

# Правила остановки методов градиентного типа для невыпуклых задач при аддитивных ошибках градиента

Ф. С. Стонякин<sup>1</sup>

9 декабря 2022 г.

---

<sup>1</sup>На базе совместных работ с Курузовым И. А., Поляком Б. Т. и Алкусой М. С.

# Введение. Постановка задачи. Невыпуклые задачи и условие градиентного доминирования

- ▶ **Наша цель** — исследование методов градиентного типа в случае приближённого значения градиента  $\tilde{\nabla}f(x)$  минимизируемой функции  $f$  в произвольной запрашиваемой точке  $x \in \mathbb{R}^n$  при некотором фиксированном  $\Delta > 0$ :

$$\nabla f(x) = \tilde{\nabla}f(x) + v(x), \quad \text{причём} \quad \|v(x)\| \leq \Delta^1. \quad (1)$$

- ▶ Упор делается на **невыпуклые задачи**. В частности, на задачи минимизации функции удовлетворяющих  $(PL)$ -условию.

---

<sup>1</sup>Здесь и далее норма евклидова.

1. Рассмотрим простой пример не сильно выпуклой функции

$$f(x) = \langle Ax, x \rangle, \quad (2)$$

где  $A = \text{diag}(L, \mu, 0)$  — диагональная матрица 3-го порядка с ровно двумя положительными элементами  $L > \mu > 0$ .

При  $v(x) = (0, 0, \Delta)$ ,  $x_0 = (0, 0, 0)$  и  $h_k > 0$  последовательность  $x_{k+1} = x_k - h_k \tilde{\nabla} f(x_k)$  стремится к бесконечности при неограниченном увеличении  $k$ .

2. Далее, можно рассмотреть функцию Розенброка двух переменных  $x = (x^{(1)}, x^{(2)})$

$$f(x) = 100 \left( x^{(2)} - \left( x^{(1)} \right)^2 \right)^2 + \left( 1 - x^{(1)} \right)^2. \quad (3)$$

При  $x_0 = (1, 1) = x_*$  и погрешности  $v(x_k)$ , такой, что  $x_k^{(2)} = \left( x_k^{(1)} \right)^2$  траектория градиентного метода может уходить довольно далеко от точного решения  $x_*$ .

## Цель работы:

1. Исследовать оценку расстояния  $\|x_N - x_0\|$  для выдаваемых градиентным методом точек  $x_N$ .
2. Предложить правило ранней остановки градиентного метода, которое могло бы гарантировать некоторый компромисс между стремлением достичь приемлемого качества точки выхода по функции и целью обеспечить достаточно умеренное удаление такой точки от выбранного начального положения.

# Правила ранней остановки: история вопроса

- ▶ По-видимому, впервые идеология раннего прерывания итераций предложена в статье <sup>2</sup>, посвящённой методике для приближённого решения возникающих при регуляризации некорректных или плохо обусловленных задач.
- ▶ Известные подходы, связанные с ранней остановкой методов первого порядка в случае использования на итерациях неточной информации о градиенте (см. <sup>3</sup>, гл. 6, параграф 1, а также к примеру недавний препринт <sup>4</sup>)

---

<sup>2</sup>Емелин И. В., Красносельский М. А. Правило останова в итерационных процедурах решения некорректных задач. // Автомат. и телемех., 1978. Вып. 12. С. 59–63.

<sup>3</sup>Поляк Б. Т. Введение в оптимизацию. М.: Наука, 1983. 384 с.

<sup>4</sup>Vasin A., Gasnikov, A. and Spokoiny, V. Stopping rules for accelerated gradient methods with additive noise in gradient. // arxiv preprint (2021), Url: <https://arxiv.org/pdf/2102.02921.pdf>.

## Правила ранней остановки: история вопроса

- ▶ Однако данные (известные нам) результаты для выпуклых (не сильно выпуклых) задач отличаются по сравнению с полученными нами, что обычно гарантируется достижение худшего уровня по функции (даже без погрешностей получается лишь сублинейная скорость сходимости) или оценок вида  $\|x_N - x_*\| \leq \|x_0 - x_*\|$  без исследования  $\|x_N - x_0\|$ , где  $\{x_k\}_{k \in \mathbb{N}}$  — образуемая методом последовательность,  $x_*$  — ближайшее к точке старта метода  $x_0$  точное решение задачи минимизации  $f$ .

# План

- ▶ Предположения о задаче
- ▶ Известные результаты: градиентный метод с аддитивными помехами
- ▶ Основные теоретические результаты работы
- ▶ Результаты вычислительных экспериментов

Stopping Rules for Gradient Methods for Non-Convex Problems with Additive Noise in Gradient. Fedor S. Stonyakin, Ilya A. Kuruzov, Boris T. Polyak. ArXiv preprint: arXiv:2205.07544 [math.OC], 2022, <https://arxiv.org/abs/2205.07544>.

Gradient-Type Methods for Optimization Problems with Polyak-Lojasiewicz Condition: Early Stopping and Adaptivity to Inexactness Parameter. Ilya A. Kuruzov, Fedor S. Stonyakin, and Mohammad S. Alkousa. // Advances in Optimization and Applications, 2022.

## Предположения

Пусть норма  $\|\cdot\|$  евклидова. Будем рассматривать задачи минимизации  $f$ , удовлетворяющие условию Липшица градиента

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \forall x, y \in \mathbb{R}^n \quad (4)$$

для евклидовой нормы, а также  $(PL)$ -условию<sup>5, 6, 7</sup>

$$f(x) - f^* \leq \frac{1}{2\mu} \|\nabla f(x)\|^2 \quad \forall x \in \mathbb{R}^n, \quad (5)$$

где  $x_*$  — одно из точных решений задачи минимизации  $f$ ,  $f^* = f(x_*)$ , а  $\mu > 0$  — некоторая константа.

