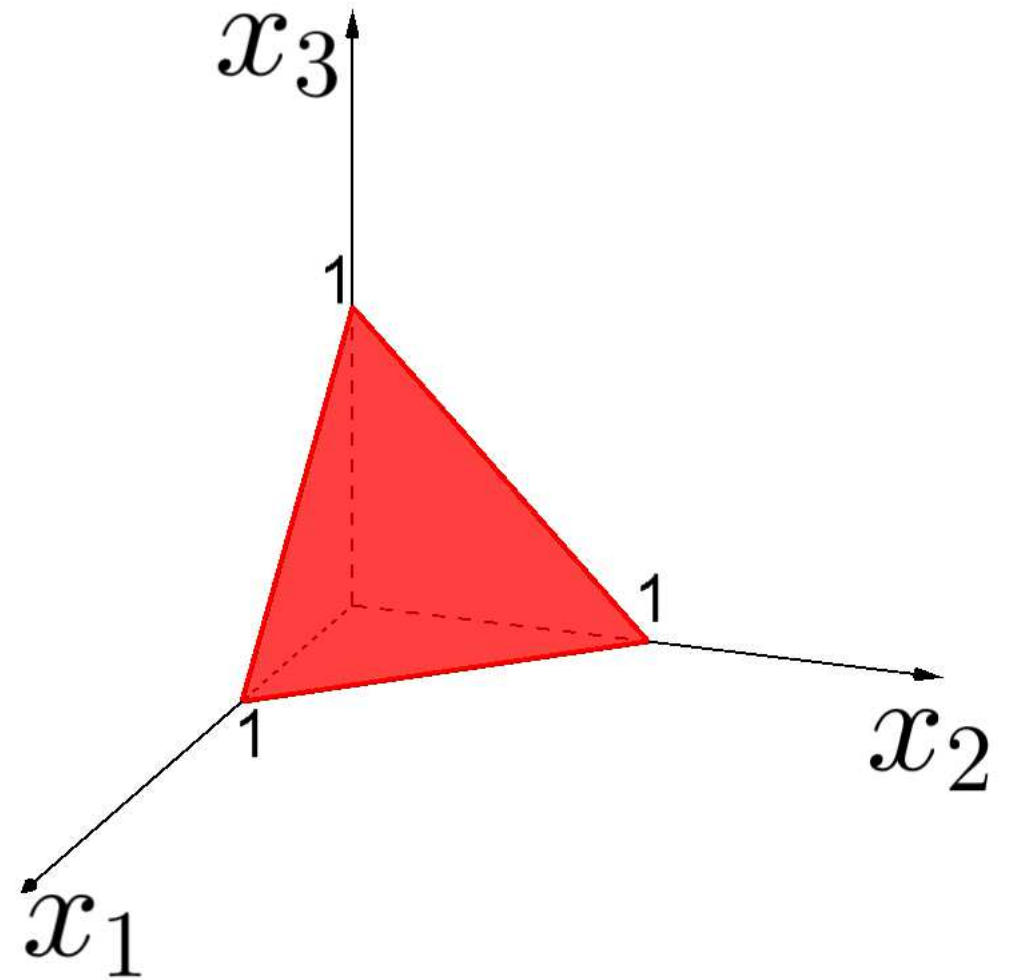
$$Q(x) := \frac{1}{2} \sum_{i=1}^n (x_i - c_i)^2 \rightarrow \min_{x \in \Lambda}. \quad (1)$$

Решение этой задачи существует и единственно. Обозначим его x^* .

