

# Распределенная оптимизация в условиях Поляка-Лоясевича\*

И. О. Автор<sup>1</sup>, И. О. Соавтор<sup>2</sup>, И. О. Фамилия<sup>1,2</sup>  
author@site.ru; co-author@site.ru; co-author@site.ru

<sup>1</sup>Организация, адрес; <sup>2</sup>Организация, адрес

В статье рассматривается новый метод децентрализованного распределенного решения больших систем нелинейных уравнений в условиях Поляка-Лоясевича. Суть метода состоит в постановке эквивалентной задачи распределенной оптимизации и последующем ее сведении сперва к задаче ограниченной оптимизации, а затем к задаче композитной оптимизации, но уже без ограничений. Предложенный метод сравнивается с градиентным спуском, ускоренным градиентным спуском, а также с последовательным и параллельным алгоритмом обратного распространения ошибки при обучении многослойной нейронной сети с нелинейной функцией активации нейрона.

**Ключевые слова:** большие нелинейные системы; распределенная оптимизация; условия Поляка-Лоясевича; многослойные нейронные сети

DOI: 10.21469/22233792

## 1 Введение

За последние несколько лет наблюдается скачок популярности задач, связанных с анализом больших данных. Это вызвано увеличением количества параметров моделей как в машинном обучении, так и в других областях прикладной математики.

**Alexander B.:** Примеры можно дополнить ссылками на статьи (cite) или просто ссылками в интернете - в колонке

Примерами таких задач могут послужить задачи обработки непрерывно поступающих данных с измерительных устройств, аудио и видеоматериалов; изучения потоков сообщений в социальных сетях или метеорологических данных; анализа данных о местонахождении абонентов сетей и оптимальное распределение мощности между вышками сотовой связи.

**Alexander B.:** Про системы тоже можно в отдельный абзац и с примерами

Суть многих из этих задач заключается в решении огромных систем нелинейных уравнений.

В этой статье мы рассматриваем новый способ решения системы нелинейных уравнений:

$$\begin{cases} f_1(x_1, \dots, x_n) = 0 \\ f_2(x_1, \dots, x_n) = 0 \\ \dots \\ f_m(x_1, \dots, x_n) = 0 \end{cases} \quad f_i : \mathbb{R}^n \rightarrow \mathbb{R}, \quad i = 1, \dots, m \quad (1)$$

Эту систему можно переписать в виде эквивалентной задачи оптимизации:

$$g(x) := \sum_{i=1}^m f_i^2(x) \rightarrow \min_{x \in \mathbb{R}^n} \quad (2)$$

\*Работа выполнена при финансовой поддержке РФФИ, проекты № №00-00-00000 и 00-00-00001.

В связи с тем, что мы имеем дело с огромными системами, то возникает мысль решать их распределенно. В этой статье мы рассмотрим децентрализованную оптимизацию. Схематически децентрализованную систему можно представить как несколько устройств или процессоров, которые связаны в сеть, при этом какие-то устройства связаны между собой, а какие-то нет, соответственно информацией могут обмениваться только те, между которыми есть канал связи. Примерами таких систем могут служить архитектуры из нескольких видеокарт или сеть из нескольких компьютеров.

В [1] задача (2) представляется в децентрализованном виде:

$$\begin{aligned} \min_{x_i \in \mathbb{R}^n} \quad & \sum_{i=1}^m f_i^2(x_i) \\ \text{s.t.} \quad & x_1 = x_2 = \dots = x_m, \end{aligned} \quad (3)$$

где  $x_i$  – это копия переменной  $x$  на каждом устройстве. Мы хотим, чтобы на каждом устройстве было одинаковое значение  $x$  – одинаковое решение, поэтому и вводится соответствующее ограничение.

Таким образом, возникает задача ограниченной оптимизации, которую авторы статьи решают методом прямодвойственного градиентного спуска. Причем, если функция  $g$  удовлетворяет условиям Поляка-Лоясиевича с константой  $\nu > 0$ , то есть:

$$\frac{1}{2} \|\nabla g(x)\|^2 \geq \nu(g(x) - g^*), \quad \forall x \in \mathbb{R}^n; \quad g^* = \min_{x \in \mathbb{R}^n} g(x) \quad (4)$$

то метод будет иметь линейную скорость сходимости.

