
ПРАВИЛА ОСТАНОВКИ ГРАДИЕНТНОГО МЕТОДА ДЛЯ СЕДЛОВЫХ ЗАДАЧ С ДВУСТОРОННИМ УСЛОВИЕМ ПОЛЯКА-ЛОЯСИЕВИЧА

ПРЕПРИНТ

Муратиди Александр Янисович
МФТИ
Москва, Россия
`muratidi.aia@phystech.edu`

Стонякин Федор Сергеевич
МФТИ
Москва, Россия
КФУ им. В. И. Вернадского
Симферополь, Россия
`fedyor@mail.ru`

Аннотация

Статья посвящена некоторым вопросам, связанным с численными методами для седловых задач с двусторонним условием градиентного доминирования Поляка–Лоясиевича [2]. Известно, что на классе достаточно гладких минимизационных задач с двусторонним условием градиентного доминирования градиентный метод сходится со скоростью геометрической прогрессии, что считается хорошим уровнем. Отметим, что в последнее время повысился интерес к этому классу задач ввиду приложений в нелинейных перепараметризованных системах глубокого обучения с перепараметризацией [3]. Существуют и примеры седловых задач такого типа, возникающие в анализе данных (специальные робастные вариации метода наименьших квадратов) [4]. В работе рассматриваются подходы к седловым задачам с двусторонним вариантом условия Поляка–Лоясиевича на базе градиентного метода с неточной информацией и предлагаются правило остановки на основе малости нормы неточного градиента внешней подзадачи. Достигение этого правила в сочетании с подходящей точностью решения вспомогательной подзадачи гарантирует достижение приемлемого качества исходной седловой задачи. Обсуждаются результаты численных экспериментов для различных седловых задач для иллюстрации эффективности предложенного метода, в том числе по сравнению с доказанными оценками скорости сходимости.

1 Постановка задачи

Будем рассматривать задачи нахождения седловой точки (x^*, y^*)

$$f(x^*, y^*) = \min_x \max_y f(x, y), \quad (1)$$

где f удовлетворяет двустороннему варианту условия Поляка–Лоясиевича (PL-условию),

$$\begin{cases} \|\nabla_x f(x, y)\|^2 \geq 2\mu_1(f(x, y) - \min_x f(x, y)) \\ \|\nabla_y f(x, y)\|^2 \geq 2\mu_2(\max_y f(x, y) - f(x, y)) \end{cases} \quad (2)$$

а также условиям Липшица градиента относительно евклидовой нормы.

$$\begin{cases} \|\nabla_x f(x_1, y) - \nabla_x f(x_2, y)\| \leq L_{11} \|x_1 - x_2\| \\ \|\nabla_x f(x, y_1) - \nabla_x f(x, y_2)\| \leq L_{12} \|y_1 - y_2\| \\ \|\nabla_y f(x, y_1) - \nabla_y f(x, y_2)\| \leq L_{22} \|x_1 - x_2\| \\ \|\nabla_y f(x_1, y) - \nabla_y f(x_2, y)\| \leq L_{12} \|x_1 - x_2\| \end{cases} \quad (3)$$

для всяких x и y .

Напомним, что для задач минимизации без ограничений

$$\min_x f(x), \quad x \in \mathbb{R}^n$$

условие градиентного доминирования Поляка-Лоясиевича (далее условие PL), которое гарантирует сходимость градиентного спуска со скоростью геометрической прогрессии для достаточно гладких задач без дополнительных предположений о выпуклости функции. При этом в оценки скорости сходимости не входит параметр размерности пространства. Точнее говоря, функция $f(\cdot)$ удовлетворяет условию PL, если она имеет непустое множество решений и конечное оптимальное значение f^* , а также существует некоторая $\mu > 0$ такая, что

$$\|\nabla f(x)\|^2 \geq 2\mu(f(x) - f^*), \quad \text{для всех } x.$$

Известно, что:

- PL-условие является более слабым по сравнению с сильной выпуклостью, т.е. оно верно для более широкого класса функций.
- В частности, PL-условие верно для невыпуклых задач, связанных с перепараметризованными системами в глубоком обучении. [3]

2 Примеры задач рассматриваемого класса

Существует некоторое количество классов задач, удовлетворяющих описанным выше условиям, однако для отдельно взятой подходящей функции вычисление параметров L_{ij} и μ_i тоже может являться достаточно сложной задачей. Давайте рассмотрим несколько примеров удовлетворяющих требуемым ограничениям.

