

МОСКОВСКИЙ ФИЗИКО-ТЕХНИЧЕСКИЙ ИНСТИТУТ
(национальный исследовательский университет)
ФИЗТЕХ-ШКОЛА ПРИКЛАДНОЙ МАТЕМАТИКИ И ИНФОРМАТИКИ
КАФЕДРА ИНТЕЛЛЕКТУАЛЬНЫХ СИСТЕМ

Крейнин Матвей Вадимович

МЕТОДЫ ПРЕДОБУСЛАВЛИВАНИЯ С ЗАТУХАНИЕМ ВЕСОВ

03.03.01 — Прикладные математика и физика

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА БАКАЛАВРА

Научный руководитель:

к.ф.-м.н. А. Н. Безносовых

Москва — 2024

Аннотация

Исследуется задача минимизации целевой функции потерь. Рассматривается проблема оптимизации целевой функции потерь градиентными методами первого порядка. Исследуется сходимость методов градиентной оптимизации с предобуславливанием, использующих регуляризацию с затуханием весов. Специально рассматриваются популярные методы оптимизации из данного класса методов такие как AdamW и OASIS. Исследуются различные альтернативы этим методам с целью изучения их скорости сходимости и точности, показываемой моделью. Предлагается новый способ добавления регуляризации в метод оптимизации Adam. Доказывается теорема о скорости сходимости данных методов при различных допущениях на функцию потерь и показывается сходимость к исходной функции потерь. Проводятся вычислительные эксперименты с различными эталонными наборами данных, моделями и проводится анализ гиперпараметров, чтобы сравнить их на реальных задачах.

Содержание

1	Введение	4
2	Обзор литературы	7
3	Основная часть	8
3.1	Обозначения	8
3.2	Затухание весов	8
3.3	Скорость сходимости методов с предобуславливанием и затуханием весов	10
3.4	Решение методов с предобуславливанием и затуханием весов . .	15
3.5	Эксперименты	16
A	Приложение	22
A.1	Доказательство лемм	22

1 Введение

Огромная часть машинного обучения основана на решении задачи оптимизации без ограничений

$$\min_{w \in \mathbb{R}^d} f(w). \quad (1)$$

Задачи вида (1) охватывают множество приложений, включая минимизацию эмпирического риска [1], глубокое обучение [2], и задачи обучения с учителем [3] такие, как наименьшие квадраты с регуляризацией [4] или логистическая регрессия [5].

Классический метод решения задачи оптимизации (1) это градиентный спуск.

$$w_{t+1} = w_t - \eta \nabla f(w_t), \quad (2)$$

Задача минимизации (1) может быть трудноразрешимой особенно, когда размер выборки крайне велик или размерность задачи велика.

В таких случаях подсчет полного градиента на каждой итерации в градиентном спуске становится очень дорогим в плане времени или вычислительных ресурсов, которые нужны для этого, особенно учитывая, что градиентному спуску часто требуется большое количество итераций для сходимости. В современном машинном обучении, особенно с появлением глубокого обучения, растет интерес к решению все более больших и сложных задач. Популярным решением для таких проблем стал стохастический градиентный спуск [6].

С течением времени методы оптимизации постоянно совершенствовались и развивались, становясь все более сложными и запутанными. Одним из важнейших аспектов в оптимизации является правильный подбор размера шага в ходе итерационного процесса. Адаптивные методы с градиентным масштабированием динамически регулируют этот размер шага на основе информации о градиенте. Такое адаптивное поведение, применяемое к каждой переменной, улучшает процесс оптимизации, эффективно перемещаясь по сложным

ландшафтам и обеспечивая оптимальный прогресс для каждой переменной [7]. В частности, эти методы приобрели значительную популярность в области машинного обучения, где преобладают высокоразмерные задачи [8, 9].

Более подробно под методами с масштабированным градиентом понимаются техники, предполагающие предобуславливание градиента задачи по определенной матрице D_t , что позволяет градиенту учитывать геометрию задачи. В общем случае шаг алгоритмов с предобуславливанием может быть выражен как следующая модернизация шага (2):

$$w_{t+1} = w_t - \eta \cdot D_t^{-1} g_t, \quad (3)$$

где g_t - несмещенный стохастический градиент.

