

---

# Стохастическая кубическая регуляризация для быстрых невыпуклых ОПТИМИЗАЦИЯ

---

Нилеш Трипуранени – Митчелл Стерн – Чи Джин Джеффри Реджир Майкл И. Джордан  
{Nilesh\_tripuraneni, Mitchell, chijin, Regier} @ berkeley.edu  
jordan@cs.berkeley.edu

Калифорнийский университет, Беркли

## Аннотация

В данной работе предлагается стохастический вариант классического алгоритма - кубически-регуляризованный метод Ньютона [Нестеров, Поляк, 2006]. Предложенный алгоритм эффективно избегает седловых точек и находит приближенные локальные минимумы для общих гладких невыпуклых функций только в  $\sim$

$O(\epsilon^{3.5})$  стохастический градиент и стохастик

Гессиано-векторные оценки произведений. Последнее может быть вычислено так же эффективно, как и стохастические градиенты. Это улучшает  $\sim$

$O(\epsilon^4)$  скорость стохастического градиента

спуск. Наша скорость соответствует наиболее известному результату для поиска локальных минимумов, не требуя деликатных методов ускорения или уменьшения дисперсии.

## 1. Введение

Рассмотрим проблему невыпуклой оптимизации в рамках стохастической аппроксимации [Роббинс и Монро, 1951]:

$$\min_{\mathbf{x}} \mathbb{E} \sum_{t=1}^T \ell(\mathbf{x}; \xi_t) = \mathbb{E} \sum_{t=1}^T \ell(\mathbf{x}; \xi_t) \quad (1)$$

В этой настройке мы имеем доступ только к стохастической функции  $\ell(\mathbf{x}; \xi_t)$ , где случайная величина  $\xi_t$  выбирается из основного распределения  $D$ . Задача состоит в том, чтобы оптимизировать ожидаемую функцию  $\ell(\mathbf{x})$ , который в целом может быть невыпуклым. Эта структура охватывает широкий спектр проблем, в том числе настройку, в которой мы минимизируем эмпирические потери при фиксированном объеме данных, и настройку в режиме онлайн, когда данные поступают последовательно. Одним из наиболее важных приложений стохастической оптимизации является крупномасштабная статистика и проблемы машинного обучения, такие как оптимизация глубоких нейронных сетей.

Классический анализ в невыпуклой оптимизации гарантирует только сходимость к стационарной точке первого порядка (т.е. точке  $\mathbf{x}$  удовлетворяющих  $\nabla \ell(\mathbf{x}) = 0$ ), который может быть локальным минимумом, локальным максимумом или седловой точкой. Эта статья идет дальше, предлагая алгоритм, который избегает седловых точек и сходится к локальному минимуму. Локальный минимум определяется как точка  $\mathbf{x}$  удовлетворяющих  $\nabla \ell(\mathbf{x}) = 0$  а также  $\nabla^2 \ell(\mathbf{x}) \succ 0$ . Поиск такой точки представляет особый интерес для большого класса статистических задач обучения, где локальные минимумы являются глобальными или почти глобальными решениями (например, Choromanska и др. [2015], Sun и др. [2016a, b], Ge et al. [2017]).

Среди алгоритмов стохастической оптимизации первого порядка стохастический градиентный спуск (SGD) является, пожалуй, самым простым и наиболее универсальным. В то время как SGD вычислительно недорог, лучшая текущая гарантия для нахождения  $\epsilon$ -приближительный локальный минимум (см. определение 1) требует  $O(\epsilon^{-4} \text{ поли}(\rho))$  итераций [Ge et al., 2015], что неэффективно в многомерном режиме.

---

<sup>\*</sup> Равный вклад

В отличие от методов второго порядка, которые имеют доступ к гессиану  $e$  может использовать кривизну для более эффективного выхода из седел и достижения локальных минимумов. Однако построение полного гессиана может быть чрезмерно дорогим в больших размерах. Таким образом, в **недавней работе было исследовано использование векторов Гессиана**.  $p_2 e(\text{Икс}) \cdot v$ , которые могут быть вычислены так же эффективно, как градиенты во многих случаях, включая нейронные сети [Pearlmutter, 1994].

Среди алгоритмов второго порядка одним из наиболее естественных расширений алгоритма градиентного спуска является кубически-регуляризованный метод Ньютона Нестерова и Поляка [2006]. Принимая во внимание, что градиентный спуск находит минимизатор локального разложения Тейлора второго порядка на каждом этапе,

$$\text{Икс}_{\text{сг}} = \underset{\text{Икс}}{\operatorname{argmin}} e(\text{Икс}) + p e'(\text{Икс}) / (2\eta) + \frac{1}{2} K_{\text{xx}} \tau K_2,$$

кубический регуляризованный метод Ньютона находит минимизатор локального разложения Тейлора третьего порядка,

$$\text{Икс}_{\text{сг}} = \underset{\text{Икс}}{\operatorname{argmin}} e(\text{Икс}) + p e'(\text{Икс}) / (2\eta) + \frac{1}{2} (2\eta)^2 p_2 e'(\text{Икс}) / (2\eta) + \frac{1}{6} K_{\text{xx}} \tau K_3 \eta,$$

Большинство предыдущих работ по кубически-регуляризованному методу Ньютона было сосредоточено на нестохастических или частично стохастических параметрах. Это заставляет нас задать центральные вопросы этой статьи: Можем ли мы разработать полностью стохастический вариант кубически-регуляризованного метода Ньютона? Если так, такой алгоритм быстрее, чем SGD?