---

<sup>5</sup>Поляк Б. Т. Градиентные методы минимизации функционалов, 1963.

<sup>6</sup>Karimi, H., Nutini, J., Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the Polyak–Lojasiewicz condition, 2016.

<sup>7</sup>Belkin, M. Fit Without Fear: Remarkable Mathematical Phenomena of Deep Learning Through the Prism of Interpolation, 2021.



# Условие PL

1. Любая  $\mu$ -сильно выпуклая функция удовлетворяет условию PL с константой  $\mu$ ;
2. Можно рассмотреть систему нелинейных уравнений<sup>8</sup>  
 $g(x) = 0$  (записанную в векторном виде) и задачу нахождения какого-нибудь решения этой системы, где  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m \leq n$ . Тогда при условии  
$$\lambda_{\min} \left( \frac{\partial g(x)}{\partial x} \cdot \left[ \frac{\partial g(x)}{\partial x} \right]^T \right) \geq \mu, \forall x$$
 функция  $f(x) = \|g(x)\|^2$  удовлетворяет условию (5) для произвольного  $x_*$  такого, что  $f(x_*) = 0$ .<sup>9</sup>

---

<sup>8</sup>Гасников А. В. Современные численные методы оптимизации. Метод универсального градиентного спуска. М.: МЦНМО, 2021. 272 с.

<sup>9</sup>Nesterov Y., Polyak B. Cubic regularization of Newton method and its global performance. // Math. Program. Ser. A. 2006. Vol. 108. P. 177–205.

# Известные результаты: градиентный метод с аддитивными помехами

## Теорема 1 (1/3)

Пусть для численного решения задачи минимизации

$$f(x) \rightarrow \min_{x \in \mathbb{R}^n}$$

в случае, если  $f$  удовлетворяет условию  $L$ -гладкости (4), используется градиентный спуск вида

$$x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k). \quad (6)$$

Тогда верно неравенство

$$\min_{k=0, \dots, N-1} \|\nabla f(x_k)\|_2 \leq \sqrt{\frac{2L(f(x_0) - f(x_*))}{N}}. \quad (7)$$

## Теорема 1 (2/2)

Пусть  $f$   $L$ -гладкая и верно  $(PL)$ -условие (5). Тогда имеет место сходимость градиентного метода со скоростью геометрической прогрессии

$$\begin{aligned} f(x_N) - f(x_*) &\leq \left(1 - \frac{\mu}{L}\right)^N (f(x_0) - f(x_*)) \leq \\ &\leq \exp\left(-\frac{\mu}{L}N\right) (f(x_0) - f(x_*)), \end{aligned} \quad (8)$$

причём

$$\|x_* - x_0\| \leq \frac{\sqrt{2L(f(x_0) - f^*)}}{\mu}, \quad (9)$$

где  $x_*$  — ближайшее к  $x_0$  решение задачи минимизации  $f$ .

Нижняя оценка: ON THE LOWER BOUND OF MINIMIZING POLYAK-ŁOJASIEWICZ FUNCTIONS. // Under review as a conference paper at ICLR 2023.

## Замечание 1

Существенно, что оценки при использовании градиентного метода (6) с неточным градиентом вида (1)

$$\begin{aligned} \min_{k=0, N-1} \|\nabla f(x_k)\| &\leq \sqrt{\Delta^2 + \frac{2L(f(x_0) - f(x_*))}{N}} \leq \\ &\leq \Delta + \sqrt{\frac{2L(f(x_0) - f(x_*))}{N}} \end{aligned} \quad (10)$$

и

$$f(x_{k+1}) - f^* \leq \left(1 - \frac{\mu}{L}\right)^{k+1} (f(x_0) - f^*) + \frac{\Delta^2}{2\mu}, \quad (11)$$

вообще говоря, не улучшаемы. К примеру, известны нижние оценки уровня точности по функции  $O\left(\frac{\Delta^2}{2\mu}\right)$  для градиентного метода с аддитивно неточным градиентом (см., например, раздел 2.11.1 пособия<sup>10</sup>, а также имеющиеся там ссылки) даже на классе сильно выпуклых функций.

<sup>10</sup>Воронцова Е. А., Хильдебранд Р. Ф., Гасников А В., Стонякин Ф. С. Выпуклая оптимизация 2021.

## Пример

Действительно, рассмотрим такой пример:

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{1}{2} \sum_{i=1}^n \lambda_i (x^i)^2, \quad (12)$$

где  $0 \leq \mu = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = L$ ,  $L \geq 2\mu$ . Решением задачи (12) будет  $x^* = 0$ . Предположим, что неточность имеется только в вычислении первой компоненты градиента. То есть вместо  $\partial f(x)/\partial x^1 = \mu x^1$  нам доступно только  $\tilde{\partial} f(x)/\partial x^1 = \mu x^1 - \Delta$ .