Идея использования матрицы шкалирования отсылает нас к методу Ньютона, где $D_t = \nabla^2 f(w)$. Однако вычисление и обращение гессиана сопряжено со значительными трудностями, что приводит нас к необходимости использования определенных эвристик в качестве замены матрицы D_t . Примерами таких эвристических методов являются Adagrad [10], Adam [11], RMSProp, OASIS [12] и так далее, где стратегии вычислений для D_t не требуют оценки гессиана. Например, в Adagrad предусловие представлено в виде:

$$D_t = \text{diag} \left\{ \sqrt{\sum_{t'=0}^t g_{t'} \odot g_{t'}} \right\},$$

где \odot - Адамарово произведение. На самом деле этот подход использует только стохастические градиенты.

RMSProp и Adam используют похожие идеи:

$$D_t^2 = \beta D_{t-1}^2 + (1 - \beta) \text{diag} \{g_t \odot g_t\}$$

где $\beta \in (0, 1)$ представляет собой степень учета предыдущих итераций [11].

В OASIS используется другой подход:

$$D_t = \text{diag} \{z_k \odot \nabla^2 f(w_t) z_k\},$$

где z_k - случайный вектор из распределения Рандамахера, т.е. каждый элемент вектора $z_k^i \in \{-1, 1\}$ с вероятностью $\frac{1}{2}$ [12]. На первый взгляд кажется, что используется матрица гессиана, но на самом деле она аппроксимируется через дифференцирование скалярной функции.

Несмотря на преимущества методов предобуславливания, они склонны к переобучению, таким образом, возникает необходимость в их совместном применении с регуляризацией. Этот подход широко применяется для решения различных задач машинного обучения, включая классификацию изображений [13], распознавание речи [14] и обработку естественного языка [15], и показал свою эффективность в улучшении обобщающей способности нейронных сетей [16].

С регуляризацией задача (1) переформулируется как

$$\min_{w \in \mathbb{R}^d} F(w) := f(w) + r(w), \quad (4)$$

где r - функция регуляризации.

В методах с предуславливанием есть несколько способов добавления регуляризации. Можно добавить регуляризатор r в подсчет g_t , и тогда он будет учитываться при вычислении D_t . Этот способ равносильно рассмотрению оптимизационной задачи (4). Или же мы можем добавить регуляризатор только на последнем шаге, уменьшая норму w [17]. Такой способ регуляризации называется затуханием весов и, как ни странно, оказывается более эффективным в практических задачах. Существует и другой способ рассмотрения регуляризатора, который будет рассмотрен далее в статье.

Несмотря на свою практическую эффективность, методы, использующие затухание весов, относительно мало изучены с точки зрения теории сходимости методов. В связи с этим возникает ряд исследовательских вопросов:

- *Сходятся ли с теоретической точки зрения методы с предобуславливанием и затуханием весов?*
- *Если сходятся, то какова скорость их сходимости?*
- *К какой задаче они сходятся?*

2 Обзор литературы

Стохастические методы имеют обширный анализ их сходимости [18, 19, 20], в то время как методы, включающие предобуславливания, являются относительно новыми и неизученными. В одной из первых работ по предобуславливанию [10] авторы провели тщательный анализ теории сходимости Adagrad. Однако в более поздних работах, например, обсуждающих RMSProp [21] или Adam [11], теоретическим аспектам уделяется мало внимания, либо существующая теория содержит неточности в доказательстве.

Со временем ошибки были исправлены, что привело к разработке надежных теорий сходимости для методов с предобуславливанием [22, 23]. В другом исследовании Лоцилов и Хуттер [17] исследовали свойства алгоритмов Adam и AdamW в терминах гиперпараметров, а также изучили методы рестартов. Чжан и др. [24] исследовали механизм заглядывания в будущее в Adam. В [25] исследователи из Nvidia предложили новый способ добавления регуляризации в алгоритм Adam, который на их экспериментах дал прирост в качестве обучения. Совсем недавно были созданы теории сходимости для современных методов, таких как OASIS [12, 26]. Кроме того, появилась теория, рассматривающая изменяющиеся во времени матрицы предобуславливания [27]. Тем не менее, многие вопросы в этой области остаются без ответа. Некоторые из них сформулированы в конце предыдущего параграфа и рассматриваются в нашей статье.