В этой работе мы дадим положительный ответ на оба вопроса, сократив разрыв между его использованием в нестохастических и стохастических условиях.

В частности, мы предлагаем стохастический вариант кубически-регуляризованного метода Ньютона. Мы предоставляем неасимптотический анализ его сложности, показывающий, что предложенный алгоритм находит

• стационарная точка второго порядка, использующая только  $\sim O(\epsilon^{3.5})$  стохастический градиент и стохастический гессиан-вектор

оценки, где  $\sim O(\cdot)$  скрывает полилогарифмические факторы. 1 Наш курс улучшается на  $\sim O(\epsilon^4)$  скорость стохастический градиентный спуск и соответствует наиболее известному результату для поиска локальных минимумов без необходимости каких-либо деликатных методов ускорения или уменьшения дисперсии (подробности см. в разделе 1.1). Мы также эмпирически показываем, что стохастический кубически-регуляризованный метод Ньютона, предложенный в этой статье, выгодно работает как на синтетических, так и на реальных невыпуклых задачах относительно современных методов оптимизации.

## 1.1 Связанная работа

В последнее время наблюдается всплеск интереса к методам оптимизации, которые могут избежать седловых точек и найти •-приблизительные локальные минимумы (см. определение 1) в различных условиях. Мы предоставляем краткое резюме этих результатов. Все сложности итераций в этом разделе изложены в терминах поиска •-приблизительных локальных минимумов и только подчеркивают зависимость от •-а также  $d$ .

### 1.1.1 Синглетон-функция

Эта линия работы оптимизируется по общей функции  $e$  без каких-либо особых структурных допущений. В этой настройке алгоритм оптимизации имеет прямой доступ к градиентным или гессианским оракулам на каждой итерации. В работе Нестерова и Поляка [2006] впервые предложен кубически-регуляризованный метод Ньютона, который требует  $O(\epsilon^{-1.5})$  градиент и гессианский оракул призывает весь  $e$ , найти •-

стационарная точка второго порядка. Позже алгоритм ARC [Cartis et al., 2011] и методы области доверия [Curtis et al., 2017] также продемонстрировали, что достигают той же гарантии при аналогичном доступе к гессианскому оракулу. Однако эти алгоритмы полагаются на доступ к полному гессиану на каждой итерации, что непомерно в больших измерениях.

Недавно, вместо использования полного Гессиана, Кармон и Дучи [2016] показали, что использование решателя градиентного спуска для подзадачи кубической регуляризации позволяет их алгоритму найти •- стационарные точки второго порядка в  $\sim$

$O(\epsilon^{-2})$  Гессиано-векторные оценки произведений. С методами ускорения,

число произведений вектора Гессиана может быть уменьшено до  $\sim O(\epsilon^{-1.75})$  [Carmon et al., 2016, Agarwal et al., 2017, Ройер и Райт, 2017].

1 Один запрос для одной реализации  $\leftarrow \leftarrow D, p e'(\text{Икс}, \cdot)$  или  $p_2 e'(\text{Икс}, \cdot) \cdot v$  (для предварительно заданного  $v$ ) упоминается как стохастический градиент или стохастическая оценка оракула по гессиану-вектору.

МЕТОД	время выполнения	Уменьшение дисперсии
Стохастический градиентный спуск [Ge et al., 2015] $O(\epsilon^{-4} \text{ поли}(r))$		не нужно
Наташа 2 [Аллен-Чжу, 2017]	$O(\epsilon^{-3.5})^2$	необходимый
Стохастическая кубическая регуляризация (эта статья)	$O(\epsilon^{-3.5})$	не нужно

Таблица 1: Сравнение наших результатов с существующими результатами для стохастической невыпуклой оптимизации с доказуемой сходимостью к приближенным локальным минимумам.

Между тем, в области результатов, совершенно не связанных с гессианом, Jin et al. [2017] показали, что простой вариант градиентного спуска может найти  $\epsilon$ -вторые стационарные точки в  $\sim$

$$O(\epsilon^{-2}) \text{ оценки градиента.}$$

Обратите внимание, что это направление работы не учитывает стохастические градиенты или стохастические гессианы. Ограничение доступа только к запросам стохастических функций делает проблему оптимизации более сложной.

### 1.1.2 Настройка конечной суммы

В настройке конечной суммы (также известной как установка бесконечной суммы) где  $e(x) = \frac{1}{N} \sum_{s=1}^N e_s(x)$ , один предполагает, что алгоритмы имеют доступ к отдельным функциям  $e_s$ . В этом случае могут использоваться методы уменьшения дисперсии [Johnson and Zhang, 2013]. Agarwal и соавт. [2017] дают алгоритм, требующий  $\sim$

$O(\frac{1}{\epsilon^2} \cdot N \cdot d \cdot \frac{1}{\epsilon^2})$  Гессенский вектор оракула называет (каждый  $k \in [N]$ ) найти  $\epsilon$ -приближенный местный минимум. Аналогичный результат достигается алгоритмом, предложенным Reddi et al. [2017].

### 1.1.3 Стохастическая аппроксимация

Каркас стохастической аппроксимации где  $e(x) = E_{\leftarrow D} [e(\text{Икс}; \leftarrow)]$  предполагает только доступ к стохастическому градиенту и гессиану через  $e(\text{Икс}; \leftarrow)$ . В этом случае Ge et al. [2015] показали, что общая сложность итераций градиента для SGD, чтобы найти  $\epsilon$ -стационарная точка второго порядка была  $O(\epsilon^{-4} \text{ поли}(r))$ .