Для простейшей динамики градиентного спуска (в этом разделе мы обозначаем номер компоненты вектора верхним индексом, потому что далее будет использоваться одновременно номер итерации и номер компоненты (13)):

$$x_{k+1} = x_k - \frac{1}{L} \tilde{\nabla} f(x_k),$$

можно получить, что при  $x_0^1 \geq 0$  и достаточно больших  $k \in \mathbb{N}$  ( $k \gg L/\mu$ )

$$x_k^1 \geq \frac{\Delta}{L} \cdot \frac{1 - (1 - \mu/L)^k}{1 - (1 - \mu/L)} \simeq \frac{\Delta}{\mu}. \quad (13)$$

Следовательно, при  $x_0^1 \geq 0$  и достаточно больших  $k \in \mathbb{N}$

$$f(x_k) - f(x_*) \gtrsim \frac{\Delta^2}{2\mu}.$$

## Основные результаты

$$x_{k+1} = x_k - \frac{1}{L} \tilde{\nabla} f(x_k) \quad (14)$$

$$f(x_{k+1}) - f(x_k) \leq \frac{\Delta^2}{L} - \frac{1}{4L} \|\tilde{\nabla} f(x_k)\|^2. \quad (15)$$

Неравенство (15) показывает, что в случае достаточно большого значения  $\|\tilde{\nabla} f(x_k)\|$  можно гарантировать, что  $f(x_{k+1}) < f(x_k)$ . Тем самым для всякого  $C > 2$  возникает альтернатива: или верно неравенство  $\|\tilde{\nabla} f(x_k)\| \leq C\Delta$ , и это само по себе гарантирует достижение приемлемого качества точки выхода  $x_k$  по функции в силу  $(PL)$ -условия, или же

$$f(x_{k+1}) - f(x_k) < -\frac{\Delta^2}{L} \left( \frac{C^2}{4} - 1 \right). \quad (16)$$

Тем самым за конечное число шагов градиентного метода (20) возможно получить  $x_k$  такое, что значение  $f(x_k)$  достаточно близко к минимальному  $f^*$ . Выберем для определённости  $C = \sqrt{6}$  (чтобы получить «удобный» коэффициент) и будем рассматривать 2 сценария:

1.  $\|\tilde{\nabla} f(x_k)\| > \Delta\sqrt{6}$ , что с учётом (15) влечёт неравенство

$$f(x_{k+1}) - f(x_k) < -\frac{\Delta^2}{2L}, \quad (17)$$

указыв. на убывание значения  $f$  при переходе от  $x_k$  к  $x_{k+1}$ .

- 2.

$$\|\tilde{\nabla} f(x_k)\| \leq \Delta\sqrt{6}, \quad (18)$$

откуда

$$f(x_k) - f^* \leq \frac{7\Delta^2}{\mu}. \quad (19)$$



## Теорема 2

Пусть на  $N$ -й итерации градиентного метода

$$x_{k+1} = x_k - \frac{1}{L} \tilde{\nabla} f(x_k) \quad (20)$$

впервые выполнен критерий остановки

$$\|\tilde{\nabla} f(x_k)\| \leq \Delta \sqrt{6}. \quad (21)$$

Тогда для точки выхода  $\hat{x} = x_N$  гарантированно будет верно неравенство

$$f(\hat{x}) - f^* \leq \frac{7\Delta^2}{\mu}. \quad (22)$$

При этом справедлива следующая оценка количества итераций до остановки

$$N < \frac{2L}{\Delta^2} (f(x_0) - f^*). \quad (23)$$

### Теорема 3

Пусть градиентный метод (20) работает либо

$$N_* = \left\lceil \frac{L}{\mu} \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right\rceil. \quad (24)$$

шагов, либо при некотором  $N \leq N_*$  на  $N$ -й итерации метода (20) впервые выполнен критерий остановки (21). Тогда для точки выхода  $\hat{x}$  ( $\hat{x} = x_N$  или  $\hat{x} = x_{N_*}$ ) метода (20) гарантированно будет верно неравенство

$$f(\hat{x}) - f^* \leq \frac{7\Delta^2}{\mu},$$

причём

$$\|\hat{x} - x_0\| \leq \frac{2\Delta}{\mu} \sqrt{1 + \frac{L}{\mu} \left\lceil \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right\rceil} + \frac{4\sqrt{L(f(x_0) - f^*)}}{\mu}. \quad (25)$$

### Замечание 3

В силу теоремы 1 и (25) величину  $\|\hat{x} - x_0\|$  можно считать сопоставимой с  $\|x_* - x_0\|$  при достаточно малом  $\Delta > 0$ , поскольку

$$\|x_* - x_0\| \leq \frac{\sqrt{2L(f(x_0) - f^*)}}{\mu}. \quad (26)$$

### Замечание 4

Ввиду (25) достаточно потребовать выполнения (4) и (5) ( $L$ -гладкость и условие PL) только в  $R$ -окрестности  $x_0$ , где

$$R = \frac{2\Delta}{\mu} \sqrt{1 + \frac{L}{\mu} \left[ \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right]} + \frac{4\sqrt{L(f(x_0) - f^*)}}{\mu}.$$

Characterization of Gradient Dominance and Regularity Conditions for Neural Networks. Yi Zhou, Yingbin Liang. ArXiv preprint: arXiv:2205.07544 [math.OC], 2017, <https://arxiv.org/abs/1710.06910>.

# Вариант градиентного метода с адаптивно подбираемым шагом

**Идея метода с адаптивным шагом:** подбирать величину  $L_k$  так, чтобы было выполнено условие

$$\tilde{f}(x_{k+1}) \leq \tilde{f}(x_k) + \langle \tilde{\nabla} f(x_k), x_{k+1} - x_k \rangle + L_k \|x_{k+1} - x_k\|^2 + \frac{\Delta^2}{2L_k} + 2\delta,$$

где  $\tilde{f}(x)$  есть неточное значение функции  $f$  в точке  $x$ .

---

## Алгоритм 1. Адаптивный вариант градиентного метода.

---

**Require:**  $L_{\min} \geq 0, L_0 \geq L_{\min}, \delta \geq 0, \Delta \geq 0$

1:  $k := 0$

2:  $L_k = \max\left(\frac{L_{k-1}}{2}, L_{\min}\right)$

3: Вычисляем новое значение  $x$ :

$$x_{k+1} = x_k - \frac{1}{2L_k} \tilde{\nabla} f(x_k) \quad (27)$$

4: Если выполнено

$$\tilde{f}(x_{k+1}) \leq \tilde{f}(x_k) + \langle \tilde{\nabla} f(x_k), x_{k+1} - x_k \rangle + L_k \|x_{k+1} - x_k\|^2 + \frac{\Delta^2}{2L_k} + 2\delta \quad (28)$$

то увеличиваем  $k := k + 1$  и переходим к шагу 2. Иначе  $L_k := 2L_k$  и переходим к шагу 3.