СТАНДАРТНЫЙ СИМПЛЕКС



$n = 2$

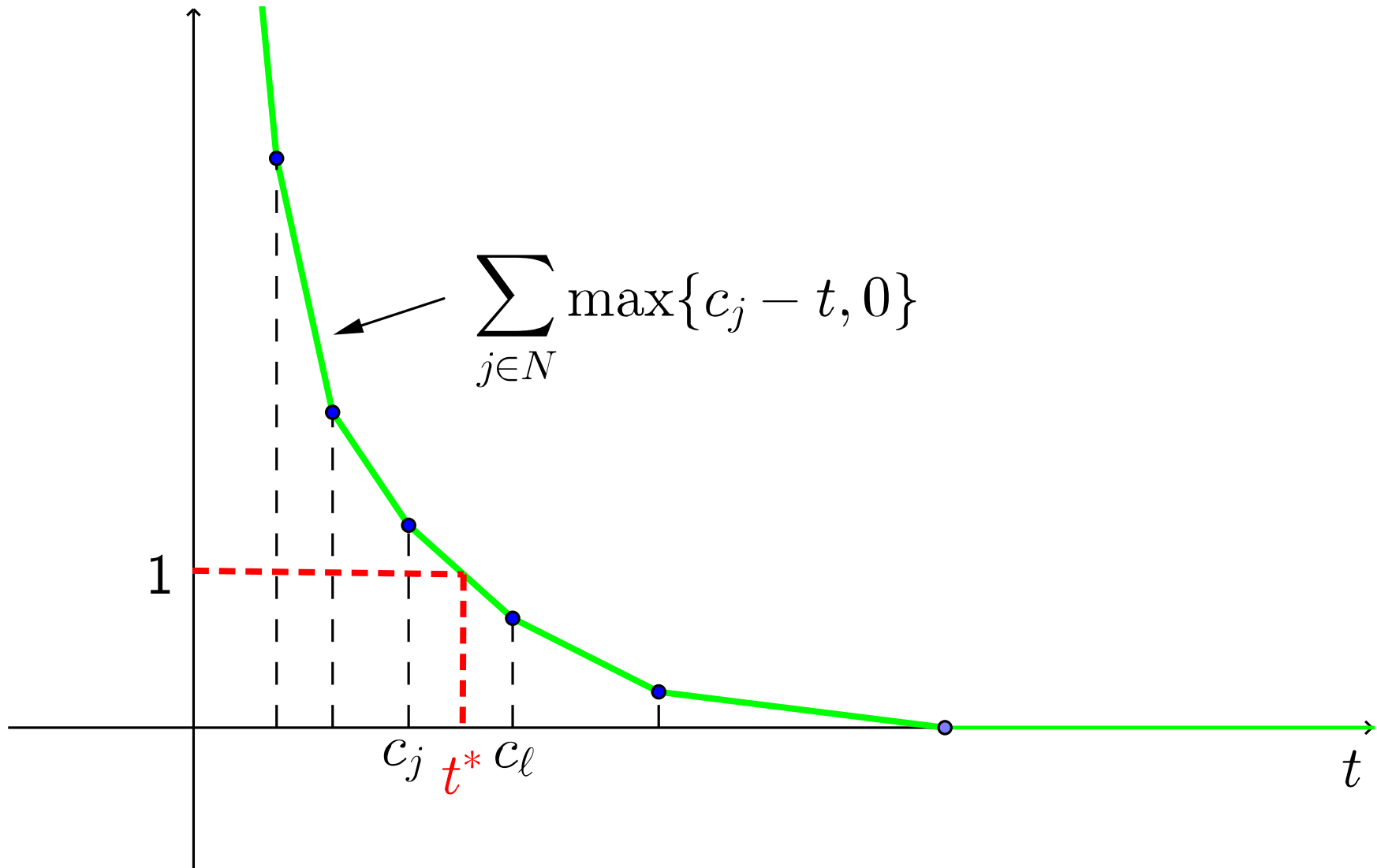


$n = 3$

КРАТКАЯ ИСТОРИЯ ВОПРОСА

1. Held M., Wolfe P., Crowder H. P. Validation of the subgradient optimization // Math. Programming. 6, pp. 62–88 (1974).
2. Brucker P. An $O(n)$ algorithm for quadratic knapsack problems // Operations Research Letters 3 (1984) 163–166.
3. Maculan, N., de Paula Jr., G.G.: A linear-time median-finding algorithm for projecting a vector on the simplex of \mathbb{R}^n // Operations Research Letters, 8 (4), pp. 219–222 (1989).
4. Малоземов В. Н., Певный А. Б. *Быстрый алгоритм проектирования точки на симплекс* // Вестник СПбГУ. Сер. 1. 1992. Вып. 1 (No 1). С. 112–113.

5. Michelot C. *A finite algorithm for finding the projection of a point onto the canonical simplex of \mathbb{R}^n* // JOTA. 1986. Vol. 50. No 1. P. 195—200.
6. Causa A., Raciti F. *A purely geometric approach to the problem of computing the projection of a point on a simplex* // JOTA. 2013. Vol. 156. No 2. P. 524—528.
7. Patriksson M. A survey on the continuous nonlinear resource allocation problem. Eur. J. Oper. Res, Vol. 185, No. 1 (Feb., 2008), pp. 1–46.
8. Малоземов В. Н., Тамасян Г. Ш. Ещё один быстрый алгоритм проектирование точки на стандартный симплекс // Семинар «DNA & CAGD». Избранные доклады. 5 сентября 2013 г.
(<http://dha.spb.ru/rep13.shtml#0905>)
9. Тамасян Г. Ш. Сравнительное изучение двух быстрых алгоритмов проектирования точки на стандартный симплекс // Семинар «CNSA & NDO». Избранные доклады. 15 мая 2014 г.
(<http://www.apmath.spbu.ru/cnsa/rep14.shtml#0515>)



АЛГОРИТМ МАЛОЗЁМОВА – ПЕВНОГО

1) Меняем знаки у компонент c_j точки c и числа $\{-c_j\}$ упорядочиваем по неубыванию. Получаем последовательность $a_1 \leq \dots \leq a_n$.

2) Проводим последовательные вычисления по рекуррентной формуле

$$\begin{aligned}\varphi_1 &= 0, \\ \varphi_{k+1} &= \varphi_k + k(a_{k+1} - a_k), \quad k = 1, \dots, n-1,\end{aligned}\tag{2}$$

пока не встретим индекс k_0 , на котором

$$\varphi_{k_0} < 1 \leq \varphi_{k_0+1}.$$

Если $\varphi_n < 1$, то полагаем $k_0 = n$.

3) Вычисляем λ^* по формуле

$$\lambda^* = a_{k_0} + \frac{1}{k_0}(1 - \varphi_{k_0}).\tag{3}$$

Компоненты x_i^* проекции точки c на стандартный симплекс имеют вид

$$x_i^* = (\lambda^* + c_i)_+, \quad i \in 1 : n,\tag{4}$$

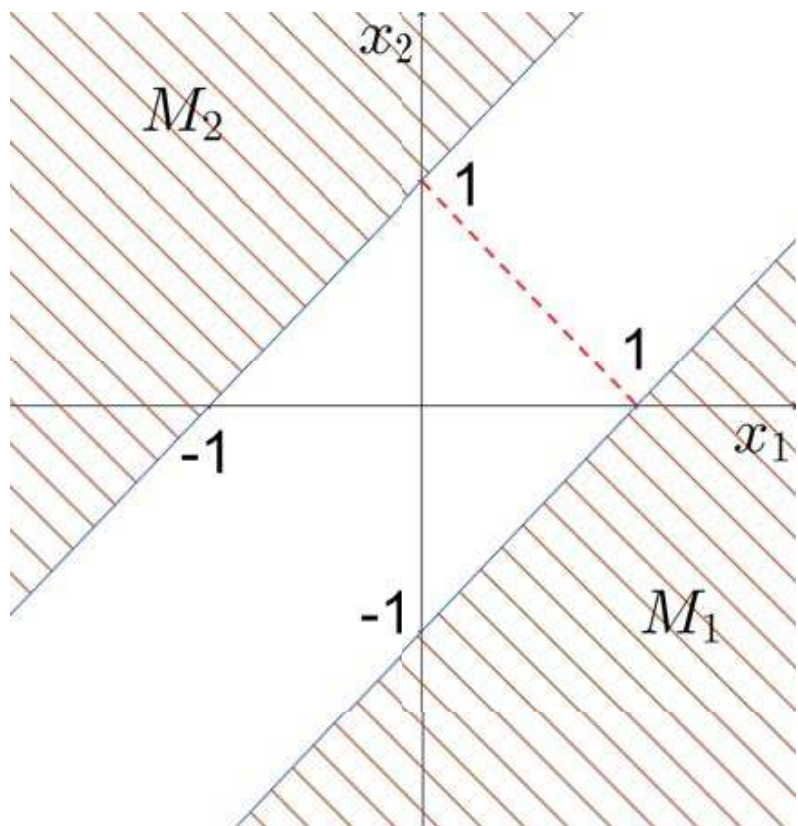
где $(u)_+ = \max\{0, u\}$.