Совсем недавно, Kohler и Lucchi [2017] рассмотрели субэмплированную версию алгоритма кубической регуляризации, но не предоставляя неасимптотический анализ для их алгоритма, чтобы найти приблизительный локальный минимум; они также предполагают доступ к точным (ожидаемым) значениям функций на каждой итерации, которые недоступны в полностью стохастическом параметре. Сюй и соавт. [2017] рассматривают случай стохастических гессианов, но также требуют доступа к точным градиентам и значениям функций на каждой итерации. Недавно Аллен-Чжу [2017] предложил алгоритм с механизмом, использующим уменьшение дисперсии, который находит стационарную точку второго порядка с  $\sim$

$$O(\epsilon^{-3.5})^2 \text{ Гессиано-векторные оценки произведений.}$$

После этой работы Аллен-Чжу и Ли [2017] и Сюй и Янг [2017] показывают, как использовать градиентные оценки для эффективной аппроксимации векторно-гессианских произведений. Используя эту технику вместе с уменьшением дисперсии, они оба обеспечивают алгоритмы достижения  $\sim$

$$O(\epsilon^{-3.5}) \text{ Оценить с использованием градиентных оценок.}$$

Мы отмечаем, что наш результат соответствует лучшим на сегодняшний день результатам, используя более простой подход без каких-либо деликатных методов ускорения или уменьшения дисперсии. См. Таблицу 1 для краткого сравнения.

## 2 предварительных

Нас интересуют задачи стохастической оптимизации вида  $\min_{\text{Икс}} E_{\leftarrow D} [e(x)] = E_{\leftarrow D} [e(\text{Икс}; \leftarrow)]$ , где  $\leftarrow$  случайная величина с распределением  $D$ . В общем, функция  $e(\text{Икс})$  может быть невыпуклой Эта формулировка охватывает как стандартную установку, где целевая функция может быть выражена в виде конечной суммы  $N$  отдельных функций  $e(\text{Икс}, \leftarrow_s)$ , а также онлайн-настройки, когда данные поступают последовательно.

Наша цель - минимизировать функцию  $e(\text{Икс})$  используя только стохастические градиенты  $\rho e(\text{Икс}; \leftarrow)$  и стохастические гессиано-векторные произведения  $\rho^2 e(\text{Икс}; \leftarrow) \cdot v$ , где  $v$  это вектор нашего выбора. Хотя на практике формирование целого гессиана является дорогостоящим и зачастую трудоемким, вычисление векторного продукта Гессиана столь же дешево, как вычисление градиента, когда наша функция представляется в виде арифметической схемы [Pearlmutter, 1994], как в случае нейронных сетей ,

Обозначения: Мы используем жирные заглавные буквы  $A, B$  для обозначения матриц и жирных строчных букв  $x, y$  для обозначения векторов. Для векторов мы используем  $\| \cdot \|_K$  обозначать  $2$ -норма, а для матриц мы используем  $\| \cdot \|_K$  обозначать

<sup>2</sup> Оригинальный документ сообщает о скорости  $\sim$

$O(\epsilon^{-3.25})$  из-за другого определения  $\epsilon$ -стационарных второго порядка точка,  $\min ( \rho^2 F(x) )$   $O(\epsilon^{-1.4})$ , который слабее стандартного определения, как в определении 1.

спектральная норма и  $\lambda_{\min}(\cdot)$  обозначить минимальное собственное значение. Если не указано иное, мы используем обозначение  $O(\cdot)$  скрыть только абсолютные константы, которые не зависят ни от какого параметра задачи, и обозначения  $\sim$

$O(\cdot)$  скрыть только абсолютные константы и логарифмические факторы.

## 2.1 Допущения

На протяжении всей статьи мы предполагаем, что функция  $e(\text{Икс})$  ограничен снизу некоторым оптимальным значением  $e_{\min}$ .

Мы также делаем следующие предположения о гладкости функции:

Предположение 1. Функция  $e(\text{Икс})$  имеет

- $\nabla e$  - Градиенты Липшица и  $\nabla^2 e$  - Липшиц гессиан: для всех  $\text{Икс}_1$  а также  $\text{Икс}_2$ ,

$$\|\nabla e(\text{Икс}_1) - \nabla e(\text{Икс}_2)\| \leq L \|\text{Икс}_1 - \text{Икс}_2\|; \|\nabla^2 e(\text{Икс}_1) - \nabla^2 e(\text{Икс}_2)\| \leq K \|\text{Икс}_1 - \text{Икс}_2\|.$$

Приведенные выше предположения утверждают, что градиент и гессиан не могут резко измениться в небольшой локальной области, и являются

стандартными в предыдущей работе по обходу седловых точек и нахождению локальных минимумов. Далее мы сделаем следующие предположения о дисперсии стохастических градиентов и стохастических гессианов:

Предположение 2. Функция  $e(\text{Икс}, \tau)$  имеет

- для всех  $\text{Икс}$ ,  $E \|\nabla e(\text{Икс}, \tau) - \nabla e(\text{Икс})\|^2 \leq \sigma^2$  и также  $\nabla e(\text{Икс}, \tau) = \nabla e(\text{Икс}) + M_1$  в виде;
- для всех  $\text{Икс}$ ,  $E \|\nabla^2 e(\text{Икс}, \tau) - \nabla^2 e(\text{Икс})\|^2 \leq \sigma^2$  и также  $\nabla^2 e(\text{Икс}, \tau) = \nabla^2 e(\text{Икс}) + M_2$  в виде

Мы отмечаем, что приведенные выше допущения не являются существенными для нашего результата и могут быть заменены любыми условиями, которые вызывают концентрацию. Более того, ограниченное гессианское предположение может быть удалено, если допустить, что  $e(\text{Икс}, \tau)$  имеет  $\sigma$ -Градиенты Липшица для всех  $\tau$ , что гарантирует концентрацию Гессе без дальнейшего предположения о дисперсии  $\nabla^2 e(\text{Икс}, \tau)$  [Tropp et al., 2015].