5: **return**  $x_k$

---

## Теорема 4 (1/2)

Пусть в алгоритме 1  $L_{\min} \geq \frac{\mu}{4}$ , а также верны  $(PL)$ -условие (5) и

$$|f(x) - \tilde{f}(x)| \leq \delta, \quad (29)$$

причем  $\Delta^2 \geq 16L\delta$ . Тогда алгоритм 1 работает либо

$$N_* = \left\lceil \frac{8L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2} \right\rceil \quad (30)$$

шагов, либо при некотором  $N \leq N_*$  на  $N$ -й итерации алгоритма 1 выполнен критерий остановки

$$\|\tilde{\nabla} f(x_k)\| \leq 2\Delta. \quad (31)$$

## Теорема 4 (2/2)

Тогда для точки выхода  $\hat{x}$  ( $\hat{x} = x_N$  или  $\hat{x} = x_{N^*}$ ) алгоритма 1 гарантированно будет верно неравенство

$$f(\hat{x}) - f^* \leq \frac{5\Delta^2}{\mu},$$

причём

$$\|\hat{x} - x_0\| \leq 8 \frac{\Delta}{\mu} \sqrt{\frac{1}{2} \gamma^2 + 4 \gamma \frac{L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2}} + 16 \frac{\sqrt{\gamma L(f(x_0) - f^*)}}{\mu}, \quad (32)$$

где  $\gamma = \frac{L}{L_{\min}}$ . Также суммарное количество обращений к подпрограмме для вычисления неточных значений функции и шага (27) не более чем  $2N + \log \frac{2L}{L_0}$ .

## Замечание 5

Условие

$$\|\tilde{\nabla} f(x_k)\| \leq 2\Delta \quad (33)$$

выполнено для любых  $L_k \geq L$  и по построению получаем, что  $L_k \leq 2L$ . В оценках выше величина  $2L$  оценивает максимальное значение параметра  $L_k$ . Оценки выше останутся верными, если заменить  $L$  на  $\frac{1}{2} \max_{j \leq k} L_j$  и  $\gamma$  на  $\frac{\max_{j \leq k} L_j}{2 \min_{j \leq k} L_j}$ . Аналогично можно заменить параметр алгоритма  $L_{\min}$  на  $\min_{j \leq k} L_j$ .



## Замечание 6

Оценка на количество итераций

$$N_* = \left\lceil \frac{8L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2} \right\rceil \quad (34)$$

в теореме 4 говорит о конечности процесса, однако является сильно завышенной. На практике более интересной оценкой является

$$N_* = \left\lceil \frac{4\hat{L}}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2} \right\rceil,$$

где  $\hat{L} = \frac{\mu}{4} \frac{1}{1 - \left( \prod_{j=0}^{N_*-1} \left( 1 - \frac{\mu}{4L_j} \right) \right)^{\frac{1}{N_*}}}$  — величина, зависящая от

подобранных параметров  $L_j$  в процессе работы алгоритма 1.

### Замечание 7

Возможно ослабить требование  $L_{\min} \geq \frac{\mu}{4}$  до  $L_{\min} > 0$ . В таком случае будет верна оценка

$$\|x_N - x_0\| \leq N\Delta \sqrt{\frac{1}{2L_{\min}^2} + \frac{4}{\mu L_{\min}}} + 16\sqrt{\frac{L}{L_{\min}}} \frac{\sqrt{L(f(x_0) - f^*)}}{\mu}$$

на  $N$ -ой итерации градиентного метода, однако уже не верна оценка

$$\|\hat{x} - x_0\| \leq 8\frac{\Delta}{\mu} \sqrt{\frac{1}{2}\gamma^2 + 4\gamma\frac{L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2}} + 16\frac{\sqrt{\gamma L(f(x_0) - f^*)}}{\mu} \quad (36)$$

из теоремы 4. В таком случае возможно оценить достаточное для гарантированного выполнения критерия  $\|\tilde{\nabla} f(x_k)\| \leq 2\Delta$  количество итераций алгоритма 1

$$N < \frac{2L}{\Delta^2 - 16L\delta} (f(x_0) - f^*), \quad (37)$$

где  $\Delta^2 > 16L\delta$ .

---

## Алгоритм 2. Адаптивный градиентный метод для неизвестных $L$ и $\Delta$ .

---

**Require:**  $x_0, L_{\min} \geq \frac{\mu}{4} > 0, L_0 \geq L_{\min}, \Delta_0 > 0, \Delta_{\min} > 0$ .

- 1:  $k := 0, L_k := \max \left\{ \frac{L_{k-1}}{2}, L_{\min} \right\}$ .
- 2: Если условие останова выполнено для  $x_k$ , то переходим к шагу 6. Иначе вычисляем  $x_{k+1} = x_k - \frac{1}{2L_k} \tilde{\nabla} f(x_k)$
- 3: Если

$$f(x_{k+1}) \leq f(x_k) + \langle \tilde{\nabla} f(x_k), x_{k+1} - x_k \rangle + \Delta_k \|x_{k+1} - x_k\| + \frac{L_k}{2} \|x_{k+1} - x_k\|^2, \quad (35)$$

то переходим к шагу 4. Иначе  $L_k := 2L_k$  and  $\Delta_k = 2\Delta_k$  и переходим к шагу 2.