3 Основная часть

3.1 Обозначения

- Мы используем x^i , где x это вектор и $i \in \overline{1, d}$ обозначает i -ю компоненту d -мерного вектора x .
- Для любых $x, y \in \mathbb{R}^d$ скалярное произведение обозначается, как $\langle x, y \rangle := \sum_{i=1}^d x^i y^i$.
- L_f – константа липшица функции f , то есть $\forall x, y \in \mathbb{R}^d \rightarrow f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L_f}{2} \|x - y\|^2$
- $\|x\| := \sqrt{\langle x, x \rangle}$, где $x \in \mathbb{R}^d$ это l_2 норма вектора x .
- $\|x\|_A^2 := x^T A x$, где $x \in \mathbb{R}^d, A \in \mathbb{R}^{d \times d}$
- Для матрицы $A \in \mathbb{R}^{d \times d}$, A^{-1} – обратная матрица.
- Мы используем $A \preceq B$ для двух матриц $A, B \in \mathbb{R}^{d \times d}$, чтобы обозначить что $x^T A x \leq x^T B x$ для любых $x \in \mathbb{R}^d$.
- $\text{diag} \{\beta_1, \dots, \beta_d\}$ – диагональная матрица, состоящая из элементов: $\beta_1, \dots, \beta_d \in \mathbb{R}$.

3.2 Затухание весов

Как было сказано выше, в методах с предобуславливанием существует несколько техник добавления регуляризации в оптимизируемую функцию. Мы рассмотрим три различных подхода, которые проиллюстрированы в Алгоритм 1 с помощью различных цветов (каждый отдельный цвет это отдельный алгоритм).

Algorithm 1 Различные способы использования предобуславливания с регуляризацией

Require: η — шаг обучения, f — оптимизируемая функция

while w не сойдется **do**

$t = t + 1$

$g_t \leftarrow$ стохастический градиент f

$g_t \leftarrow g_t + \nabla r(w_t)$ обычная регуляризация

$D_t \leftarrow$ матрица предобуславливания с помощью g_t

$w_t \leftarrow w_{t-1} - \eta \cdot D_t^{-1} g_t$ обычная регуляризация,

$w_t \leftarrow w_{t-1} - \eta \cdot D_t^{-1} (g_t + \nabla r(w_t))$ масштабированное затухание весов,

$w_t \leftarrow w_{t-1} - \eta \cdot D_t^{-1} g_t - \eta \cdot \nabla r(w_t)$ затухание весов,

end while

Если говорить более конкретно, то первая техника регуляризации, показанная в синим, заключается в простом добавлении регуляризационного члена к оптимизируемой функции. Этот регуляризатор включается в стохастический градиент и учитывается при вычислении D_t . По сути, этот подход предполагает применение базового метода оптимизации с предобуславливанием к задаче (4). Вторая техника регуляризации, показанная на рисунке оранжевым, является новым подходом. Хотя член регуляризатора не влияет на вычисление матрицы предобуславливания D_t , он добавляется перед применением D_t . Это означает, что скорость обучения принимается одинаковой для градиента и регуляризатора. Последний рассматриваемый нами подход к регуляризации известен как затухание веса, в алгоритме он подсвечивается цветом красным в алгоритме 1. Как и во втором методе, матрица D_t вычисляется без использования регуляризатора, а в этом методе регуляризатор включен на шаге обновления весов, что позволяет избежать влияния регуляризации на матрицу предобуславливания.

Важно учитывать влияние регуляризации при разработке алгоритмов

оптимизации, и я надеюсь, что моё исследование окажется полезным для исследователей в этой области.

3.3 Скорость сходимости методов с предобуславливанием и затуханием весов

Давайте попробуем оценить скорость сходимости методов с предобуславливанием и затуханием весов.

Хотя шаг оптимизации весов модели может показаться простым, он может быть рассмотрен с другой стороны. Давайте вынесем матрицу D_t^{-1} за скобки, что даёт нам следующий шаг:

$$w_{t+1} = w_t - \eta \cdot D_t^{-1}(\nabla f(w_t) + D_t \nabla r(w_t)). \quad (5)$$

Это подталкивает нас к тому, чтобы вести новую функцию \tilde{r} , такую, что $\nabla \tilde{r}_t(w) = D_t \nabla r(w)$ и новую целевую функцию

$$\tilde{F}_t(w) := f(w) + \tilde{r}_t(w), \quad (6)$$

, где новая целевая функция \tilde{F}_t меняется на каждом оптимизационном шаге, так как D_t тоже обновляется на каждом оптимизационном шаге.