## 2.2 Кубическая регуляризация

Наша цель в этой статье - найти  $\hat{x}$ -стационарная точка второго порядка, которую мы определяем следующим образом:

Определение 1. Для  $\nabla e$ -Функция Гессиана Липшица  $e$ , мы говорим, что  $\hat{x}$  является  $\epsilon$ -стационарная точка второго порядка (или  $\epsilon$ -приблизительный локальный минимум) если

$$\|\nabla e(\hat{x})\| \leq \epsilon \text{ а также } \lambda_{\min}(\nabla^2 e(\hat{x})) \geq \frac{\epsilon}{L} \quad (2)$$

$\hat{x}$ -Стационарная точка второго порядка не только имеет небольшой градиент, но также имеет гессиан, близкий к положительному полуопределенному. Таким образом, это часто также упоминается как  $\epsilon$ -приблизительный локальный минимум. В детерминированной среде кубическая регуляризация [Нестеров, Поляк, 2006] является классическим алгоритмом поиска стационарной точки второго порядка  $\nabla e$ -Функция Гессиана-Липшица  $e(\text{Икс})$ . В каждой итерации сначала формируется локальная верхняя граница функции с использованием разложения Тейлора третьего порядка вокруг текущей итерации.  $\text{Икс}_t$ :

$$m_t(x) = e(\text{Икс}_t) + \nabla e(\text{Икс}_t)^T (x - \text{Икс}_t) + \frac{1}{2} (x - \text{Икс}_t)^T \nabla^2 e(\text{Икс}_t) (x - \text{Икс}_t) + \frac{K}{6} \|x - \text{Икс}_t\|^3.$$

Это называется *кубической подмоделью*. Затем кубическая регуляризация минимизирует эту кубическую подмодель для получения следующей итерации:  $\text{Икс}_{t+1} = \arg\min_{\text{Икс}} m_t(\text{Икс})$ . Когда кубическая подмодель может быть решена точно, кубическая регуляризация требует  $O(\frac{1}{\epsilon^{1.5}})$  итераций, чтобы найти  $\epsilon$ -стационарная точка второго порядка.

$\epsilon^{1.5}$

Чтобы применить этот алгоритм в стохастической установке, необходимо решить три проблемы: (1) у нас есть доступ только к стохастическим градиентам и гессианам, а не к истинному градиенту и гессиану; (2) наше единственное средство взаимодействия с гессианом - через произведения гессен-векторов; (3) кубическая подмодель не может быть решена точно на практике, только до некоторого допуска. Мы обсуждаем, как преодолеть каждое из этих препятствий в нашей статье.

## 3 Основные результаты

Мы начнем с мета-алгоритма стохастической кубической регуляризации общего назначения в алгоритме 1, который использует подпрограмму черного ящика для решения стохастических кубических подзадач. На высоком уровне, в

---

Алгоритм 1 Стохастическая кубическая регуляризация (мета-алгоритм)

---

Входные данные: размеры мини-партии  $N_1, N_2$ , инициализация  $\text{Икс}_0$ , количество итераций  $T_{\text{max}}$ , и окончательная терпимость  $\epsilon$ .

```

1: за  $\tau = 0, \dots, T_{\text{max}}$  делать
2:   Образец  $S_1 \leftarrow \{ \leftarrow S_1, y=1 \}, S_2 \leftarrow \{ \leftarrow S_2, y=1 \}$ 
3:    $\Gamma_\tau = \frac{1}{|S_1|} \Pi_{S_1} \leftarrow \frac{1}{2} S_1 p e(\text{Икс}_\tau, \leftarrow y)$ 
4:    $B_\tau[\cdot] = \frac{1}{|S_2|} \Pi_{S_2} \leftarrow \frac{1}{2} S_2 p_2 e(\text{Икс}_\tau, \leftarrow y[\cdot])$ 
5:    $m \leftarrow \text{Cubic-Subsolver}(\Gamma_\tau, B_\tau[\cdot], \epsilon)$ 
6:    $\text{Икс}_{\tau+1} = \text{Икс}_\tau + Q \leftarrow m$ 
7:   если  $m \leq \frac{1}{100}$  тогда
8:      $\text{Cubic-Finalsolver}(\Gamma_\tau, B_\tau[\cdot], \epsilon)$ 
9:      $\text{Икс} \leftarrow \text{Икс}_{\tau+1}$ 
10:    переменная
11:    конец, если
12: конец для вывода: Икс – если условие досрочного завершения было достигнуто, в противном случае финальная итерация  $\text{Икс}_{T_{\text{max}}+1}$ .