- 4: Переопределим  $\Delta_k$  так, чтобы было верно (35), а также  $\Delta_k \geq \Delta_{\min}$  и  $\Delta_k \geq \max_{1 \leq j < k} \Delta_j$ .
  - 5: Обновляем  $k := k + 1$  и переходим к шагу 1.
  - 6: **return**  $x_k$ .
-

# Метод

## Критерий Останова

Выполняется одна из следующих альтернатив:

1. Алгоритм 2 работает  $N_*$  шагов
2. Для некоторого  $k_* \leq N_*$  на  $k_*$ -ой итерации выполняется критерий останова

$$\|\tilde{\nabla} f(x_{k_*})\| \leq 2 \max_{j \leq k} \Delta_j, \quad (38)$$

$$N_* = \left\lceil \frac{8L_{\max}}{\mu} \log \left( \frac{\mu(f(x_0) - f^*)}{4\Delta_{\max}^2} \right) \right\rceil.$$

# Теоретические результаты для алгоритма 2

## Теорема 5

Пусть либо алгоритм 2 работает

$$N_* = \left\lceil \frac{8L_{\max}}{\mu} \log \left( \frac{\mu(f(x_0) - f^*)}{4\Delta_{\max}^2} \right) \right\rceil \quad (39)$$

шагов, либо для некоторого  $k_* \leq N_*$  на  $k_*$ -ой итерации, выполнен критерий останова

$$\|\tilde{\nabla} f(x_{k_*})\| \leq 2 \max_{j \leq k} \Delta_j.$$

Тогда для точки выхода  $\hat{x}$  ( $\hat{x} = x_{k_*}$  or  $\hat{x} = x_{N_*}$ ) алгоритма 2, верно следующее соотношение

$$f(\hat{x}) - f^* \leq \frac{5 \max_{j \leq k} \Delta_j^2}{\mu} \leq \frac{5\Delta_{\max}^2}{\mu}. \quad (40)$$

# Теоретические результаты

## Теорема 6

Более того,

$$\|\hat{x} - x_0\| \leq \frac{8\Delta_{\max}}{\mu} \sqrt{\frac{\gamma^2}{2} + \frac{2\gamma L_{\max}}{\mu}} \log \left( \frac{\mu(f(x_0) - f^*)}{4\Delta_{\max}^2} \right) + \frac{16\sqrt{\gamma L_{\max}(f(x_0) - f^*)}}{\mu},$$

где  $\gamma = \frac{4L_{\max}}{\mu}$ ,  $\Delta_{\max} = 2\Delta \max \left\{ \frac{L}{L_{\min}}, 1 \right\}$ ,  $L_{\max} = L \max \left\{ \frac{\Delta}{\Delta_{\min}}, 1 \right\}$ .

Также, общее количество вызовов подпрограммы для вычисления функции  $f$  будет не больше, чем

$$N_* \log_2 \left( \frac{4L}{L_{\min}} \max \left\{ \frac{L}{L_{\min}}, \frac{\Delta}{\Delta_{\min}} \right\} \right).$$

# Теоретические результаты

1. Величина  $L_{\max}$  оценивает максимальное значение параметра  $L_k$ . Оценки сверху остаются корректными, если заменить  $L_{\max}$  на  $\max_{j \leq k} L_j$ .
2. Пусть в каждой точке  $x$  мы имеем доступ к модели  $(\tilde{f}, \tilde{\nabla} f)$  функции  $f$ , такой что следующие условия выполняются:

$$\tilde{f}(x) \leq \tilde{f}(y) + \langle \tilde{\nabla} f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2 + \Delta \|x - y\| + \delta,$$

$$\tilde{f}(x) - f^* \leq \frac{1}{\mu} \left( \|\tilde{\nabla} f(x)\| + \Delta^2 \right) + \delta.$$

В таком случае, оценки на число итераций и точность решения, и расстояние от стартовой точки до точки выхода остаются верными, если изменить параметры  $\Delta_{\max}$  и  $L_{\max}$  соответствующим образом.

# Численные эксперименты

## 1. Задача минимизации квадратичной формы

Сравним количество итераций, требуемых для достижения условия (21) и оценку  $N_*$  из теоремы 3.

- ▶ Задача минимизации следующей квадратичной формы

$$\min_{x \in \mathbb{R}^n} \sum_{j=p+1}^n d_j x_j^2, \quad (41)$$

где  $p$  — количество нулевых собственных значений матрицы квадратичной формы,  $d_j$  — некоторые положительные константы. Т.е. имеем квадратичную форму с неотрицательно определенной вырожденной диагональной матрицей.

- ▶ Случайно генерируемая неточность:  $v(x) \sim \mathcal{U}(S_1^n(0))$
- ▶  $\mu = \min_{j=\overline{p+1,n}} (d_j), L = \max_{j=\overline{p+1,n}} (d_j)$
- ▶  $X^* = \{x | x_j = 0, j = \overline{p+1,n}\}$



# Численные эксперименты

## Параметры эксперимента:

- ▶  $L = 1$
- ▶  $\mu$  варьируется от 0 до 1
- ▶ Параметры  $d_j$  будут браться равномерно из отрезка  $[\mu, L]$
- ▶  $n = 100$
- ▶  $p = 10$

## Численные эксперименты

$\mu$	$\Delta$	$N$	$N_*$
0.01	$10^{-7}$	1528	3817
	$10^{-4}$	841	2436
	$10^{-1}$	155	1054
0.1	$10^{-7}$	169	406
	$10^{-4}$	104	267
	$10^{-1}$	40	129
0.9	$10^{-7}$	10	48
	$10^{-4}$	8	33
	$10^{-1}$	5	17
0.99	$10^{-7}$	6	44
	$10^{-4}$	5	30
	$10^{-1}$	3	16

Таблица 1: Сравнение номера итерации  $N$  для достижения условия (21) и оценки  $N_*$  из теоремы 3.