2.1 Пример 1

$f(x,y) = F(Ax, By)$, где $F(\cdot, \cdot)$ сильно-выпукло-сильно-вогнутая и A, B — произвольные матрицы, удовлетворяет двустороннему PL-условию.

Пример взят из статьи [2].

2.2 Пример 2

Невыпукло-невогнутая $f(x,y) = x^2 + 3\sin^2 x \sin^2 y - 4y^2 - 10\sin^2 y$, удовлетворяет двустороннему PL-условию с $\mu_1 = 1/16, \mu_2 = 1/14$ и условиям Липшица градиента $L_{11} = 8, L_{22} = 28$.

Пример взят из статьи [2].

2.3 Пример 3

Задача Robust least squares (RLS) состоит в минимизации суммы квадратов расхождения между предсказанными значениями и фактическими значениями целевой переменной, но с учетом возможных выбросов в данных. Другими словами, RLS ищет оптимальную модель, которая максимально точно описывает данные, но при этом устойчива к наличию шумов и выбросов в используемых данных. Задача RLS часто используется в машинном обучении и статистике для построения устойчивых моделей регрессии и классификации.

Можно рассмотреть следующий пример задачи RLS [4]:

$$F(x,y) := \|Ax - y\|_M^2 - \lambda\|y - y_0\|_M^2, \quad \text{где } \|\cdot\|_M^2 = x^T M x,$$

где M положительно полуопределенна и $\lambda > 1$. В статье [4] показано, что такая задача удовлетворяет двустороннему PL-условию.

3 Цель статьи

Для решения таких задач часто используется градиентный метод с возможностью использования неточной информации о градиенте, который будет подробно рассмотрен в следующей главе. Теоретические оценки количества итераций, необходимых для достижения приемлемого качества решения с помощью такого градиентного метода в общем случае неулучшаемы. Однако, по-видимому, на практике они предполагают значительное завышение затрат в сравнении с реальной необходимостью. Мы отправляемся от статьи [5], в которой уже было предложено правило ранней остановки для задач минимизации. Однако для седловых задач такого еще нет. Поэтому, основной целью данной работы стала разработка правил ранней остановки градиентных методов для седловых задач, которые могли бы гарантировать достижение приемлемого качества точки выхода по функции, а также провести эксперименты на различных седловых задачах и сравнить численные результаты предложенного метода с доказанными оценками, чтобы убедиться в эффективности предложенных методов.

4 Градиентный метод с неточной информацией и его использование в седловых задачах

Считаем, что значения параметров $L_{11}, L_{12}, L_{21}, L_{22}, \mu_1, \mu_2, \gamma$, которые будут задавать точность решения внутренней подзадачи, известны и положительны.

Будем сводить поставленную задачу к задаче минимизации вспомогательной функции вида

$$g(x) = \max_y f(x, y). \quad (4)$$

Тогда согласно лемме A.5 [1] получаем

$$\|\nabla g(x_1) - \nabla g(x_2)\| \leq L \|x_1 - x_2\| \quad (5)$$

при

$$L = L_{11} + \frac{L_{12}^2}{\mu_2}. \quad (6)$$

При этом обычный градиентный спуск в этом случае не применим. Это связано с тем, что на практике возможно, как правило, лишь приближенно для всякого x найти значение градиента $g(x)$. Поэтому считаем, что в любой точке x нам доступно значение неточного градиента $\tilde{\nabla}g(x) = \nabla_x f(x, y)$, причем $\|\tilde{\nabla}g(x) - g(x)\| \leq \Delta$, при некотором фиксированном положительном $\Delta = L_{12}\gamma$.