КОММЕНТАРИИ К АЛГОРИТМУ

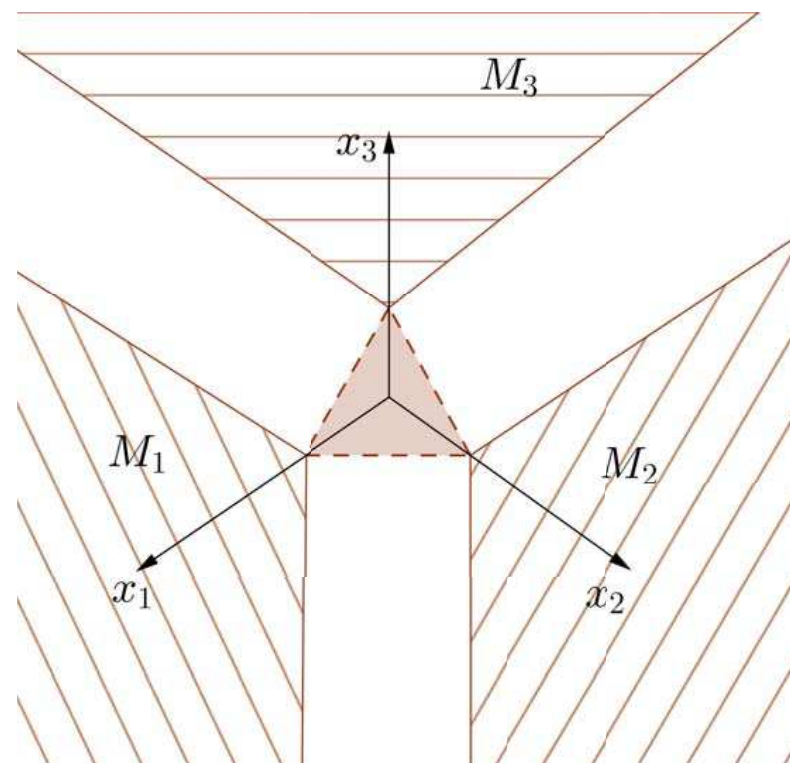
- Алгоритм сходится не более, чем за n шагов.
- В основе алгоритма лежит рекуррентное соотношение (2) для скалярных величин.
- В общем случае алгоритм требует сортировку элементов массива длиной n и не более $4n - 2$ арифметических операций.