# Численные эксперименты

## Квадратичная функция

$\mu$	$\Delta$	Alg. constant	Adaptive $L$	Full Adaptive
0.01	$10^{-7}$	2.29	2.03	1.84
	$10^{-4}$	2.26	2.31	2.43
	$10^{-1}$	2.27	1.95	1.07
0.1	$10^{-7}$	2.17	1.97	0.87
	$10^{-4}$	2.18	1.68	0.80
	$10^{-1}$	2.14	2.26	0.83
0.9	$10^{-7}$	0.92	1.58	0.69
	$10^{-4}$	0.91	0.96	0.79
	$10^{-1}$	0.96	0.95	0.72
0.99	$10^{-7}$	0.96	0.95	0.71
	$10^{-4}$	0.95	0.92	0.68
	$10^{-1}$	1.05	0.93	0.69

Таблица 2: Сравнение достигнутой величины  $\frac{\|\tilde{\nabla} f(x_N)\|}{\Delta}$ .

# Численные эксперименты

## 2. Задача минимизации функции логистической регрессии

Рассмотрим задачу минимизации логистической регрессии:

$$f(x) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle w_i, x \rangle)), \quad (42)$$

где  $y = (y_1 \dots y_m)^\top \in [-1, 1]^m$  — вектор целевой переменной,  $W = [w_1 \dots w_m] \in \mathbb{R}^{n \times m}$  — матрица признаков, где вектор  $w_i \in \mathbb{R}^n$  из того же пространства, что и оптимизируемый вектор весов  $x$ .

Заметим, что в общем случае эта задача может не иметь конечного решения, поэтому мы создадим такой искусственный датасет, что существует конечный вектор  $x^*$ , минимизирующий заданную функцию.

# Численные эксперименты

## Параметры эксперимента:

- ▶ Размерность задачи:  $n = 200$
- ▶ Количество элементов в сумме:  $m = 700$
- ▶ Функция удовлетворяет  $(PL)$ -условию локально
- ▶ Случайно генерируемая неточность:  $v(x) \sim \mathcal{U}(S_1^n(0))$

# Численные эксперименты

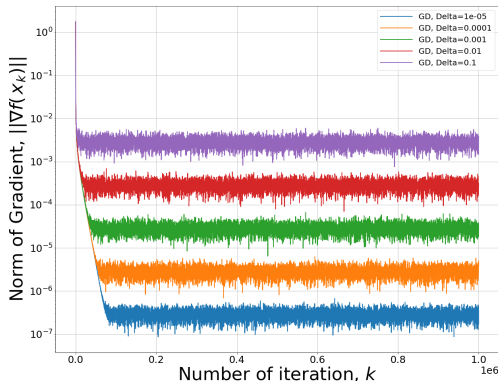
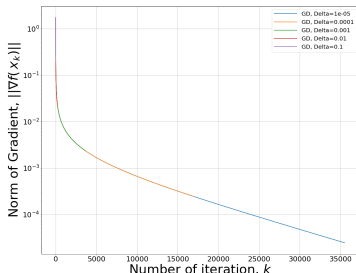


Рис. 1: Скорость сходимости градиентного метода по норме градиента для различных значений неточностей градиента  $\Delta$  на задаче минимизации логистической регрессии на первых  $N = 10^5$  итераций без использования правила останова.

# Численные эксперименты

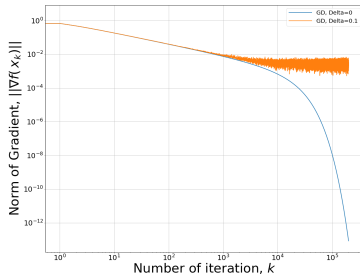


**Рис. 2:** Скорость сходимости градиентного метода по норме градиента для различных значений неточностей градиента  $\Delta$  на задаче минимизации логистической регрессии при использовании правила остановки (21).

Как можно видеть, метод останавливается, когда метод достигает точности  $\|\tilde{\nabla} f(x_k)\| \sim \Delta$ . Также можно заметить, что на данном примере траектории методов практически совпадают до достижения соответствующей точности.

# Численные эксперименты

(a)



(b)

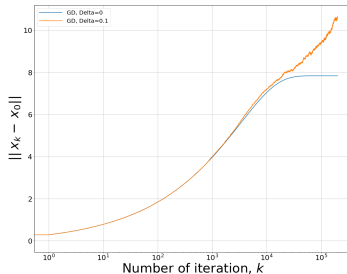


Рис. 3: Результаты для градиентного метода по норме градиента без использования критерия останова для  $\Delta = 0.1$  на задаче минимизации логистической регрессии для  $v = -\frac{\Delta}{\|\nabla f(x)\|} \nabla f(x)$ : (a) скорость сходимости по норме градиента; (b) расстояние от начальной точки до  $x_k$ .



# Численные эксперименты

## Логистическая регрессия

	Alg. constant			Adaptive $L$		
$\Delta$	Iters	Time, ms	$\frac{\ \tilde{\nabla} f(\hat{x})\ }{\Delta}$	Iters	Time, ms	$\frac{\ \tilde{\nabla} f(\hat{x})\ }{\Delta}$
$10^{-5}$	—	—	—	902	2856.98	2.22
$10^{-4}$	9700	2605.02	2.42	472	1678.02	2.16
$10^{-2}$	83	49.86	2.29	17	68.36	2.10

	Full Adaptive		
$\Delta$	Iters	Time, ms	$\frac{\ \tilde{\nabla} f(\hat{x})\ }{\Delta}$
$10^{-5}$	23	449.68	3.37
$10^{-4}$	25	370.09	3.62
$10^{-2}$	17	161.27	0.84

Таблица 3: Результаты алгоритмов для (42) с различными величинами  $\Delta$  при использовании правила  $\|\tilde{\nabla} f(x)\| \leq \sqrt{6}\Delta$ .