Тогда (5) означает, что

$$g(x) - g^* \leq \frac{1}{\mu_1} (\|\tilde{\nabla}g(x)\|^2 + \Delta^2) \quad \forall x \in \mathbb{R}^n, \quad (7)$$

поэтому для всякого x верно

$$\|\tilde{\nabla}g(x)\|^2 \geq \mu_1(g(x) - g^*) - \Delta^2.$$

К задаче минимизации g будем применять градиентный метод вида

$$x_{k+1} = x_k - \frac{1}{L} \tilde{\nabla}g(x_k) \quad (8)$$

для «внешней задачи» и

$$y_{l+1} = y_l + \frac{1}{L_{22}} \nabla_y f(x_k, y_l) \quad (9)$$

для каждой из «внутренних» подзадач, необходимых для нахождения приближенного значения градиента для внешней задачи с заданной точностью решения по аргументу.

Если $\|y_k - y^*\| \leq \gamma$ и γ достаточно мал, а также

$$g(x_k) = \max_y f(x_k, y) = f(x_k, y_k^*),$$

тогда имеем:

$$\|\tilde{\nabla}g(x_k) - \nabla g(x_k)\| = \|\nabla f(x_k, y_k) - \nabla_x f(x_k, y_k^*)\| \leq L_{12} \|y_k - y_k^*\| \leq L_{12}\gamma. \quad (10)$$

Таким образом, в (10) $\tilde{\nabla}g(x_k)$ представляет собой аддитивно неточный градиент g в точке x_k с $\Delta = L_{12}\gamma$, а γ зависит от ошибки решения вспомогательной проблемы для y .

Лемма А3 в [2] (стр. 16) утверждает, что

$$\|\nabla g(x)\|^2 \geq 2\mu_1(g(x) - g(x^*)),$$

т.е. g удовлетворяет обычному условию Поляка-Лоясиевича.

Получим оценку на количество итераций для «внутренней» подзадачи, гарантирующее приемлемую её точность.

Из известного результата сходимости со скоростью геометрической прогрессии по аргументу для L-гладких функций, удовлетворяющих PL-условию (а значит, и квадратичному росту) получаем, что

$$\begin{aligned} \|y_k - y_k^*\|^2 &\leq \frac{2}{\mu_2}(f(x_k, y_k^*) - f(x_k, y_k)) \leq \\ &\leq \frac{2}{\mu_2}(1 - \frac{\mu_2}{L_{22}})^p (f(x_k, y_k^*) - f(x_k, y_0)) \leq \frac{2}{\mu_2}(1 - \frac{\mu_2}{L_{22}})^p C_2 \end{aligned}$$

для некоторого $C_2 > 0$.

Если потребовать, чтобы

$$L_{12}^2\gamma^2 \geq \frac{2}{\mu_2}(1 - \frac{\mu_2}{L_{22}})^p C_2,$$

тогда

$$(1 - \frac{\mu_2}{L_{22}})^p \leq \exp^{-\frac{\mu_2 p}{L_{22}}} \leq \frac{L_{12}^2\gamma^2\mu_2}{2C_2},$$

и при

$$p \geq \left\lceil \frac{L_{22}}{\mu_2} \log \frac{2C_2}{L_{12}^2\gamma^2\mu_2} \right\rceil \quad (11)$$

получаем требуемую погрешность по аргументу.

Продолжим для «внешней» подзадачи. Ввиду (5) для метода (8) получим следующее соотношение:

$$\begin{aligned} g(x_{k+1}) &\leq g(x_k) + \langle \nabla g(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 = \\ &= g(x_k) - \frac{1}{L} \langle \nabla g(x_k), \tilde{\nabla}g(x_k) \rangle + \frac{1}{2L} \|\tilde{\nabla}g(x_k)\|^2 = \\ &= g(x_k) + \frac{1}{2L} (\nabla g(x_k)^2 - 2\langle \nabla g(x_k), \tilde{\nabla}g(x_k) \rangle + \tilde{\nabla}g(x_k)^2) - \frac{\|\nabla g(x_k)\|^2}{2L} = \\ &= g(x_k) + \frac{1}{2L} \|\nabla g(x_k) - \tilde{\nabla}g(x_k)\|^2 - \frac{\|\nabla g(x_k)\|^2}{2L} \leq \\ &\leq g(x_k) + \frac{\Delta^2}{2L} - \frac{1}{2L} \|\nabla g(x_k)\|^2, \end{aligned}$$