ЛУЧШИЙ И ХУДШИЙ СЛУЧАИ АЛГОРИТМА МАЛОЗЁМОВА – ПЕВНОГО

минимальная трудоёмкость — $\varphi_2 \geq 1$



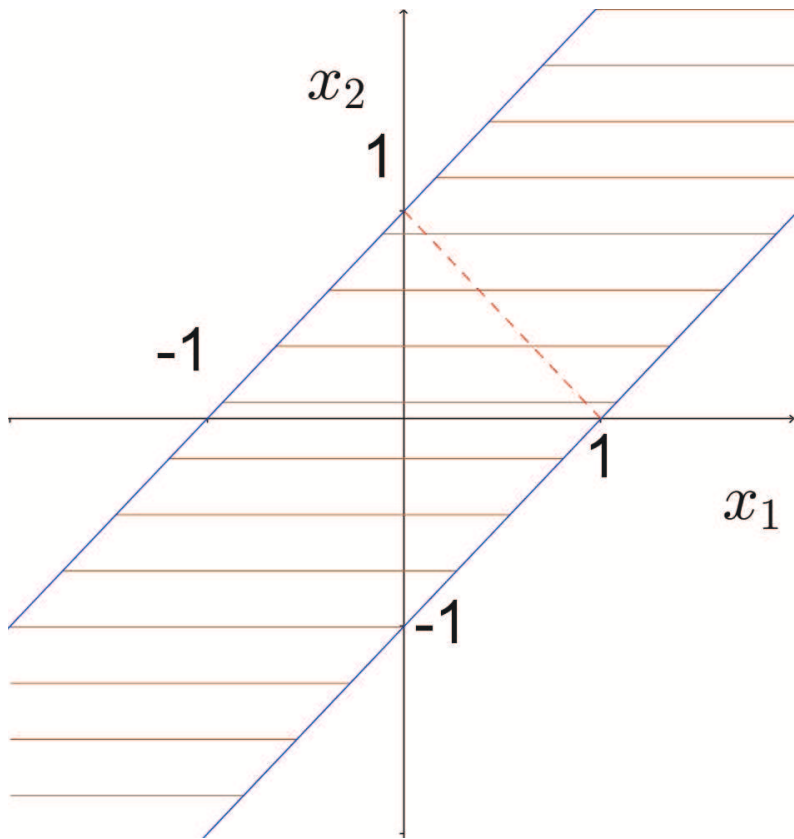
$n = 2$



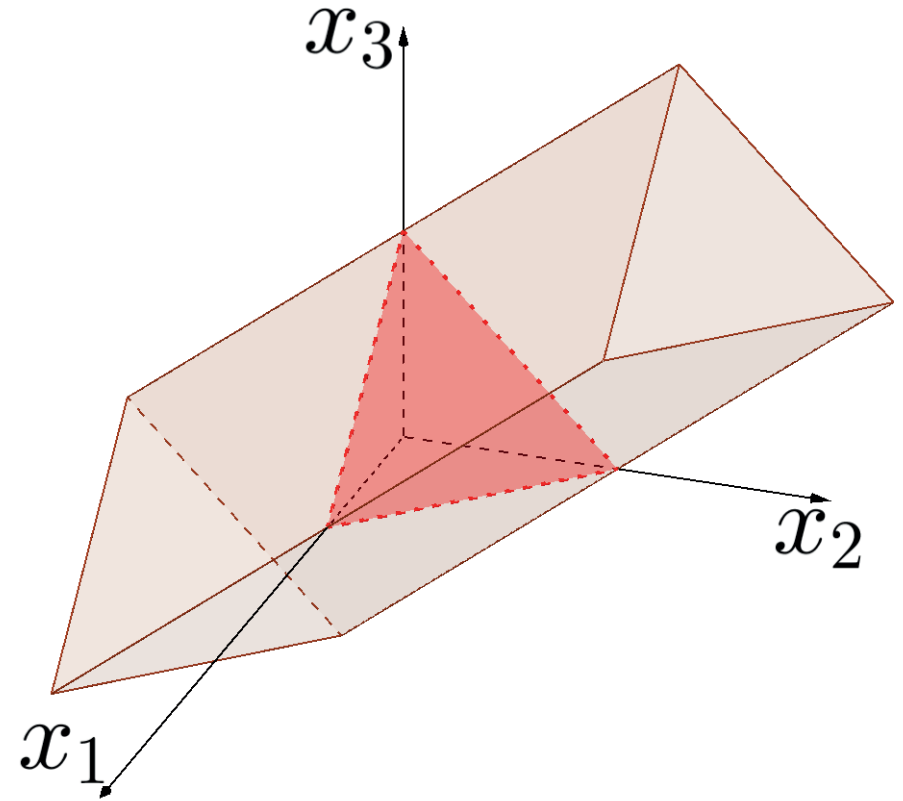
$n = 3$

ЛУЧШИЙ И ХУДШИЙ СЛУЧАИ АЛГОРИТМА МАЛОЗЁМОВА – ПЕВНОГО

максимальная трудоёмкость — $\varphi_{n-1} < 1$



$n = 2$



$n = 3$

МОДИФИКАЦИИ АЛГОРИТМА МАЛОЗЁМОВА – ПЕВНОГО

1. Очередь с приоритетом (наивная реализация, MalozemovNaivePQ).
 - предварительная подготовка не требуется;
 - нахождение очередного значения $O(n)$;
 - в худшем случае $O(n^2)$.
2. Очередь с приоритетом на основе сортирующего дерева (двоичной кучи, MalozemovPQ).
 - предварительное построение дерева $O(n)$;
 - нахождение очередного значения $O(\log(n))$;
 - в худшем случае $O(n \log(n))$.

МЕДИАНА МНОЖЕСТВА

Определение. Пусть $S = \{a_1, \dots, a_\ell\}$, где

$$a_1 \leq \dots \leq a_\ell.$$

Под медианой M множества S понимается величина $a_m \in S$, которая стоит на позиции $m = \left[\frac{\ell+1}{2} \right]$, т. е.

$$M = a_m,$$

где $[z]$ — целая часть числа z .

АЛГОРИТМ MACULAN – de PAULA (MACULAN).

Инициализация. $S := \{c_1, \dots, c_n\}, \quad J := N,$
 $v := 0, \quad p := 0, \quad q := 0.$

Определим функцию

$$\Phi(t, x, L, v, p, q) := \sum_{j \in L} (x_j - t) + v + p(q - t), \quad (5)$$

где $x = (x_1, \dots, x_n), \quad L \subseteq N.$

Шаг 1. Рассмотрим множество $S = \{c_j \mid j \in J\}$. Найдём медиану M множества S и построим следующие индексные множества

$$L_ = := \{j \in J \mid c_j = M\}, \quad L_ > := \{j \in J \mid c_j > M\}, \quad L_ < := \{j \in J \mid c_j < M\}.$$

Пусть $m \in L_ =$, т. е. $c_m = M$. Вычислим

$$z = \Phi(M, c, L_ >, v, p, q). \quad (6)$$

Если $z \geq 1$, то переходим к **Шагу 2**; иначе к **Шагу 3**.

Шаг 2. Положим

$$J := L_{>} \cup \{m\}, \quad S := \{c_j \mid j \in J\},$$

где m — индекс медианы S найденный на **Шаге 1**.

Если $|J| \geq 3$, то переходим к **Шагу 1**;
иначе при $v = z$ и $q = M$ к **Шагу 4**.

Шаг 3. Положим

$$J := L_{<} \cup \{m\}, \quad v := z, \quad q := M, \\ p := p + |L_{=}| - 1 + |L_{>}|.$$

Если $|J| \geq 2$, то переходим к **Шагу 1**; иначе к **Шагу 4**.