# Функция Розенброка

Исследуем скорость сходимости алгоритма 1 на двумерной функции Розенброка:

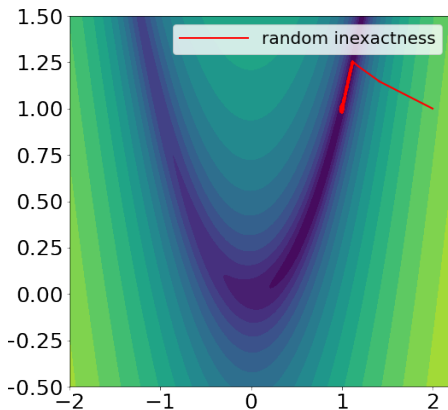
$$f(x_1, x_2) = 100(x_2 - (x_1)^2)^2 + (x_1 - 1)^2.$$

В наших экспериментах мы будем варьировать  $\Delta$  и брать  $\delta = \Delta^2$ . Стартовая точка для всех параметров  $x_1 = 1, x_2 = 2$ . Расстояние от начальной точки до оптимальной **1** равно 1. В таблице 4 представлены результаты сходимости для различных типов шумов. В качестве постоянного шума брался вектор  $v = (1, 0)^T$ . В данном и следующем эксперименте мы будем пользоваться критерием останова

$$\|\tilde{\nabla} f(x_k)\| \leq 2\Delta. \quad (43)$$

Inexactness	$\Delta$	Iters	Time, ms	$\ x_N - x_0\ $	$\frac{\ \nabla f(x_N)\ }{\Delta}$	$f(x_N) - f_*$
Random	$10^{-4}$	7266	2273.26	0.999	2.70	$0.89 \cdot 10^{-7}$
	$10^{-3}$	5412	2594.61	0.994	2.90	$1.00 \cdot 10^{-5}$
	$10^{-2}$	3690	3163.32	0.940	2.61	$0.89 \cdot 10^{-3}$
Antigradient	$10^{-4}$	7188	2615.61	0.999	2.99	$0.11 \cdot 10^{-6}$
	$10^{-3}$	5493	2490.56	0.993	3.00	$0.11 \cdot 10^{-4}$
	$10^{-2}$	3536	3031.05	0.931	2.99	$0.12 \cdot 10^{-2}$
Constant	$10^{-4}$	7491	2301.32	1.000	1.54	$0.27 \cdot 10^{-7}$
	$10^{-3}$	5697	2490.32	0.997	1.87	$0.24 \cdot 10^{-5}$
	$10^{-2}$	3965	3485.89	0.965	1.93	$0.30 \cdot 10^{-3}$

Таблица 4: Результаты работы адаптивного градиентного спуска для двумерной функции Розенброка при использовании условия останова (43).



**Рис. 4:** Траектория градиентного метода со случайной неточностью в погрешности градиента для функции Розенброка без условия останова. Градиентный метод, начиная с некоторой итерации, начинает осцилировать вокруг некоторой точки, и дальнейшие итерации бессмысленны.

## Функция Нестерова-Скокова

Рассмотрим систему нелинейных уравнений  $g(x) = 0$ , где  $g_1 = \frac{1}{2}(x_1 - 1)$ ,  $g_i = x_i - 2x_{i-1}^2 + 1$ ,  $i = \overline{2, n}$ . Задача решения этой системы эквивалентна минимизации функции Нестерова-Скокова<sup>11</sup>

$$f(x) = \frac{1}{4}(1 - x_1)^2 + \sum_{i=1}^{n-1} (x_{i+1} - 2x_i^2 + 1)^2. \quad (44)$$

Данная функция является обобщением функции Розенброка. Она так же является невыпуклой и удовлетворяет условию Липшица для градиента только локально. Также функция (44) обладает глобальным минимумом в точке  $(1, 1 \dots 1, 1)^T$  и оптимальным значением  $f^* = 0$ .

---

<sup>11</sup>Нестеров Ю. Е., Скоков В. А. К вопросу о тестировании алгоритмов безусловной оптимизации. — в книге “Численные методы математического программирования” — М.: ЦЭМИ, 1980. — С. 77–91.

Как было показано в предыдущих экспериментах, предложенный нами критерий останова одинаково хорошо работает для всех опробованных типов шумов. В данном эксперименте мы будем использовать только шум, равномерно распределенный на сфере. Все наши эксперименты мы будем начинать со стартовой точки  $(-1, 1, \dots, 1, 1)^T$ . Таким образом,  $\|x_0 - x_*\| = 2$ . Мы будем варьировать точность градиента  $\Delta$  и размерность задачи  $n$ .

$n$	$\Delta$	Iters	Time, ms	$\ x_N - x_0\ $	$\frac{\ \nabla f(x_N)\ }{\Delta}$	$f(x_N) - f^*$
3	$10^{-4}$	14097	230.58	1.996	2.86	$0.20 \cdot 10^{-4}$
	$10^{-3}$	2477	247.64	2.155	2.93	$0.11 \cdot 10^{-2}$
	$10^{-2}$	606	383.63	2.650	2.19	$0.87 \cdot 10^{-2}$
5	$10^{-4}$	73028	275.03	2.930	2.93	$0.30 \cdot 10^{-3}$
	$10^{-3}$	15765	292.39	3.312	2.65	$0.49 \cdot 10^{-2}$
	$10^{-2}$	6	200.87	0.036	1.45	0.98
7	$10^{-4}$	2898	316.51	0.049	2.69	0.98
	$10^{-3}$	103	164.23	0.036	2.07	0.98
	$10^{-2}$	17	104.77	0.036	1.42	0.98

Таблица 5: Результаты работы адаптивного градиентного спуска для функции Нестерова-Скокова при использовании условия останова (43).