```

---

Чтобы разобраться со случайными градиентами и гессианами, мы выберем две независимые миниатчи  $S_1$  а также  $S_2$  на каждой итерации. Обозначая средний градиент и средний гессиан

$$\Gamma_\tau = \frac{1}{|S_1|} \Pi_{S_1} p e(\text{Икс}_\tau, \leftarrow y), B_\tau = \frac{1}{|S_2|} \Pi_{S_2} p_2 e(\text{Икс}_\tau, \leftarrow y). \quad (3)$$

это подразумевает *стохастическая кубическая подмодель*:

$$M_\tau(x) = e(\text{Икс}_\tau) + (\nabla M_\tau(\text{Икс}_\tau))^T (x - \text{Икс}_\tau) + \frac{1}{2} (x - \text{Икс}_\tau)^T B_\tau (x - \text{Икс}_\tau) + \frac{1}{6} K_{xx} \tau K_3. \quad (4)$$

Хотя подзадача зависит от  $B_\tau$ , отметим, что наш мета-алгоритм никогда явно не формулирует эту матрицу, предоставляя только доступ к подсолверу  $B_\tau$  через гессиано-векторные произведения, которые мы обозначим  $B_\tau[\cdot]: \mathbb{R}^d \rightarrow \mathbb{R}^d$ . Следовательно, мы предполагаем, что субсолвер выполняет градиентную оптимизацию для решения подзадачи, так как  $p M_\tau(\text{Икс})$  зависит от  $B_\tau$  только через  $B_\tau[x - \text{Икс}]$ .

После выборки мини-пакетов для градиента и гессиана алгоритм 1 вызывает черный кубический субсолвер для оптимизации стохастической подмодели  $M_\tau(\text{Икс})$ . Подсолвер возвращает изменение параметра  $m$ , примерный минимизатор подмодели вместе с соответствующим изменением в значении подмодели,  $m = M_\tau(\text{Икс}_{\tau+1}) - M_\tau(\text{Икс}_\tau)$ . Затем алгоритм обновляет параметры, добавляя  $m$  к текущей итерации и проверяет,  $m$  удовлетворяет условию остановки. Более подробно, подпрограмма Cubic-Subsolver принимает вектор  $\Gamma$  и функцию для вычисления гессиано-векторных произведений  $B[\cdot]$ , затем оптимизирует полином третьего порядка  $\sim$

$M(\text{знак равно } \Gamma + \frac{1}{2} B[\cdot - \Gamma] + \frac{1}{6} K_{xx} \tau K_3)$  обозначим минимизатор этого многочлена. В общем,  $\text{subsolver}$  не может вернуть точное решение  $?$ . Следовательно, мы терпим определенную степень субоптимальности:

**Состояние 1.** Для любой фиксированной малой константы  $c$ , Cubic-Subsolver  $(\Gamma, B[\cdot], \epsilon)$  заканчивается в  $T(\epsilon)$  градиентные итерации (которые могут зависеть от  $c$ ), и возвращает удовлетворяющий по крайней мере одному из следующего:

1. Максимум  $M(\cdot), f(\text{Икс}_{\tau+1}) \in (\text{Икс}_\tau, \text{Икс}_{\tau+1})$ . (Дело 1)
2.  $k k \cdot k \cdot k \cdot c Q \leftarrow$  и если  $k \cdot k$  тогда  $\sim M(\cdot) \sim M(?) + c \frac{1}{12} \cdot \rightarrow k \cdot k$ . (Дело 2)

Первое условие выполняется при изменении параметра  $\text{Икс}$  приводит к уменьшению подмодели и функции оба достаточно большие (Дело 1). Если это не выполняется, второе условие гарантирует, что оно не слишком велико по отношению к истинному решению. и что кубическая подмодель решается с точностью  $c \cdot \rightarrow k \cdot k$  когда  $k \cdot k$  большой (Случай 2).

Как упомянуто выше, мы предполагаем, что субсолвер использует градиентную оптимизацию для решения подзадачи так, чтобы он получал доступ к гессиану только через произведения вектора гессиана. Соответственно, это может быть любой стандартный алгоритм первого порядка, такой как градиентный спуск, ускоренный градиент Нестерова.



запрашивает градиент  $\rho \sim m(\cdot)$ . В большинстве  $T(\epsilon)$  такие запросы будут сделаны по определению. Таким образом, каждый итерация занимает  $\sim O\left(\frac{1}{\epsilon^{2+2}} \cdot T(\epsilon)\right)$  стохастический градиент / оценка гессиано-векторного произведения при  $\epsilon$  является маленький (см. замечание 1).

Наконец, отметим, что строки 8-11 Алгоритма 1 дают условие завершения нашего мета-алгоритма. При уменьшении значения подмодели  $m$  слишком мала, наша теория гарантирует Икс  $\tau$  является стационарная точка второго порядка, где? является оптимальным решением кубической подмодели. Тем не менее, Cubic-Subsolver может дать только неточное решение удовлетворяющее Условию 1, которое не достаточно для Икс  $\tau$  быть стационарная точка второго порядка. Поэтому мы называем Cubic-Finalsolver (который просто использует градиентный спуск и подробно описан в алгоритме 3) для решения подзадачи с более высокой точностью.

### 3.1 Градиентный спуск как кубический субсолвер

Одним конкретным примером кубического субсолвера является простой вариант градиентного спуска (Алгоритм 2), изученный в работе Кармона и Дучи [2016]. Два основных различия относительно стандартного градиентного спуска: (1) строки 1–3: когда  $g$  большой, подмодель  $m(\cdot)$  может быть плохо обусловлен, поэтому вместо делая градиентный спуск, итерация перемещается только на один шаг в  $g$  направление, которое уже гарантирует достаточный спуск; (2) строка 6: алгоритм добавляет небольшое возмущение к  $g$  чтобы избежать некоторого «трудного» случая для кубической подмодели. Мы отсылаем читателей к Carmon and Duchi [2016] для получения более подробной информации об алгоритме 2.