т. е.

$$g(x_{k+1}) - g(x_k) \leq \frac{\Delta^2}{2L} - \frac{1}{2L} \|\nabla g(x_k)\|^2. \quad (12)$$

Из пункта 3.4 [5]:

$$\begin{aligned} g(x_k) - g(x^*) &\leq (1 - \frac{\mu_1}{L})^k (g(x_0) - g(x^*)) + \frac{\Delta^2}{2\mu_1} = \\ &= (1 - \frac{\mu_1}{L})^k (g(x_0) - g(x^*)) + \frac{L_{12}^2\gamma^2}{2\mu_1} = \\ &= (1 - \frac{\mu_1\mu_2}{L_{11}\mu_2 + L_{12}^2})^k (g(x_0) - g(x^*)) + \frac{L_{12}^2\gamma^2}{2\mu_1} = \end{aligned}$$

$$= \left(1 - \frac{\mu_1\mu_2}{L_{11}\mu_2 + L_{12}^2}\right)^k C_1 + \frac{L_{12}^2\gamma^2}{2\mu_1} \quad (13)$$

для некоторого $C_1 > 0$, тогда

$$\begin{cases} C_1(1 - \frac{\mu_1}{L})^k \leq \exp^{-\frac{\mu_1}{L}k} C_1 \leq \frac{\varepsilon}{2}, \\ \frac{L_{12}^2\gamma^2}{2\mu_1} \leq \frac{\varepsilon}{2}. \end{cases} \quad (14)$$

Отсюда получаем оценку на количество итераций для внешней подзадачи:

$$\begin{aligned} \exp^{-\frac{\mu_1}{L}k} &\leq \frac{\varepsilon}{2C_1}, \\ -\frac{\mu_1}{L}k &\leq -\log \frac{2C_1}{\varepsilon}, \\ k &\geq \left\lceil \frac{L}{\mu_1} \log \frac{2C_1}{\varepsilon} \right\rceil. \end{aligned}$$

А также $L_{12}^2\gamma^2 \leq \varepsilon\mu_1$, и тогда можно переписать оценку количества итераций для решения внутренней подзадачи:

$$p \geq \left\lceil \frac{L_{22}}{\mu_2} \log \frac{2C_2}{\varepsilon\mu_1\mu_2} \right\rceil.$$

Итого при

$$N = k \cdot p > \left\lceil \frac{L}{\mu_1} \log \frac{2C_1}{\varepsilon} \right\rceil \cdot \left\lceil \frac{L_{22}}{\mu_2} \log \frac{2C_2}{\varepsilon\mu_1\mu_2} \right\rceil$$

мы достигнем условия $g(x_k) - g(x^*) \leq \varepsilon$.

Такой подход к решению для седловых задач на базе градиентного метода с неточной информацией является известным, однако нам неизвестны случаи, в которых аккуратно были бы прописаны оценки сложности этой схемы для седловых задач с двусторонним аналогом PL-условия.

5 Правила остановки

В предыдущей главе мы рассмотрели подход к решению седловых задач такого типа на базе градиентного метода с неточной информацией. В этой же будет предложен критерий ранней остановки на базе малости нормы неточного градиента внешней подзадачи, который гарантирует достижение приемлемого качества решения исходной седловой задачи. А также критерий остановки для внутренних подзадач, гарантирующий достижение требуемого качества решения по аргументу. Будет дана оценка количества итераций, необходимых для их реализации.