Шаг 4. Вычисляем t^* по формуле

$$t^* := q - \frac{1 - v}{1 + p}. \quad (7)$$

Компоненты x_i^* проекции точки c на стандартный симплекс имеют вид

$$x_i^* = (t^* + c_i)_+, \quad i \in 1 : n. \quad (8)$$

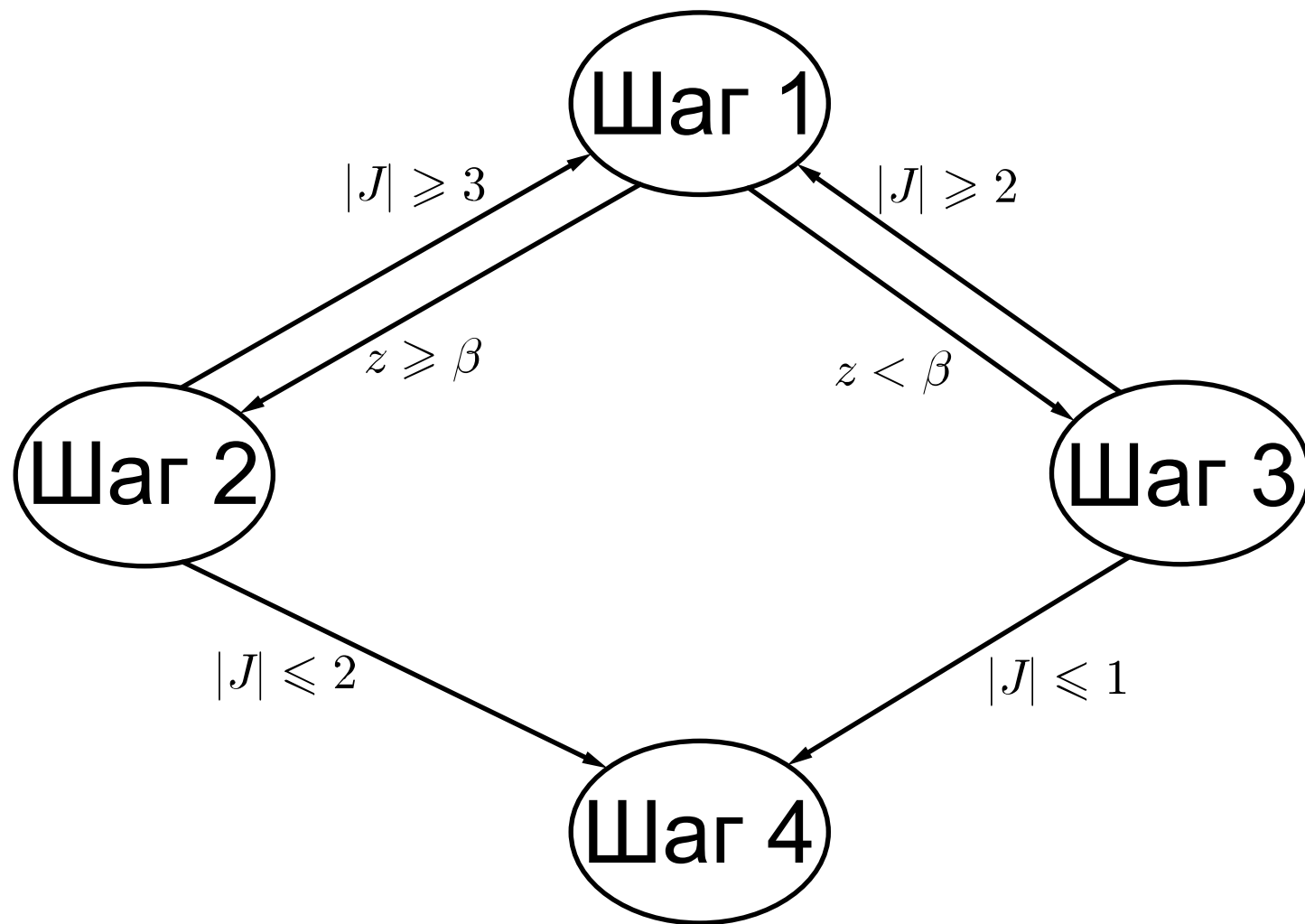


Схема алгоритма Maculan.

КОММЕНТАРИИ К АЛГОРИТМУ

Замечание 1. Величина z вычисляемая на **Шаге 1** (см. (6)) является значением функции Φ в точке t равной медиане M .

Замечание 2. Несложно понять, что мощность индексного множества J от итерации к итерации убывает. Это гарантирует конечность процесса.

О лучших и худших случаях для алгоритма Maculan – de Paula.

- Лучший случай: при $c_i = c_j, \forall i \neq j$, алгоритм завершается за одну итерацию.
- Худший случай: при $c_i \neq c_j, \forall i \neq j$, количество итераций определено размерностью задачи с точностью до константы.

$$\log_2(\hat{n}) - 0.5849626 < \hat{k} < \log_2(\hat{n}) + 2 \quad (9)$$

$$k_{\max}(\hat{n}) - k_{\min}(\hat{n}) \leq 2 \quad (10)$$

Минимальная и максимальная размерности для данного числа итераций

Пусть $d(k)$ — минимальная, а $D(k)$ — максимальная размерность задачи, которую алгоритм Maculan решает за k итераций.

$$k_{\min}(\hat{n}) = \underset{k}{\operatorname{argmin}}\{D(k) \mid D(k) \geq \hat{n}\}, \quad (11)$$

$$k_{\max}(\hat{n}) = \underset{k}{\operatorname{argmax}}\{d(k) \mid d(k) \leq \hat{n}\}. \quad (12)$$

k	1	2	3	4	5	6	7	8	9	10	11	12	13
$d(k)$	1	3	4	6	10	18	34	66	130	258	514	1026	2050
$D(k)$	3	6	12	24	48	96	192	384	768	1536	3072	6144	12288

Например, задача размерности 100 потребует

$k_{\min}(100) = 7$ или $k_{\max}(100) = 8$ итераций алгоритма Maculan.

Численные эксперименты

10000 точек в n -мерном евклидовом пространстве при n равном 10 , 10^2 , 10^3 , 10^4 , 10^5 и 10^6 .

Координаты точек $c = (c_1, \dots, c_n)$ формировались пятью способами:

(А) генерировались по непрерывному равномерному распределению на отрезке $[-10000, 10000]$;