Адаптируя их результат для нашей установки, мы получаем следующую лемму, которая утверждает, что субсолвер градиентного спуска удовлетворяет нашему условию 1.

**Лемма 1.** Существует абсолютная константа  $c_0$ , такой, что при одних и тех же предположениях  $e$  (Икс) и такой же выбор параметров  $N_1, N_2$  как и в теореме 1, алгоритм 2 удовлетворяет условию 1 с вероятностью не менее 1  $O(T(\epsilon)) \sim O(n^{\frac{1}{2+2}})$ .

В нашем следующем следствии применяется градиентный спуск (алгоритм 2) в качестве приближенного кубического субсолвера в нашем мета-алгоритме (алгоритм 1), который сразу же дает общее количество оценок градиента и вектора Гессиана для полного алгоритма.

**Следствие 1.** При тех же настройках, что и теорема 1, если  $\epsilon \ll \min \frac{N_2}{c_1 M_1}, \frac{N_2}{c_{22} M_2}$ , то с вероятностью, большей чем 1  $O(\sqrt{\frac{1}{\epsilon^{1.5}}})$  подпрограмма Cubic-Subsolver с алгоритмом 2, с вероятностью, большей чем 1  $O(\frac{1}{\epsilon^{2+2}})$  выведет  $\epsilon$ -стационарная точка второго порядка  $e$  (Икс) в пределах  $O(\frac{1}{\epsilon^{2+2}})$   $P \rightarrow \epsilon$   $\diamond \diamond$   $(7)$

оценки суммарного стохастического градиента и вектора Гессиана.

Из следствия 1 мы видим, что доминирующий член в решении подмодели  $\frac{1}{\epsilon}$  достаточно мал, что дает общую сложность итерации  $\sim O(\frac{1}{\epsilon^{3.5}})$  когда другие проблемно-зависимые параметры постоянны. Это улучшает  $O(\frac{1}{\epsilon^4})$  поли( $r$ ) сложность достигается SGD.

Эскиз 4 доказательства

В этом разделе описаны основные шаги, необходимые для понимания и доказательства нашей основной теоремы (теорема

1). Мы описываем наш высокоуровневый подход и приводим эскизный пример в Приложении А в стохастической настройке, предполагающий доступ оракула к точному субсолверу. В случае неточного субсолвера и подробных доказательств мы отложим добавление Б. Напомним, что при итерации  $T$  алгоритма 1, стохастическая кубическая подмодель  $m_T$  построен вокруг текущей итерации Икс  $\tau$  с формой, приведенной в следующем уравнении,  $m_T(x) = e(\text{Икс } \tau) + (xx - \tau) \triangleright \tau$

$12(x - \tau) \triangleright B_T(x - \tau) + \frac{1}{6} K(x - \tau) K(x - \tau)$ , где  $g$  та же  $B$  усредненные стохастические градиенты и Гессенцы. На высоком уровне мы покажем, что для каждой итерации выполняются следующие два утверждения:

**Претензия 1** Если Икс  $\tau+1$  это не  $\epsilon$ -стационарная точка второго порядка  $e$  (Икс), кубическая субмодель имеет большой спуск  $m_T(\text{Икс } \tau+1) \triangleright m_T(\text{Икс } \tau)$

**Претензия 2** Если кубическая субмодель имеет большой спуск  $m_T(\text{Икс } \tau+1) \triangleright m_T(\text{Икс } \tau)$ , тогда истинная функция также имеет большой спуск  $e(\text{Икс } \tau+1) \triangleright e(\text{Икс } \tau)$

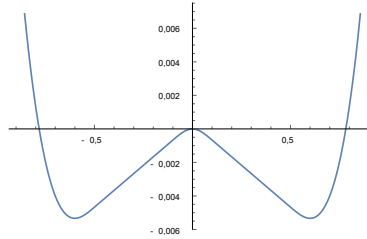


Рисунок 1: кусочно-кубическая функция  $f(x)$  используется вдоль одного из измерений в синтетическом эксперименте. Другое измерение использует масштабированный квадратик.

Учитывая эти претензии, очень просто поспорить за правильность нашего подхода. Мы знаем, что если мы наблюдаем большое уменьшение значения кубической субмодели  $m_T(\text{Икс}_{T+1}) - m_T(\text{Икс}_T)$  в течение алгоритма 1, затем по пункту 2 функция также будет иметь большой спуск. Но с тех пор  $\epsilon$  ограничен снизу, это не может происходить бесконечно, поэтому в конечном итоге мы должны встретить итерацию с малым кубическим подмодельным спуском. Когда это произойдет, мы можем сделать вывод по пункту 1, что  $\text{Икс}_{T+1}$  является  $\epsilon$ -стационарной точкой второго порядка. Отметим, что утверждение 2 особенно важно в стохастической установке, поскольку у нас больше нет доступа к истинной функции, а есть только подмодель. Утверждение 2 гарантирует, что прогресс в  $m_T$  по-прежнему указывает на прогресс в  $\epsilon$ , позволяя алгоритму завершиться в нужное время. Более подробный эскиз и полные доказательства можно найти в Приложении.

## 5 экспериментов

В этом разделе мы представляем эмпирические результаты на синтетических и реальных наборах данных, чтобы продемонстрировать эффективность нашего подхода. Все эксперименты выполняются с использованием TensorFlow [Abadi et al., 2016], который позволяет эффективно вычислять произведения вектора Гессияна с использованием метода, описанного Pearlmutter [1994].