5.1 Критерий остановки для внешней подзадачи

Из (12) видно

$$\|\nabla g(x_k)\|^2 \geq \frac{\|\tilde{\nabla}g(x_k)\|^2}{2} - L_{12}^2\gamma^2$$

верно

$$f(x_{k+1}) - f(x_k) \leq \frac{L_{12}^2\gamma^2}{2L} - \frac{1}{2L} \left(\frac{\|\tilde{\nabla}g(x_k)\|^2}{2} - L_{12}^2\gamma^2 \right),$$

откуда имеем

$$g(x_{k+1}) - g(x_k) \leq \frac{L_{12}^2\gamma^2}{L} - \frac{1}{4L} \|\tilde{\nabla}g(x_k)\|^2. \quad (15)$$

Из неравенства (15) видно, что если значение $\|\tilde{\nabla}g(x_k)\|$ достаточно велико, то можно гарантировать, что $g(x_{k+1}) < g(x_k)$, что указывает на конечность процесса. Тем самым, для всякого $C > 2$ возникает

альтернатива: или верно неравенство $\|\nabla g(x_k)\| \leq CL_{12}\gamma$, и это гарантирует достижение приемлемого качества точки выхода x_k по функции в силу PL-условия, или же

$$g(x_{k+1}) - g(x_k) < -\frac{L_{12}^2\gamma^2}{L} \left(\frac{C^2}{4} - 1 \right).$$

Тем самым, за конечное число шагов градиентного метода (8) возможно получить x_k такое, что значение $g(x_k)$ достаточно близко к минимальному $g(x^*)$. Выберем для определённости $C = \sqrt{6}$ (чтобы получить «удобный» коэффициент) и будем рассматривать 2 сценария:

1. $\|\tilde{\nabla}g(x_k)\| > L_{12}\gamma\sqrt{6}$, откуда с учетом (15) получаем

$$g(x_{k+1}) - g(x_k) < -\frac{L_{12}^2\gamma^2}{2L}. \quad (16)$$

- 2.

$$\|\tilde{\nabla}g(x_k)\| \leq L_{12}\gamma\sqrt{6}, \quad (17)$$

откуда ввиду (7) имеем

$$g(x_{k+1}) - g(x^*) \leq \frac{7L_{12}^2\gamma^2}{\mu_1}. \quad (18)$$

Будем считать оценку (18) приемлемой и договоримся обрывать процесс (8) в случае, если выполнено (17).

5.2 Критерий остановки для внутренних подзадач

Обозначим $t(y) = f(*, y)$ — внутренняя подзадача для фиксированного x ,

$$\begin{aligned} L_{22}(y_{k+1} - y_k) &= \Delta t(y_k), \\ t(y_{k+1}) &\geq t(y_k) + \langle \Delta t(y_k), y_{k+1} - y_k \rangle - \frac{L_{22}}{2} \|y_{k+1} - y_k\|^2, \\ t(y_{k+1}) - t(y_k) &\geq \frac{1}{L_{22}} \|\nabla t(y_k)\|^2 - \frac{1}{2L_{22}} \|\nabla t(y_k)\|^2, \\ t(y_{k+1}) - t(y_k) &\geq \frac{1}{2L_{22}} \|\nabla t(y_k)\|^2. \end{aligned}$$

Пусть $\|\nabla t(y_k)\| \leq C_3$, тогда из PL-условия и условия квадратичного роста:

$$\frac{\mu_2}{2} \|y_k - y^*\|^2 \leq t^* - t(y_k) \leq \frac{1}{2\mu_2} C_3^2. \quad (19)$$

Ввиду требуемого $\|y_k - y^*\| \leq \gamma$ возьмем $C_3 = \mu_2 \cdot \gamma$.

Таким образом, возникает альтернатива: или градиент становится достаточно мал, и ввиду PL-условия нам удается достигнуть необходимого качества решения по аргументу, либо продолжаем делать шаги градиентного метода.

Также можно использовать \hat{y} с предыдущего шага как начальную точку для следующего. Тогда на практике для задач, в которых требуется большая точность, критерий остановки будет срабатывать «почти сразу».