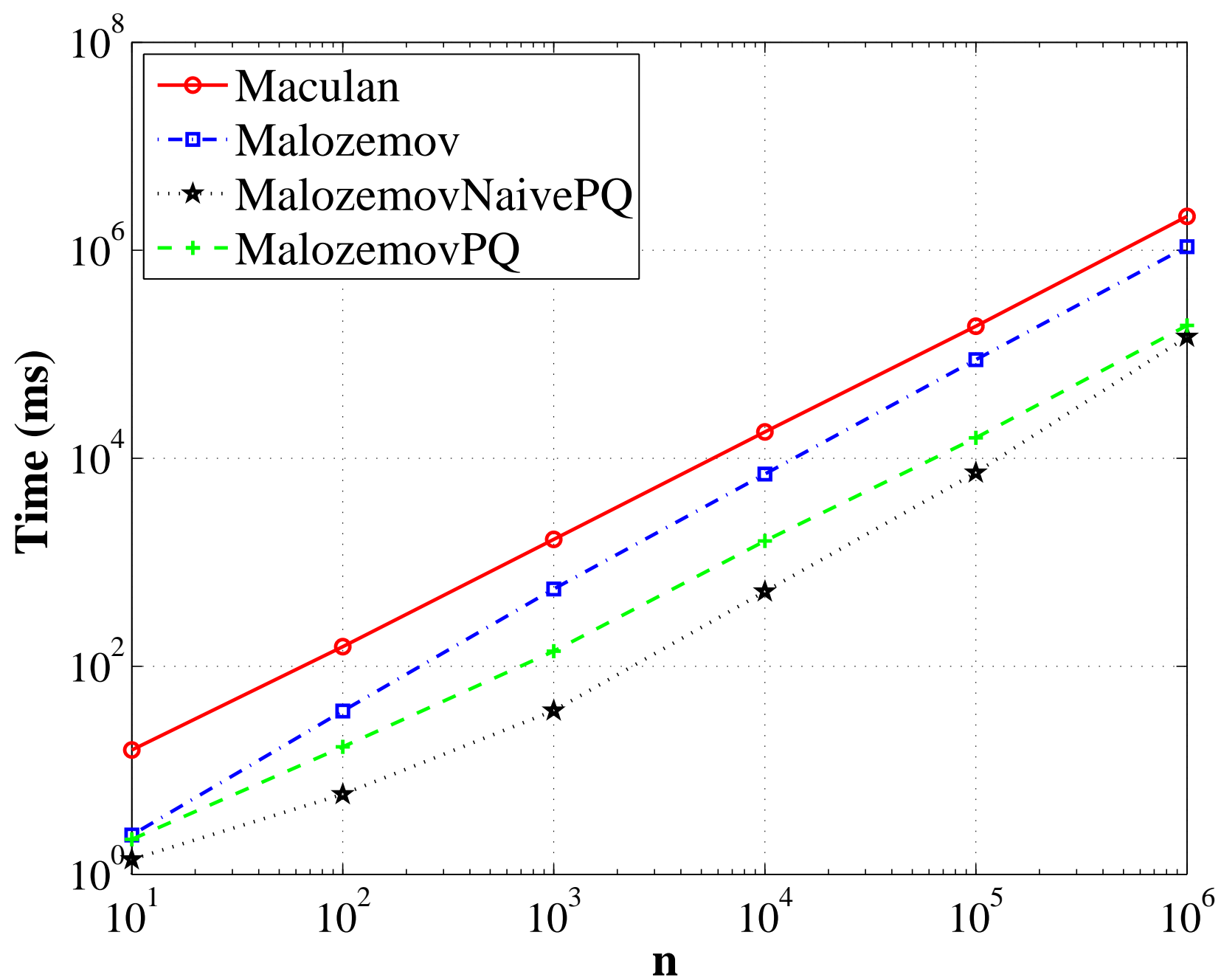
(В) $c_j = \xi + v_j$, $j \in 1 : n$, где (v_1, \dots, v_n) — любой элемент из Λ , ξ — произвольное фиксированное число из отрезка $[-10000, 10000]$;

(С) фиксировалось произвольное целое число k от 1 до n и вещественное c_k , тогда c_j при $j \neq k$ выбирались удовлетворяющие неравенству

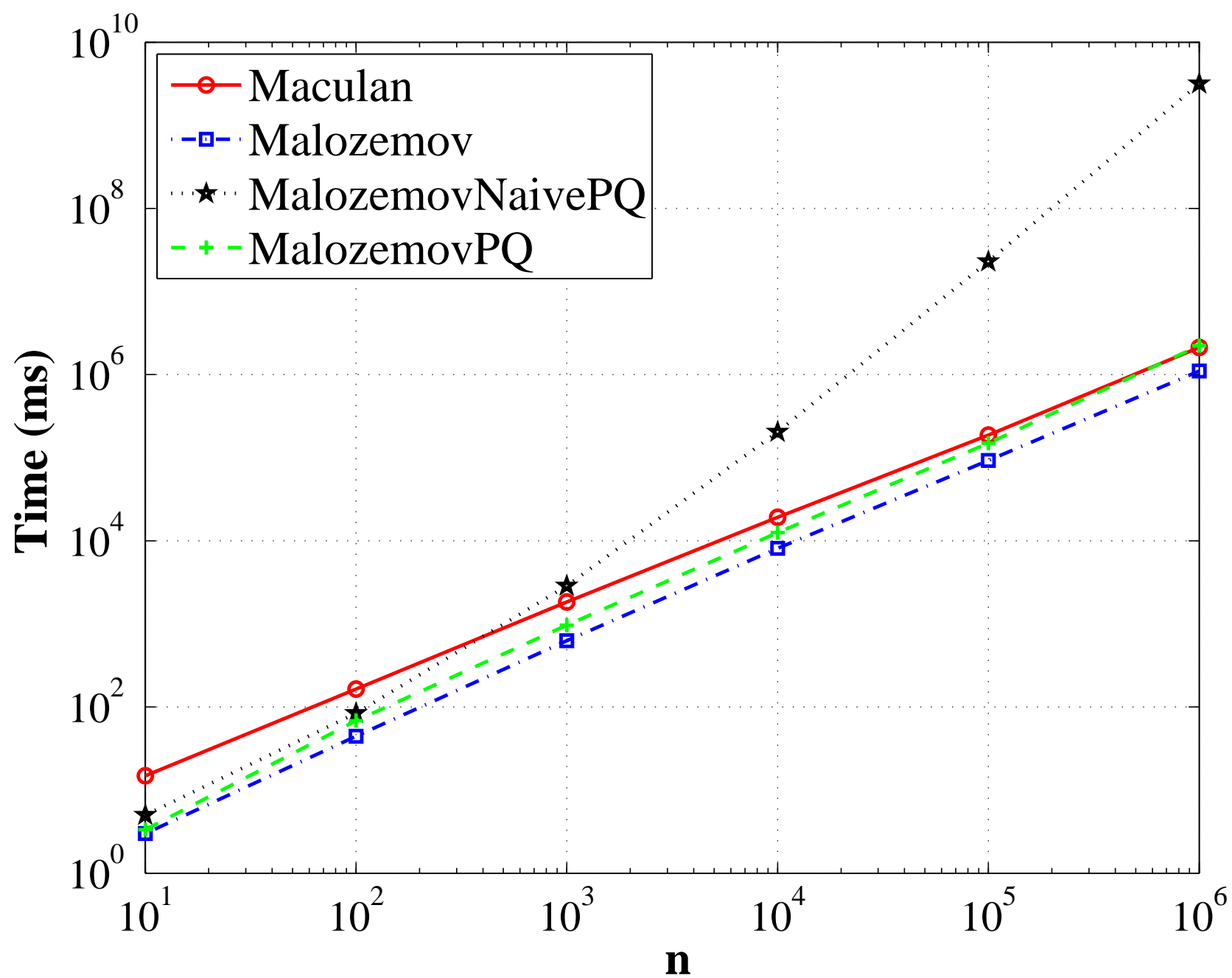
$$c_j \leq c_k - 1;$$

(D) $c_j \neq c_\ell$ для всех $1 \leq j \neq \ell \leq n$, т. е. все компоненты точки c различны;