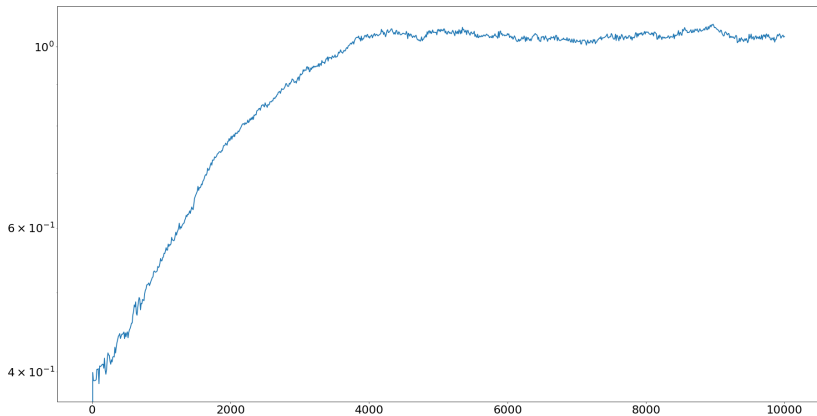
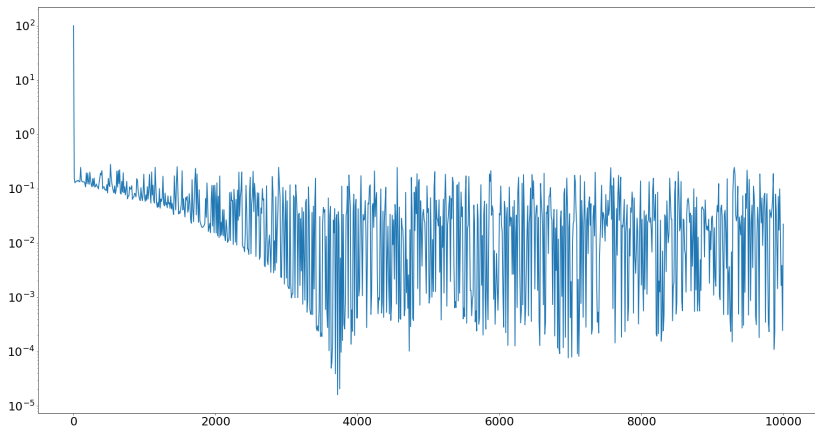


Рис. 5: Расстояние  $\|x_N - x_0\|$  без критерия останова выходит в какой-то момент просто на плато.





**Рис. 6:** Это плато приблизительно одинаково для всех уровней шума, и различается только насколько сильно на этом плато значение колеблется. Похожим образом ведет себя и величина  $f(x_N) - f^*$ .

# Численные эксперименты

## Нелинейная система

Рассмотрим задачу решения следующей системы нелинейных уравнений

$$g_i(x) = \sum_{j=1}^n A_{ij} \sin(x_j) + B_{ij} \cos(x_j) := E_i, \quad i = 1, \dots, m, \quad (45)$$

где  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$  и

$A_{ij}, B_{ij} \in \mathbb{R}, \forall i = 1, \dots, m; j = 1, \dots, n$ . Задача выше может быть переписана в виде следующей задачи оптимизации:

$$\min_{x \in \mathbb{R}^n} f(x) := \sum_{i=1}^m (g_i(x) - E_i)^2. \quad (46)$$

# Численные эксперименты

## Нелинейная система

$m$	$\frac{L}{\mu}$	$\Delta$	Alg. constant		Adaptive $L$		Full Adaptive	
			Dist	Grad	Dist	Grad	Dist	Grad
8	$2.1 \cdot 10^5$	$10^{-4}$	5.1	2.41	5.1	2.25	5.1	3.07
		$10^{-1}$	4.8	2.34	4.8	2.09	4.9	2.28
32	$5.0 \cdot 10^6$	$10^{-4}$	7.4	2.41	7.5	2.27	7.5	3.31
		$10^{-1}$	7.4	2.37	7.4	2.24	7.4	2.72
128	$7.7 \cdot 10^8$	$10^{-4}$	14.4	2.43	14.4	2.33	14.4	15.67
		$10^{-1}$	14.3	2.43	14.3	2.32	14.4	1.84

Таблица 6: Результаты для задачи (46) для различных значений  $\Delta$ .

Мы сравнили расстояние  $\|x_n - x_0\|$  и величину  $\frac{\|\widehat{\nabla} f(x_N)\|}{\Delta}$ .

# Пока нерешённые задачи

- ▶ Стохастические варианты градиентного метода и правила остановки
- ▶ Бесконечномерные задачи
- ▶ Тестирование на реальных задачах обучения

## Заключение

1. Предложены критерии ранней остановки для градиентного метода, гарантирующее с нашей точки зрения компромисс между стремлением достичь приемлемого качества точки выхода по функции и целью обеспечить не столь существенное удаление этой точки от выбранного начального положения.
2. Рассмотрен как градиентный метод с постоянным шагом, так и его вариант с адаптивной регулировкой шага, что полезно в ситуации неизвестного (или даже бесконечного) значения параметра  $L$ .
3. Получены теоретические оценки количества итераций для достижения приемлемого качества приближённого решения.
4. Проведены численные эксперименты, демонстрирующие работу критерия останова на практике для задач логистической регрессии, минимизации функции Розенброка и её многомерного аналога (функции Нестерова-Скокова).

Спасибо за внимание!