Мы, в свою очередь, сводим задачу ограниченной оптимизации к задаче композитной оптимизации, смягчив жесткие условия на совпадение  $x_{i=1}^m$  в задаче (3). И предлагаем решать полученную задачу аналогами метода подобных треугольников или слайдинга [4].

Наш метод сравнивается с методами градиентного спуска и ускоренного градиентного спуска, описанными в [2] и [5]. Также мы сравнили наш метод с последовательным и параллельным вариантами самого распространенного на данный момент алгоритма обратного распространения ошибки для обучения нейронных сетей с нелинейной функцией активации нейрона, предложенными в [3]. Сравнение производится в ходе вычислительного эксперимента при обучении нейронных сетей с различным количеством слоев и функцией активации нейрона – сигмной. Обучение производится на классических данных (CIFAR, MNIST, IMAGNET).

## 2 Постановка задачи

### 2.1 Определения и обозначения

Для скалярного произведения двух векторов  $x, y \in \mathbb{R}^n$  мы будем использовать  $\langle x, y \rangle := \sum_{i=1}^n x_i y_i$ . Скалярное произведение порождает вторую норму  $\ell_2$ -норма в  $\mathbb{R}^n$  в следующем виде  $\|x\|_2 := \sqrt{\langle x, x \rangle}$ . Определим произвольную норму  $\ell_p$  как  $\|x\|_p := (\sum_{i=1}^n |x_i|^p)^{1/p}$  для  $p \in (1, \infty)$ , а для  $p = \infty$  мы используем  $\|x\|_\infty := \max_{1 \leq i \leq n} |x_i|$ . Для максимального и минимального собственного значения положительно определенной матрицы  $A \in \mathbb{R}^{n \times n}$  мы вводим следующие обозначения  $\lambda_{\max}(A)$  и  $\lambda_{\min}^+(A)$  соответственно и под  $\chi(A) := \lambda_{\max}(A) \lambda_{\min}^+(A)$  мы понимаем число обусловленности матрицы  $A$ . Для обозначения произведения Кронекера двух матриц  $A \in \mathbb{R}^{m \times m}$  и  $B \in \mathbb{R}^{n \times n}$  мы используем  $A \otimes B \in \mathbb{R}^{nm \times nm}$ . Единичная матрица размера  $n \times n$  имеет обозначение  $I_n$ . Для неориентированного графа на множестве вершин  $V$  с ребрами  $E$  мы используем  $G(V, E)$ .

**Определение 1 (матрица Кирхгофа).**  $L$  – матрица Кирхгофа графа  $G(V, E)$ , если

$$L_{ij} = \begin{cases} -1, & \text{if } (i, j) \in E, \\ \deg(i), & \text{if } i = j, \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

Введем так же следующую матрицу  $W = L \otimes I_n$

Рассмотрим задачу (3) и для удобства введем обозначения:

Определим  $X$ , как столбец, составленный из векторов аргументов функций  $\{f_i\}_{i=1}^m$ , т.е.

$$X = \text{col}(x_1, \dots, x_m) \quad (6)$$

Обозначим за  $F(X)$  – целевая функция задачи (3):

$$F(X) = \sum_{i=1}^m f_i^2(x_i) \quad (7)$$

**Alexander B.:** Тут я бы сделал таким образом, субсекция "Определения и обозначения в ней стоит ввести основные объекты и обозначения, которые будут использованы в статье, посмотри например <https://arxiv.org/pdf/1911.10645.pdf>

**Alexander B.:** Ввел обозначения, дальше можно начать новую субсекцию и вернуться к проблеме, которую ты заявил в введении. Мол вот такая проблема, но теперь уже не просто описываем насколько она важная и идем от математики, начинаем описывать задачу, например  $f_i$  выпукла (у тебя не так, только пример).

**Alexander B.:** Есть задача оптимизации с ограничениями, да хорошо, но мы хотим убрать ограничения, а как? а вот так, тут можно более подробно, момент не самый простой, поэтому стоит это описать понятно и подробно, пример скинул ВК

## 2.2 Сведение к задаче композитной оптимизации

Перед нами стоит задача решения большой системы нелинейных уравнений (1), которую мы переписали в виде задачи децентрализованной оптимизации (3). Теперь будем минимизировать каждую из функций  $f_i(x_i)$  на отдельном процессоре нашей децентрализованной системы. А связи между этими процессорами будут обеспечивать равенство решений  $\{x_i\}_{i=1}^m$ .