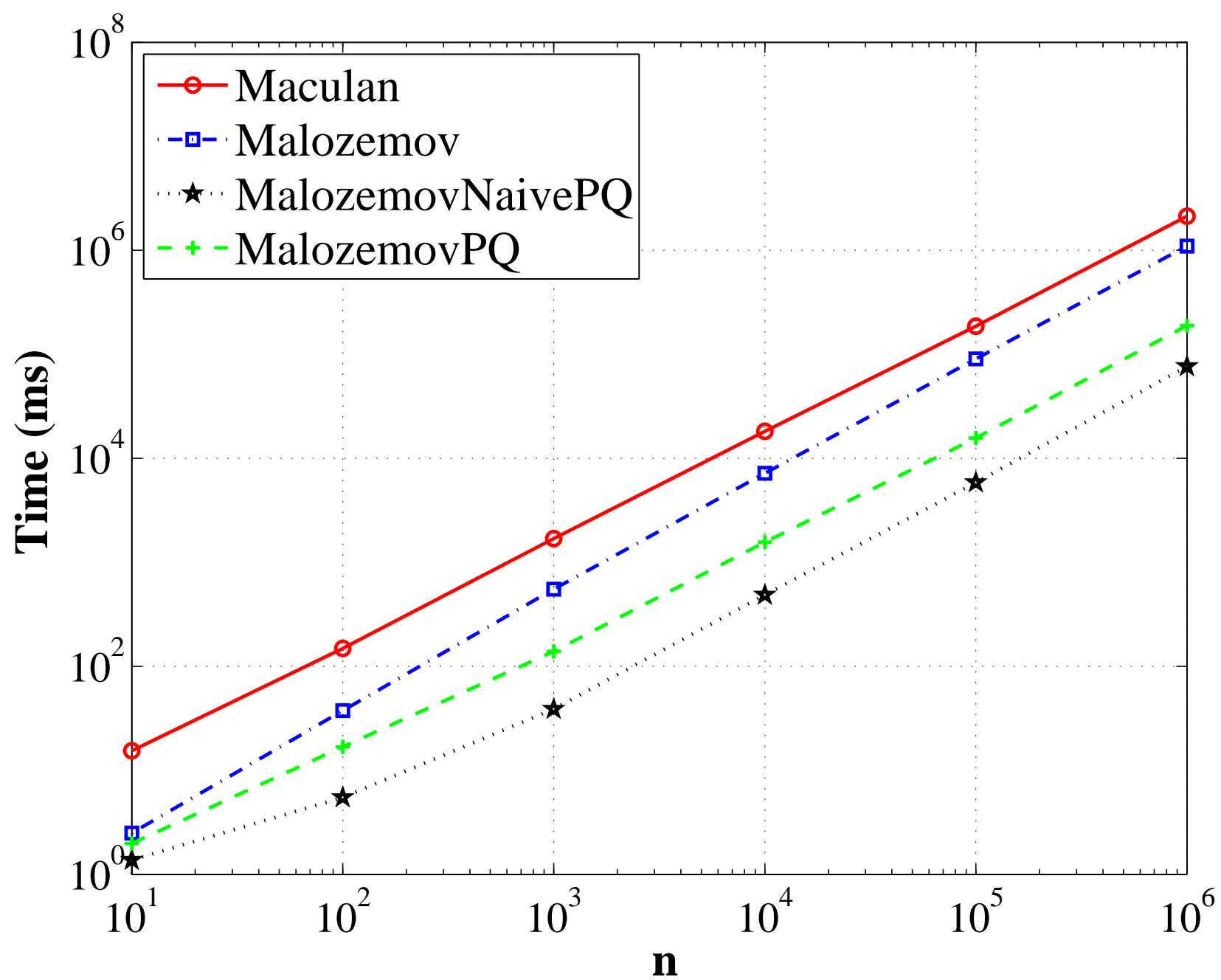
(Е) $c_j = \xi$ для всех $1 \leq j \leq n$, где ξ — произвольное фиксированное число из отрезка $[-10000, 10000]$.