Новый адаптивный регуляризатор \tilde{r}_t в общем случае к сожалению не существует. Поэтому мы наложим ограничения на начальный регуляризатор и структуру предобусловливателя, которые будут оформлены в виде следующих предположений на функцию регуляризации.

Предположение 1. *(Структура регуляризатора) Регуляризатор r сепарабелен, то есть он может быть представлен в следующем виде:*

$$r(w) = \sum_{i=1}^d r_i(w^i),$$

где $r_i(x) \geq 0$ для $i \in \overline{1, d}$ и $x \in \mathbb{R}$.

Предположение 2. (Структура матрицы предобуславливания) Матрица предобуславливания D_t может быть представлена в следующем виде:

$$D_t = \text{diag} \{d_t^1 \dots, d_t^d\}.$$

Хотя эти предположения являются достаточно сильными, но они выполняются для упомянутых ранее методов с предобуславливанием и затуханием весов, также это верно для таких популярных функций регуляризации как регуляризация Тиханова и LASSO регуляризация. Скорость сходимости обычно исчисляется количеством итераций, которые необходимы для достижения определенного уровня погрешности. Чтобы получить оценки количества итераций, необходимых для сходимости к заданной ошибке, мы должны наложить определенные предположения на оптимизируемую функцию потерь. На протяжении всего последующего анализа я предполагаю, что $f : \mathbb{R}^d \rightarrow \mathbb{R}$ является L -гладким и дважды дифференцируемым.

Предположение 3. (L -гладкость)

- Градиент функции f является L_f -гладким, то есть существует такая константа $L_f > 0$ такая, что $\forall x, y \in \mathbb{R}^d$,

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L_f}{2} \|x - y\|^2.$$

- Градиент функции r является L_r -гладким, то есть существует такая константа $L_r > 0$ такая, что $\forall x, y \in \mathbb{R}^d$,

$$r(x) \leq r(y) + \langle \nabla r(y), x - y \rangle + \frac{L_r}{2} \|x - y\|^2.$$

Для того чтобы работать в невыпуклом случае, необходимо ввести ограничение на значения функции регуляризации, это описано в 4.

Предположение 4. (Ограниченность регуляризатора) Регуляризатор ограничен, то есть существует константа $\Omega > 0$ такая, что $\forall w \in \mathbb{R}^d$

$$|r(w)| \leq \Omega.$$

Мы используем обычное ограничение на матрицу предобуславливания, которое сформулировано в предположении 5.

Предположение 5. *(Ограниченность предобуславливателя) Существуют константы $\alpha, \Gamma \in \mathbb{R} : 0 < \alpha < \Gamma$ такие, что*

$$\alpha I \preceq D_t \preceq \Gamma I \Leftrightarrow \frac{I}{\Gamma} \preceq D_t^{-1} \preceq \frac{I}{\alpha}.$$

Это было доказано в [27], что это предположение справедливо для всех современных и популярных алгоритмов с предобуславливанием, таких как Adam, Adagrad, OASIS.

В нашем анализе мы рассматриваем два способа обновления матрицы предобуславливания. В первом методе матрица обновляется через квадраты:

$$(D_{t+1})^2 = \beta(D_t)^2 + (1 - \beta)(H_t)^2, \quad (7)$$

, где H_t - матрица, содержащая новую информацию, а $\beta [0, 1]$ - параметр импульса. Этот подход используется в Adam, а также в более старых методах, таких как RMSProp и AdaHessian. Второй способ является более современным и предполагает использование первых степеней матриц, сохраняя форму преобразования

$$D_{t+1} = \beta D_t + (1 - \beta)H_t, \quad (8)$$

Этот подход используется в OASIS, недавно придуманным методом. В обоих случаях параметр импульса β обычно подбирается близким к 1, что означает, что D_t незначительно меняется в ходе обучения, что может быть формально сформулировано в лемме 1.

Выполнение предположения 5 имеет решающее значение для сходимости и теоретического анализа, и поэтому в алгоритмах часто используется метод императивного выбора для нижней границы матрицы D_t

$$\hat{D}_{t+1}^{ii} = \max\{\alpha, D_t^{ii}\}. \quad (9)$$

Лемма 1. *(Эволюция D_t , Безносиков) Предположим, что для начальной матрицы D_0 выполнены предположения 2 и 5, H_t диагональна с максимальным значением меньше или равным Γ на каждом временном шаге t , и D_t*

эволюционирует в соответствии с (7), (9) или (8), (9), тогда справедливы следующие утверждения:

1. 2 и 5 выполняется для \hat{D}_t для всех t ;
2. $\|\hat{D}_{t+1} - \hat{D}_t\|_\infty \leq \frac{(1-\beta)\Gamma^2}{2\alpha}$ for (7);
3. $\|\hat{D}_{t+1} - \hat{D}_t\|_\infty \leq 2(1-\beta)\Gamma$ for (8).