5.3 Итоговая схема

К седловой задаче для f применяются методы вида

$$x_{k+1} = x_k - \frac{\mu_2}{L_{11}\mu_2 + L_{12}^2} \tilde{\nabla}g(x_k) = x_k - \frac{\mu_2}{L_{11}\mu_2 + L_{12}^2} \nabla_x f(x_k, y_m) \quad (20)$$

для «внешней» задачи, и для вычисления $\tilde{\nabla}g(x_k)$ —

$$y_{m+1} = y_m + \frac{1}{L_{22}} \nabla_y f(x_k, y_m) \quad (21)$$

для каждой из «внутренних» подзадач.

Для каждого из методов (20), (21) предлагаются правила ранней остановки (22), (23) соответственно:

$$\|\nabla_x f(x_k, y_m)\| \leq L_{12} \gamma \sqrt{6}, \quad (22)$$

$$\|\nabla_y f(x_k, y_m)\| \leq \mu_2 \gamma. \quad (23)$$

В таком случае справедливы следующие теоремы.

5.3.1 Теорема 1

Теорема 1. Пусть на p -й итерации градиентного метода (21) впервые выполнен критерий остановки (23). Тогда для точки выхода $\hat{y} = y_p$ гарантированно будет верно неравенство

$$\|y^*(x_k) - \hat{y}\| \leq \gamma. \quad (24)$$

При этом справедлива следующая оценка количества итераций до полной остановки:

$$p \leq \left\lceil \frac{L_{22}}{\mu_2} \log \left(\frac{2C_2}{L_{12}^2 \gamma^2 \mu_2^2} \right) \right\rceil. \quad (25)$$

$$\forall x \mapsto f(x, y^*) - f(x, y_0) \leq C_2 \text{ для некоторого } C_2 \geq 0.$$

Что следует из (11), а также предложенного критерия остановки для внутренней подзадачи.

5.3.2 Теорема 2

Теорема 2. Для градиентного метода (20) критерий остановки (22) выполнен при

$$N \leq \left\lceil 2C_1 \frac{L_{11} + \frac{L_{12}^2}{\mu_2}}{L_{12}^2 \gamma^2} \right\rceil. \quad (26)$$

$$g(x_0) - g(x^*) \leq C_1 \text{ для некоторого } C_1 \geq 0.$$

Пусть критерий остановки (22) не выполнен для всех k от 0 до $N-1$. Тогда $\tilde{\nabla}g(x_k) > \sqrt{6}L_{12}\gamma$. Но тогда, просуммировав выражения (16) для всех k , получим следующее соотношение:

$$g(x_0) - g(x^*) \geq g(x_0) - g(x_N) = \sum_{k=0}^{N-1} (g(x_k) - g(x_{k+1})) > \frac{NL_{12}^2\gamma^2}{2L}.$$

Подставив L из (6), получаем требуемое соотношение.

Эта оценка является сильно завышенной, однако ее удобно использовать при отсутствии информации о параметре μ_1 .

5.3.3 Теорема 3

Теорема 3. Пусть градиентный метод (20) работает либо

$$N_* = \left\lceil \frac{L_{11}\mu_2 + L_{12}^2}{\mu_1\mu_2} \log \left(\frac{C_1\mu_1}{6L_{12}^2\gamma^2} \right) \right\rceil \quad (27)$$

шагов, либо при некотором $N \leq N_*$ на N -й итерации метода (20) впервые выполнен критерий остановки (22). Тогда для точки выхода $\hat{x} = x_N$ гарантированно будет верно неравенство

$$\|\hat{x} - x^*\| \leq \frac{L_{12}\gamma\sqrt{14}}{\mu_1}. \quad (28)$$

$$\text{Для некоторых } C_1, C_2.$$

Пусть также критерий остановки (22) не выполнен для всех k от 0 до $N-1$ и $\tilde{\nabla}g(x_k) > \sqrt{6}L_{12}\gamma$. Тогда ввиду (18) и (13) достаточно потребовать, чтобы

$$\left(1 - \frac{\mu_1\mu_2}{L_{11}\mu_2 + L_{12}^2} \right)^k C_1 \leq 6 \frac{L_{12}^2\gamma^2}{\mu_1},$$

откуда получаем требуемую оценку на N_* . Далее, при выполнении критерии верно (18). И из условия квадратичного роста:

$$\|x_N - x^*\|^2 \leq \frac{2}{\mu_1} (g(x_N) - g(x^*)) \leq 14 \frac{L_{12}^2\gamma^2}{\mu_1^2}.$$

Извлекая квадратный корень из выражений, получим требуемое.