### 5.1 Синтетическая невыпуклая проблема

Мы начнем с построения невыпуклой задачи с седловой точкой, чтобы сравнить предложенный нами подход со стохастическим градиентным спуском. Позволять  $f(x)$  быть  $W$ -образной скалярной функцией, изображенной на рисунке 1, с локальным максимумом в начале координат и двумя локальными минимумами с каждой стороны. Мы откладываем точную форму  $f(x)$  к Приложению D. Мы стремимся

решить проблему

$$\min_{\text{Икс} \in \mathbb{R}^2} W(X_1) + 10 \text{ Икс}_{22}^2,$$

с независимым шумом, взятым из  $N(0, 1)$  добавляется отдельно к каждому компоненту каждого градиента и оценки вектора Гессияна. По построению целевая функция имеет седловую точку в начале координат с собственными значениями Гессе -0,2 и 20, что обеспечивает простой, но сложный тестовый пример. Мы применили наш метод и SGD к этой проблеме, нанося на график объективное значение в зависимости от количества вызовов оракула. Размеры партий и скорости обучения для каждого метода настраиваются отдельно для обеспечения справедливого сравнения; см. Приложение D для деталей. Мы видим, что наш метод способен избежать седловой точки в начале координат и сходиться к одному из глобальных минимумов быстрее, чем SGD.

### 5.2 Глубокий автоэнкодер

В дополнение к описанной выше синтетической проблеме мы также исследуем эффективность нашего подхода к более практической проблеме глубокого обучения, а именно к обучению глубокому автоэнкодеру на MNIST [LeCun and Cortes, 2010]. Наша архитектура состоит из полностью подключенного кодера с размерами (28 → 28) 512! 256! 128! 32 вместе с симметричным декодером. Мы используем нелинейность softplus (определенную как  $\text{Softplus}(x) = \log(1 + \exp(x))$ ) для каждого скрытого слоя - поэлементная сигмоида для конечного слоя и попиксельно  $\ell_2$  потеря между входом и восстановленным выходом как наша целевая функция. Результаты этой задачи автоматического кодирования представлены на рисунке 3. Помимо обучения модели с помощью нашего метода и SGD, мы также включаем результаты с использованием AdaGrad, популярного адаптивного метода первого порядка с высокими эмпирическими характеристиками [Duchi et al., 2011], вместе с результатами для метода, объединяющего уменьшение дисперсии, и информацией второго порядка, предложенной Reddi et al. [2017]. Более подробную информацию о нашей экспериментальной установке можно найти в Приложении D.



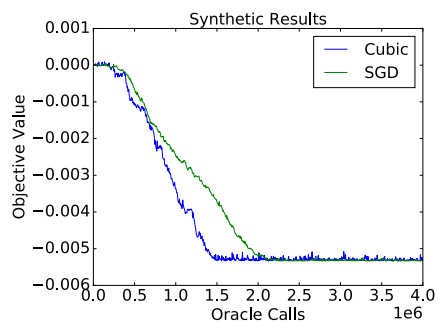


Рисунок 2: Результаты нашей синтетической невыпуклой задачи оптимизации. Стохастическая кубическая регуляризация избегает седловой точки в начале координат и сходится к глобальному оптимуму быстрее, чем SGD.

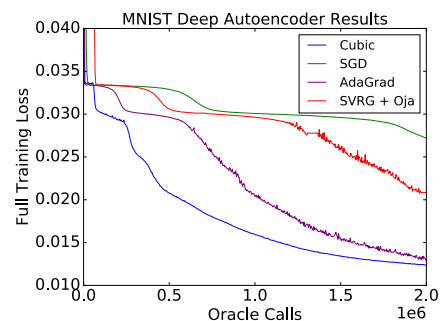


Рисунок 3: Результаты выполнения задачи глубокого автокодирования MNIST. В задаче оптимизации присутствуют несколько седловых точек. Стохастическая кубическая регуляризация способна ускользнуть от них быстрее, что позволяет достичь локального минимума быстрее, чем другие методы.

Мы наблюдаем, что стохастическая кубическая регуляризация быстро выходит за пределы двух седловых точек и опускается к локальному оптимуму, в то время как AdaGrad требуется в два-три раза дольше, чтобы избежать каждой седловой точки, а SGD еще медленнее. Это демонстрирует, что включение информации о кривизне с помощью векторно-гессианских продуктов может помочь на практике избежать седловых точек. Гибридный метод, состоящий из SVRG, чередующийся с алгоритмом Ожы для гессианского спуска, также улучшает SGD, но не соответствует производительности нашего метода или AdaGrad.

## 6. Заключение

В этой статье мы представили стохастический алгоритм, основанный на классическом кубически-регуляризованном методе Ньютона для невыпуклой оптимизации. Наш алгоритм доказуемо находит  $\epsilon$ -приближительные локальные минимумы в

$O(\epsilon^{-3.5})$  оценки общего градиента и вектора Гессе, улучшая  $\sim O(\epsilon^{-4})$  поли(  $n$ )

курс SGD. Наши эксперименты показывают благоприятную эффективность нашего метода относительно SGD как в синтетической, так и в реальной проблеме.

## Ссылки

Мартин Абади, Пол Бархам и др. Тензорный поток: система для крупномасштабного машинного обучения. В 12 *USENIX Symposium по разработке и внедрению операционных систем*, стр. 265–283, 2016. Наман Агарвал, Зейуан Аллен-Чжу,

Брайан Буллинс, Эльад Хазан и Тэньюма. Нахождение приближительного локальные минимумы быстрее градиентного спуска. В *Материалы 49-го ежегодного симпозиума ACM SIGACT по теории вычислений*, 2017.