Эта лемма доказана в ??, где мы опираемся на [27].

Чтобы проводить стохастический анализ, мы должны включить ограничения на стохастический градиент функции. Это формализуется в следующем предположении

Предположение 6. (Ожидания) g_t являются несмещенными и имеют ограниченную вариацию на любом шаге, то есть

$$\mathbb{E}[g_t] = \nabla f(w_t), \mathbb{E}[\|g_t - \nabla f\|^2] \leq \sigma^2. \quad (10)$$

Чтобы получить дополнительные оценки на сходимость методов с предобуславливанием и затуханием весов мы накладываем сильную выпуклость 7 на целевую функцию потерь.

Предположение 7. (Сильная выпуклость) Существует μ_f такая, что $\forall x, y \in \mathbb{R}^d$ выполняется:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu_f}{2} \|x - y\|_2^2$$

С помощью предположений 1 и 2 мы можем доказать существование \tilde{r} и, следовательно, \tilde{F} . Мы оформим это в лемме 2. Мы показываем только существование, но не единственность функции, но в наших оценках \tilde{F} может быть найдена до константы.

Лемма 2. (Существование \tilde{r}) Предполагая, что 1, 2 выполняются, функция \tilde{r} существует и имеет следующую форму:

$$\tilde{r}_t(w) = \sum_{i=1}^d d_t^i r_i(w_i)$$

Используя введенное предположение 3, мы можем гарантировать гладкость для \tilde{r} и оценить его константу Липшица, что формально сформулировано и доказано в лемме 3.

Лемма 3. (*L-гладкость \tilde{r}*) Предполагая, что 1, 2, 3, 5 выполняются, градиент \tilde{r} является $L_{\tilde{r}}$ -непрерывным, то есть существует константа $L_{\tilde{r}} > 0$ такая, что $\forall x, y \in \mathbb{R}^d$,

$$\tilde{r}_t(x) \leq \tilde{r}_t(y) + \langle \nabla \tilde{r}_t(y), x - y \rangle + \frac{L_{\tilde{r}}}{2} \|x - y\|^2,$$

и $L_{\tilde{r}} = \Gamma L_r$.

Используя введенные предположения, мы доказали сходимостъ методов с предобуславливанием и затуханием в общем виде. Наши результаты оформлены в Теорему 1 и Теорему 2. Доказательства теорем можно найти в Приложении ??.

Теорема 1. Предполагая, что 1, 2, 3, 4, 5 выполняются, положим ошибку $\varepsilon > 0$ и шаг обучения удовлетворяют условию:

$$\eta < \frac{2\alpha}{L_f + \Gamma L_r},$$

где L_f, L_r - константа Липшица функций f и r . Пусть существует начальная матрица предобуславливания, которая обновляется в соответствии с условиями леммы 1. Тогда количество итераций, выполняемых алгоритмами с предусловием и убывающим весом, начиная с начальной точки $w_0 \in \mathbb{R}^d$ с $\Delta_0 = \tilde{F}_0(w_0) - f^*$, где \tilde{F}_t определено в (6) и f^* решением задачи (1), необходимое для ε -приближения нормы градиента к 0, может быть ограничено количеством шагов

$$T = \mathcal{O} \left(\frac{\Delta_0 \Gamma}{\left(\eta - \frac{\tilde{L} \eta^2}{2\alpha} \right) \left(\varepsilon - \frac{\delta \Gamma}{\eta - \frac{\tilde{L} \eta^2}{2\alpha}} \right)} \right),$$

где $\tilde{L} = L_f + \Gamma L_r$ и δ может выбрано сколь угодно малым с помощью выбора гиперпараметров α, β, Γ

1. $\delta = \frac{(1-\beta)\Gamma^2}{2\alpha}$ for (7);

2. $\delta = 2(1 - \beta)\Gamma$ for (8).