6 Вычислительные эксперименты

Будем проводить эксперимент на следующем примере:

$$f(x,y) = x_1^2 + x_2^2 + x_3^2 + 3 \sin^2 x_1 \sin^2 y_1 - 4y_1^2 - 3y_2^2 - 2y_3^2 - 10 \sin^2 y_1 \quad (29)$$

Функция (29) удовлетворяет необходимым условиям с коэффициентами

- $L_{11} = 8$,
- $L_{22} = 28$,
- $\mu_1 = 1/16$,

- $\mu_2 = 1/14$.

Константы $L_{12}, L_{21} \leq L_{22} = 28$. Также зафиксируем константы $C_1 = C_2 = 100$. Для данной задачи оптимальным решением является точка $\vec{0}$.

Результаты при случайному выборе y .

γ	$f(\hat{x}, \hat{y}) - f^*$	N	$\sum p_k$	avg_p	N^*	p
10^{-3}	$1.1749 \cdot 10^{-3}$	40921	3921481	95.83	1751204	6950
10^{-5}	$1.1752 \cdot 10^{-7}$	69525	8737867	125.67	3369866	10560
10^{-8}	$1.1755 \cdot 10^{-13}$	103487	17644048	170.49	5797859	15976

Результаты при выборе $y_0 = y_{opt}$ на предыдущем шаге.

γ	$f(\hat{x}, \hat{y}) - f^*$	N	$\sum p_k$	avg_p	N^*	p
10^{-3}	$1.1752 \cdot 10^{-3}$	45266	80367	1.77	1751204	6950
10^{-5}	$1.1753 \cdot 10^{-7}$	65430	129758	1.98	3369866	10560
10^{-8}	$1.1755 \cdot 10^{-13}$	106811	237818	2.23	5797859	15976

Здесь p_k — это количество шагов градиентного метода внутренней подзадачи на k -ой итерации внешней, а avg_p — среднее шагов внутренних на одной итерации внешней.

По итогам экспериментов видно, что градиентный метод с критериями остановки намного эффективнее, чем выполнение теоретически необходимого количества итераций. Если считать количество суммарных операций взятия градиента самой ресурсоемкой частью алгоритма, то теоретически потребовалось бы провести $N^* \cdot (p+1)$ таких операций, тогда как на практике требуемая точность достигается при $N + \sum p_k$, что более чем в 10^4 раз меньше в случае со случайным y_0 и в 10^5 раз — при использовании информации с предыдущего запуска.

Список литературы

- [1] Nouiehed M., Sanjabi M., Huang T., Lee J. D., Razaviyayn M. Solving a Class of Non-Convex Min-Max Games Using Iterative First Order Methods // Advances in Neural Information Processing Systems, 2019. <https://arxiv.org/pdf/1902.08297.pdf>
- [2] Yang J., Kiyavash N., He N. Global Convergence and Variance-Reduced Optimization for a Class of Nonconvex-Nonconcave Minimax Problems // Advances in Neural Information Processing Systems, 2020. <https://arxiv.org/pdf/2002.09621.pdf>
- [3] Belkin M. Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation // arXiv preprint (2021). <https://arxiv.org/pdf/2105.14368.pdf>
- [4] Garg K., Baranwal M. Fixed-Time Convergence for a Class of Nonconvex-Nonconcave Min-Max Problems // Eighth Indian Control Conference (ICC), 2022. <https://arxiv.org/pdf/2207.12845.pdf>
- [5] Stonyakin F., Kuruzov I., Polyak B. Stopping rules for gradient methods for non-convex problems with additive noise in gradient // Journal of Optimization Theory and Applications, 2023. <https://link.springer.com/article/10.1007/s10957-023-02245-w>