Таким образом, задачу (3) можно представить в виде эквивалентной задачи условной оптимизации [1]:

$$\begin{aligned} \min_{X \in \mathbb{R}^{mn}} \quad & F(X) \\ \text{s.t.} \quad & W^{1/2}X = 0 \end{aligned} \quad (8)$$

Здесь условие  $W^{1/2}X = 0$  эквивалентно условию  $WX = 0$ , которое, в свою очередь, гарантирует совпадение решений на различных процессорах. Такая замена была произведена авторами статьи для доказательства линейной скорости сходимости прямодвойственного градиентного спуска для этой задачи в условиях Поляка-Лоясиевича.

Мы же предлагаем убрать жесткие условия и свести задачу (8) к задаче композитной оптимизации:

$$\min_{X \in \mathbb{R}^{m \times n}} F(X) + R\|W^{1/2}X\| \quad (9)$$

Здесь  $R$  – это некоторая правильно подобранная положительная константа.

Алгоритм решения поставленной задачи основан на методе подобных треугольников, описанном в [4].

### 3 Численный эксперимент

#### 3.1 Скорость сходимости

В первом эксперименте мы решаем линейную систему уравнений:

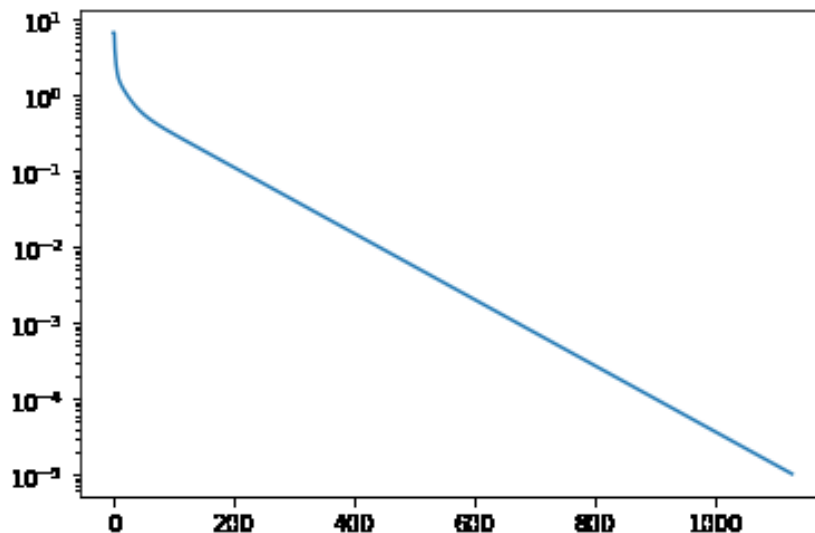
$$Ax = b \quad (10)$$

В виде задачи оптимизации она переписывается следующим образом:

$$\min_x \|Ax - b\|^2 \quad (11)$$

Мы рассматриваем случаи, когда матрица  $A$  – симметричная положительно определенная, семмитричная положительно полуопределенная и произвольная прямоугольная. Для генерации случайных симметричных положительно определенных/полупоредельнных матриц мы используем формулу  $A = Q^T D Q$ , соответственно мы создаем диагональную матрицу  $D$  с нужными нам собственными числами, а с помощью  $QR$ -разложения получаем ортогональную матрицу  $Q$ .

В данном эксперименте на каждом из вычислителей мы использовали градиентный спуск. График сходимости для матрицы размера  $10 \times 10$  представлен на Рис. 1.



**Рис. 1** Сходимость градиентного спуска для распределенного децентрализованного решения задачи (11).

#### 3.2 Анализ ошибки

Для анализа ошибок мы хотим посмотреть, как ведет себя член  $R\|W^{1/2}X\|$  в ходе оптимизации, который отвечает за синхронизацию решений на каждом устройстве с другими. При хорошей работе метода ошибка синхронизации должна быть маленькой (см. Рис. 2).