Теорема 2. *Предполагая, что 1, 2, 3, 4, 5, 7 выполняются, положим ошибку $\varepsilon > 0$ и шаг обучения удовлетворяют условию: $\eta < \frac{\alpha}{4L_f}$, гиперпараметры удовлетворяют условиям: $\lambda < \frac{\alpha\beta^2}{8L_f\Omega_0^2}$, $\beta \geq 1 - \frac{\eta(\mu_f+\lambda)\alpha}{2\Gamma^2}$, $\Omega_0^2 \geq \frac{\alpha^2\beta^2}{8L_f^2}$. Получаем оценку на необходимое количество шагов для сходимости алгоритма к заданной точности:*

$$T = \max \left\{ \log \left(\frac{4\varepsilon(1-\beta)L_f^2}{\alpha\beta^2} \right) \cdot \log \frac{2}{1+\beta}; \frac{4}{\frac{\alpha}{4L_f} \left(\mu_f + \frac{\alpha\beta^2}{8L_f\Omega_0^2} \right)} \log \left(\frac{2\|w_0 - w_*\|_2^2}{\varepsilon} \right) \right\}$$

Эти теоремы устанавливают сходимость методов с предобуславливанием и затуханием весов различных предположениях, а также определяют необходимое количество итераций для заданной точности. Для наших задач простой факт сходимости этих методов имеет огромное значение.

Однако характеристики решения \tilde{w}^* задачи

$$\min_{w \in \mathbb{R}^d} \tilde{F}(w) = f(w) + \tilde{r}(w), \quad (11)$$

к которым сходится этот метод, требуют более глубокого исследования, которое будет рассмотрено в следующем разделе.

3.4 Решение методов с предобуславливанием и затуханием весов

В предыдущем подразделе мы доказали сходимость методов с предобуславливанием, однако выше мы указали, что методы с затуханием весов сходятся к исходному решению задачи оптимизации (4) w^* , а к исходному решению \tilde{w}^* задачи (1), это достигается за счет следующего. Мы получили новую целевую функцию потерь, в которой величина регуляризации динамически изменяется от шага к шагу, на основании матрицы предобуславливания, учитывая, что матрица составляется на основе стохастических градиентов, полученных в

ходе обновления весов модели, получается, что регуляризация получается тем больше, чем больше градиент по данному весу модели, и наоборот тем меньше, чем меньше градиент по весу модели. То есть регуляризация не штрафует веса модели, где градиент вышел на значения близкие к нулю. То есть мы стараемся выйти делать больший шаг там, где стохастический градиент не приблизился к каким-то околнулевым значениям, то есть пока мы не оказались в окрестности какого-то экстремума. За счёт этого получается более разнообразная траектория обновления весов модели, которая позволяет нам получать лучшую сходимость на практике. Эти рассуждения подтверждаются экспериментами, подробно описанными в разделе 3.5.

Оценим разницу между решениями задач (4) и (6). Это ограничение основано на предположениях (3) и свойствах матрицы D_t .

Лемма 4. (*Lower bound*) *Предполагая, что 1, 2 и 3 выполняются, также предполагая, что задачи (6) и (4) имеют соответствующие решения \tilde{w}^* и w^* , тогда разница между решениями может быть ограничена снизу:*

$$\|\tilde{w}^* - w^*\|_{L_F} \geq \|\nabla r(\tilde{w}^*)(I - D_t)\|.$$

Следовательно, можно заметить, что использование регуляризации весов не в прямом подсчете градиента, которое влечет и учет их в матрице предобуславливания приводит к сходимости к решению исходной задачи, в то время, как прямое использование функции регуляризации для подсчета стохастического градиента приводит нас к альтернативному решению. Расхождение между этими решениями зависит от нормы разности между D_t и матрицей тождества ($\|D_t - I\|$). В результате анализ распределения элементов $D := \lim_{t \rightarrow \infty} D_t$ может дать представление о сходимости метода с затуханием весов.

3.5 Эксперименты

Мы рассмотрим два алгоритма OASIS [12] и Adam [11], а также их вариации. Их основное отличие заключается в вычислении матрицы предобуславливания.

В Adam это диагональная матрица, состоящая из квадратов производных, в OASIS - стохастический гессиан, который вычисляется через случайную величину из распределения Рандемахера. Я показываю три варианта регуляризации для Adam и OASIS в Алгоритме 2 и Алгоритме 3 соответственно.