Зеюань Аллен-Чжу. Наташа 2: Быстрее невыпуклая оптимизация, чем SGD. *Препринт arXiv Arxiv: 1708.08694*, 2017.

Зеюань Аллен-Чжу и Юаньчжи Ли. Neon2: поиск локальных минимумов с помощью оракулов первого порядка. *Архив Препринт arXiv: 1711.06673*, 2017.

Яир Кармон и Джон Дючи. Градиентный спуск эффективно находит кубически-регуляризованные невыпуклые Шаг Ньютона. *Препринт arXiv arXiv: 1612.00547*, 2016.

Яир Кармон, Джон Дючи, Оливер Хиндер и Аарон Сидфорд. Ускоренные методы для невыпуклых оптимизация. *Препринт arXiv arXiv: 1611.00756*, 2016.

Коралия Картис, Николас Гулд и Филипп Тоинт. Адаптивные методы кубической регуляризации для оптимизация без ограничений. Часть II: сложность оценки функций и производных в худшем случае. *Математическое программирование*, 130 (2): 295–319, 2011.

- Анна Чороманска, Микаэль Хенафф, Майкл Матье, Жерар Бен Арус и Янн ЛеКун.  
Поверхности потерь многослойных сетей. В *Материалы восемнадцатой Международной конференции по искусственному интеллекту и статистике*, 2015.
- Фрэнк Э Кертис, Даниэль П Робинсон и Мохаммадреза Самади. Алгоритм области доверия с сложностью итерации в худшем случае  $O(n^3)$  для невыпуклой оптимизации. *Математическое программирование*, 162 (1-2): 1–32, 2017.
- Джон К. Дучи, Элад Хазан и Йорам Сингер. Адаптивные субградиентные методы для онлайн-обучения и стохастическая оптимизация. *Журнал исследований машинного обучения*, 12: 2121–2159, 2011. Ронг Ге, Фуронг Хуан,
- Чи Джин и Ян Юань. Спасаясь от седловых точек - онлайн стохастический градиент для тензорного разложения. В *Материалы 28-й конференции по теории обучения*, 2015.
- Ронг Гэ, Чи Джин и Йи Чжан. Никаких ложных локальных минимумов в невыпуклых задачах низкого ранга: А единый геометрический анализ. В *Материалы 34-й Международной конференции по машинному обучению*, 2017.
- Прадик Джайн, Чи Джин, Шам М. Какаде, Пранит Нетрапалли и Аарон Сидфорд. Поток PCA: Матрица соответствия Бернштейна и почти оптимальная конечная выборка гарантируют алгоритм Оя. В *Конференция по теории обучения*, страницы 1147–1164, 2016.
- Чи Джин, Ронг Ге, Пранит Нетрапалли, Шам М. Какаде и Майкл И. Джордан. Как убежать от седловых точек эффективно. В *Материалы 34-й Международной конференции по машинному обучению*, 2017.
- Ри Джонсон и Тонг Чжан. Ускорение стохастического градиентного спуска с использованием прогнозирующей дисперсии снижение. В *Достижения в нейронных системах обработки информации*, 2013.
- Йонас Мориц Колер и Аурелиен Луччи. Подвыборка кубической регуляризации для невыпуклых оптимизации. В *Материалы 34-й Международной конференции по машинному обучению*, 2017.
- Ян ЛеКун и Коринна Кортес. MNIST База рукописных цифр, 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Юрий Нестеров. *Вводные лекции по выпуклой оптимизации: базовый курс*, Том 87. Спрингер Наука и Бизнес Медиа, 2013.
- Юрий Нестеров и Борис Т Поляк. Кубическая регуляризация метода Ньютона и его глобальные характеристики. *Математическое программирование*, 108 (1): 177–205, 2006. Barak A Pearlmutter. Быстрое точное умножение на гессиан. *Нейронные вычисления*, 6 (1): 147–160, 1994.
- Sashank J Reddi, Manzil Zaheer, et al. Общий подход к выходу из седловых точек. *Препринт arXiv Arxiv: 1709.01434*, 2017.
- Герберт Роббинс и Саттон Монро. Метод стохастической аппроксимации. *Анналы математики Статистика*, страницы 400–407, 1951.
- Клемент В. Ройер и Стивен Дж. Райт. Анализ сложности алгоритмов поиска строк второго порядка для плавной невыпуклой оптимизации. *Препринт arXiv arXiv: 1706.03131*, 2017.
- Джу Сун, Цин Цюй и Джон Райт. Полное восстановление словаря по сфере I: обзор и геометрическая картина. *IEEE Труды по теории информации*, 2016a. Джу Сун, Цин Цюй и Джон Райт. Геометрический анализ фазы поиска. В *IEEE International Symposium on Information Theory*. IEEE, 2016b.
- Джозел А. Тропп и соавт. Введение в матричные неравенства концентрации. *Основы и тенденции в машинном обучении*, 8 (1-2): 1–230, 2015.
- Пэн Сюй, Фарбод Руста-Хорасани и Майкл В.М. Методы ньютоновского типа для невыпуклых оптимизация под неточной гессенской информацией. *Препринт arXiv arXiv: 1708.07164*, 2017.
- И Сюй и Тяньбао Ян. Стохастические алгоритмы первого порядка для выхода из седловых точек почти линейное время *Препринт arXiv arXiv: 1711.01944*, 2017.