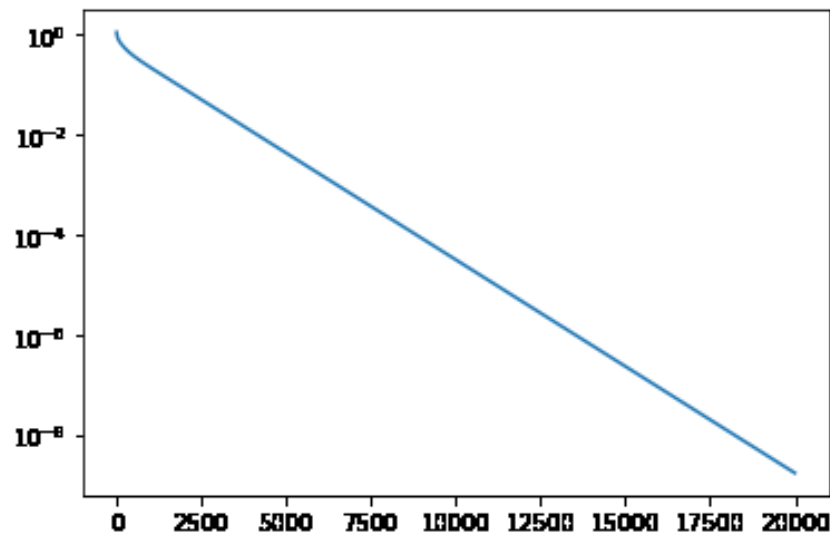


Рис. 2 .

109 Также мы хотим проанализировать как ведет себя метод при разных  $R$  (см. Рис.2).

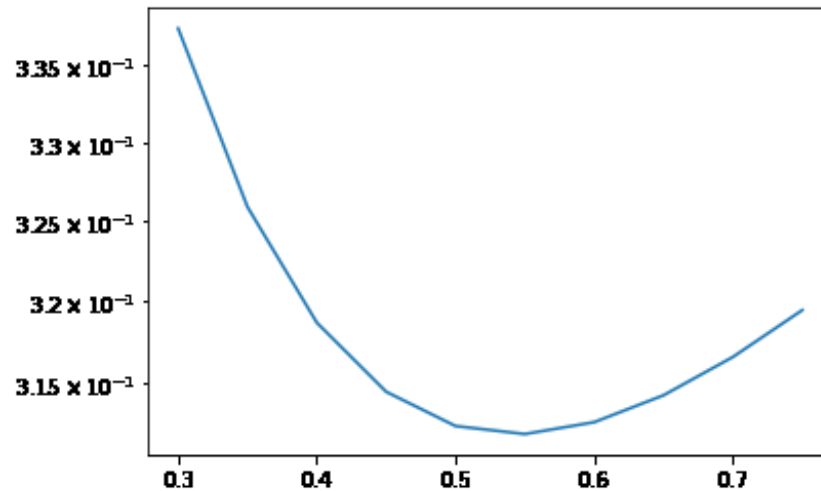


Рис. 3 .

110 На графике показано качество решения через 1000 итераций в зависимости от  $R$ . Мы  
111 видим, что при маленьких  $R$  решение плохое, это связано с плохой синхронизацией, т.к.  
112 член за нее отвечающий слишком мал. При больших  $R$  также наблюдается ухудшение  
113 решения.

## 114 4 Заключение

115 Желательно, чтобы этот раздел был, причём он не должен дословно повторять ан-  
116 нотацию. Обычно здесь отмечают, каких результатов удалось добиться, какие проблемы  
117 остались открытыми.

## 118 Литература

119 [1] Hamed Karimi, Julie Nutini, and Mark W. Schmidt. Linear convergence of gradient and proximal-  
120 gradient methods under the polyak-łojasiewicz condition. *CoRR*, abs/1608.04636, 2016.

- 121 [2] Hamed Karimi, Julie Nutini, and Mark W. Schmidt. Linear convergence of gradient and proximal-  
122 gradient methods under the polyak-łojasiewicz condition. *CoRR*, abs/1608.04636, 2016.
- 123 [3] G. Sandhya Prafulla. Speaker independent vowel recognition using backpropagation neural network  
124 on master-slave architecture jv.s. srinivas,, October 02 2013.
- 125 [4] А. В. Гасников. Универсальный метод для задач стохастической композитной оптимизации.  
126 2016.
- 127 [5] А. В. Гасников. Современные численные методы оптимизации, метод универсального гради-  
128 ентного спуска. 2018.

129 *Поступила в редакцию 01.01.2017*