Algorithm 2 Различные способы добавления регуляризации для Adam

Require: $\eta, \beta_1, \beta_2, \epsilon, f, r$

while θ не сойдется **do**

$t = t + 1$

$g_t = \nabla f(w_{t-1}) + \nabla r(w_{t-1})$

AdamL2

$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$

$v_t = \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$

$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} + \nabla r(w_{t-1})$

AdamWH

$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$

$w_t = w_{t-1} - \eta \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} - \eta \nabla r(w_{t-1})$

AdamW

end while

Algorithm 3 Различные способы добавления регуляризации для OASIS

Require: $w_0, \eta_0, D_0, \theta_0 = +\infty$

$w_1 = w_0 - \eta \hat{D}_0^{-1} \nabla f(w_0)$

for $k = 1, 2, \dots$ **do**

$g_k = \nabla f(w_k) + \nabla r(w_{t-1})$

OASISL2

$D_k = \beta D_{k-1} + (1 - \beta_2) \cdot \text{diag} (z_k \odot \nabla^2 (f(w_k) + r(w_k)) z_k)$

OASISWH

$(\hat{D}_k)_{ii} = \max\{|D_k|_{i,i}; \alpha\}, \forall i = \overline{1, d}$

$\eta_k = \min\{\sqrt{1 + \theta_{k-1}} \cdot \eta_{k-1}; \frac{\|w_k - w_{k-1}\|_{\hat{D}_k}}{2\|\nabla f(w_k) - \nabla f(w_{k-1})\|_{\hat{D}_k}^*}\}$

$w_{k+1} = w_k - \eta_k g_k D_k^{-1} - \eta \nabla r(w_{t-1})$

OASISW

$\theta_k = \frac{\eta_k}{\eta_{k-1}}$

end for

В этом разделе приводятся численные эксперименты для вышеупомяну-

тых методов оптимизации. Эксперименты проводились на процессоре *x86* и использованием графического ускорителя NVIDIA GeForce RTX 3090, эксперименты были воспроизведены на 8-ми ядерном процессоре на архитектуре ARM-64.