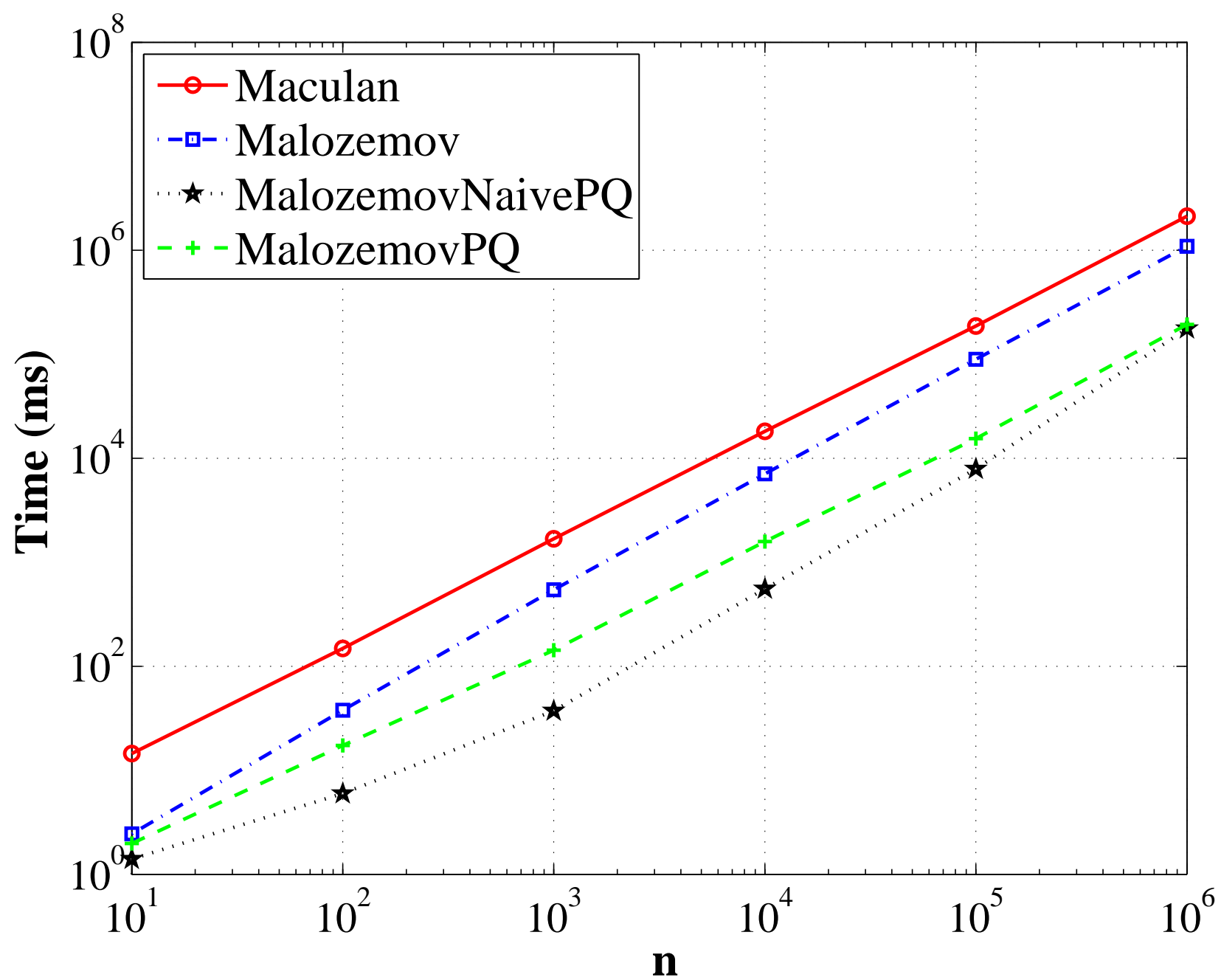
Случай (А) — равномерное распределение



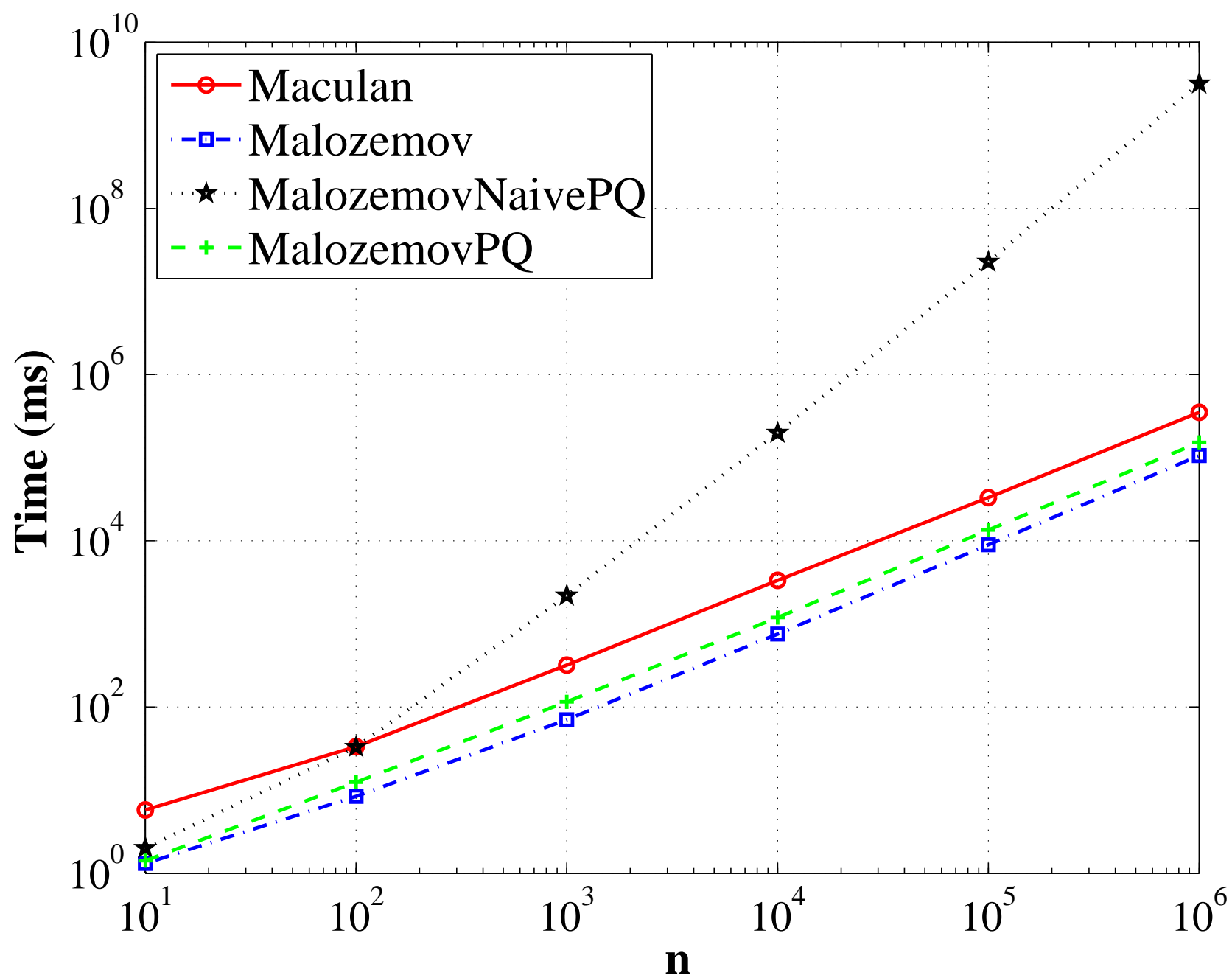
Случай (B) — худший случай для алгоритма Малоземова–Певного



Случай (С) — лучший случай для алгоритма Малоземова–Певного



Случай (D) — худший случай для алгоритма Maculan–de Paula



Случай (E) — лучший случай для алгоритма Maculan–de Paula