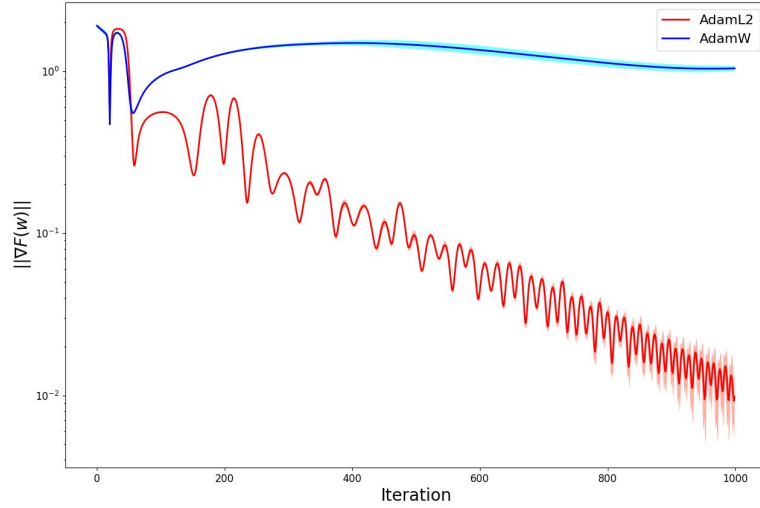


Рис. 1: Adam и AdamW по классическому критерию

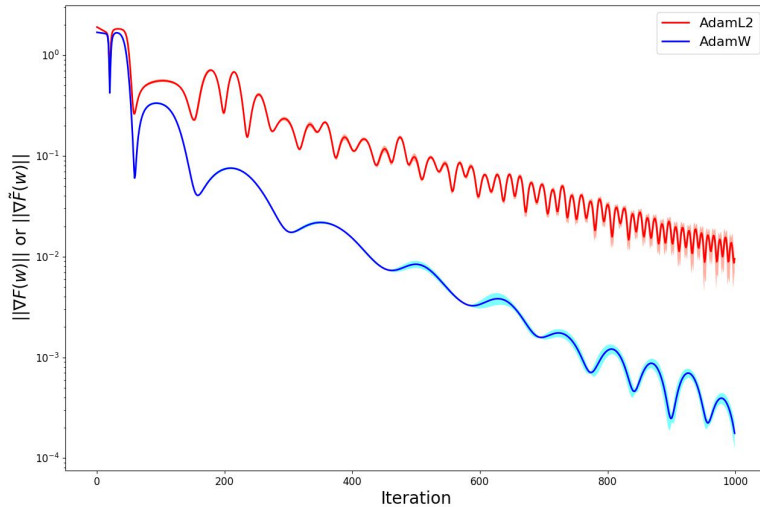


Рис. 2: Adam и AdamW с модифицированным критерием

Список литературы

- [1] Olivier Chapelle, Jason Weston, Léon Bottou, and Vladimir Vapnik. Vicinal risk minimization. *Advances in neural information processing systems*, 13, 2000.
- [2] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [3] Pádraig Cunningham, Matthieu Cord, and Sarah Jane Delany. Supervised learning. *Machine learning techniques for multimedia: case studies on organization and retrieval*, pages 21–49, 2008.
- [4] Ryan M Rifkin and Ross A Lippert. Notes on regularized least squares. 2007.
- [5] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [6] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [7] Elad Hazan, Alexander Rakhlin, and Peter Bartlett. Adaptive online gradient descent. *Advances in Neural Information Processing Systems*, 20, 2007.
- [8] Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse. Three mechanisms of weight decay regularization. *arXiv preprint arXiv:1810.12281*, 2018.
- [9] Zhewei Yao, Amir Gholami, Sheng Shen, Mustafa Mustafa, Kurt Keutzer, and Michael W. Mahoney. Adahessian: An adaptive second order optimizer for machine learning, 2021.
- [10] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.

- [11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [12] Andrew Goldberg, Xiaojin Zhu, Alex Furger, and Jun-Ming Xu. Oasis: Online active semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25, pages 362–367, 2011.
- [13] Feng Zhu, Hongsheng Li, Wanli Ouyang, Nenghai Yu, and Xiaogang Wang. Learning spatial regularization with image-level supervisions for multi-label image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5513–5522, 2017.
- [14] Yingbo Zhou, Caiming Xiong, and Richard Socher. Improved regularization techniques for end-to-end speech recognition. *arXiv preprint arXiv:1712.07108*, 2017.
- [15] Tingting Wu, Xiao Ding, Minji Tang, Hao Zhang, Bing Qin, and Ting Liu. Stgn: an implicit regularization method for learning with noisy labels in natural language processing. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7587–7598, 2022.
- [16] Federico Girosi, Michael Jones, and Tomaso Poggio. Regularization theory and neural networks architectures. *Neural computation*, 7(2):219–269, 1995.
- [17] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [18] Johannes Schneider and Scott Kirkpatrick. *Stochastic optimization*. Springer Science & Business Media, 2007.
- [19] Daniel P Heyman and Matthew J Sobel. *Stochastic models in operations research: stochastic optimization*, volume 2. Courier Corporation, 2004.
- [20] James C Spall. Implementation of the simultaneous perturbation algorithm for stochastic optimization. *IEEE Transactions on aerospace and electronic systems*, 34(3):817–823, 1998.

- [21] T. Tieleman and G. Hinton. Lecture 6.5 - rmsprop: Divide the gradient by a running average of its recent magnitude. *Lecture 6.5 - rmsprop: Divide the gradient by a running average of its recent magnitude*, 2012.
- [22] Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond, 2019.
- [23] Alexandre Défossez, Léon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence proof of adam and adagrad. *arXiv preprint arXiv:2003.02395*, 2020.
- [24] Michael R. Zhang, James Lucas, Geoffrey Hinton, and Jimmy Ba. Lookahead optimizer: k steps forward, 1 step back, 2019.
- [25] Boris Ginsburg, Patrice Castonguay, Oleksii Hrinchuk, Oleksii Kuchaiev, Vitaly Lavrukhin, Ryan Leary, Jason Li, Huyen Nguyen, Yang Zhang, and Jonathan M. Cohen. Stochastic gradient methods with layer-wise adaptive moments for training of deep networks, 2020.
- [26] Abdurakhmon Sadiev, Aleksandr Beznosikov, Abdulla Jasem Almansoori, Dmitry Kamzolov, Rachael Tappenden, and Martin Takáč. Stochastic gradient methods with preconditioned updates. *arXiv preprint arXiv:2206.00285*, 2022.
- [27] Aleksandr Beznosikov, Aibek Alanov, Dmitry Kovalev, Martin Takáč, and Alexander Gasnikov. On scaled methods for saddle point problems. *arXiv preprint arXiv:2206.08303*, 2022.

А Приложение

А.1 Доказательство лемм

Скоро будет перетехано с листочка...